



HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY
School of Information and Communication Technology

CTF CHALLENGE WRITE-UP

Web Exploitation

BKSEC: Maze Maze Mazeee

DucTapCodeDao 202416788

Mục lục

1	Phân tích Challenge	3
2	Sửa HTTP Header bằng Burp Suite và phát hiện ra hướng đi sai	3
3	Đổi target resource của HTTP request và thành công tìm được hướng đi đúng	4
4	Chạy thử đường PATH được Decode và tìm ra cách giải Round 1 của Maze	5
5	Tự động hoá việc lần theo các trang bằng Python	6
5.1	Chạy file code Python và thành công bước vào round 2	8
6	Giải mã Round 2	10

1 Phân tích Challenge

Challenge cho chúng ta một đường link:

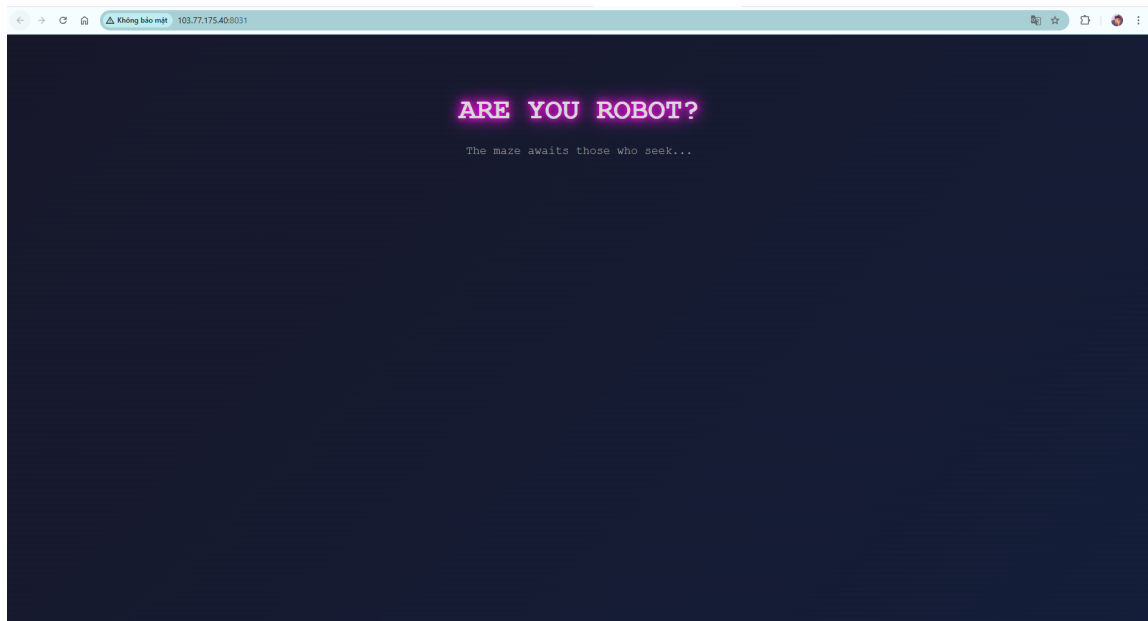
`http://103.77.175.40:8031`

cùng với một đoạn gợi ý:

"Warming up by escaping the maze is very common in CTF competitions."

Vậy nhiệm vụ của chúng ta là **thoát khỏi mê cung** (*maze*) để tìm ra **flag**.

Trước tiên, **Duc** mở trang web bằng trình duyệt *Google Chrome*, và đây là giao diện của trang web:



Hình 1: Giao diện khi mở đường link ở trình duyệt thông thường (Google Chrome).

Nhìn vào giao diện, **Duc** chỉ thấy 2 dòng chữ là “ARE YOU ROBOT?” và “The maze awaits those who seek...”.

Vì thế, **Duc** đoán rằng có thể challenge này yêu cầu người giải thay đổi phần *HTTP Header* của đường link sao cho server nhận dạng người sử dụng là *Robot*.

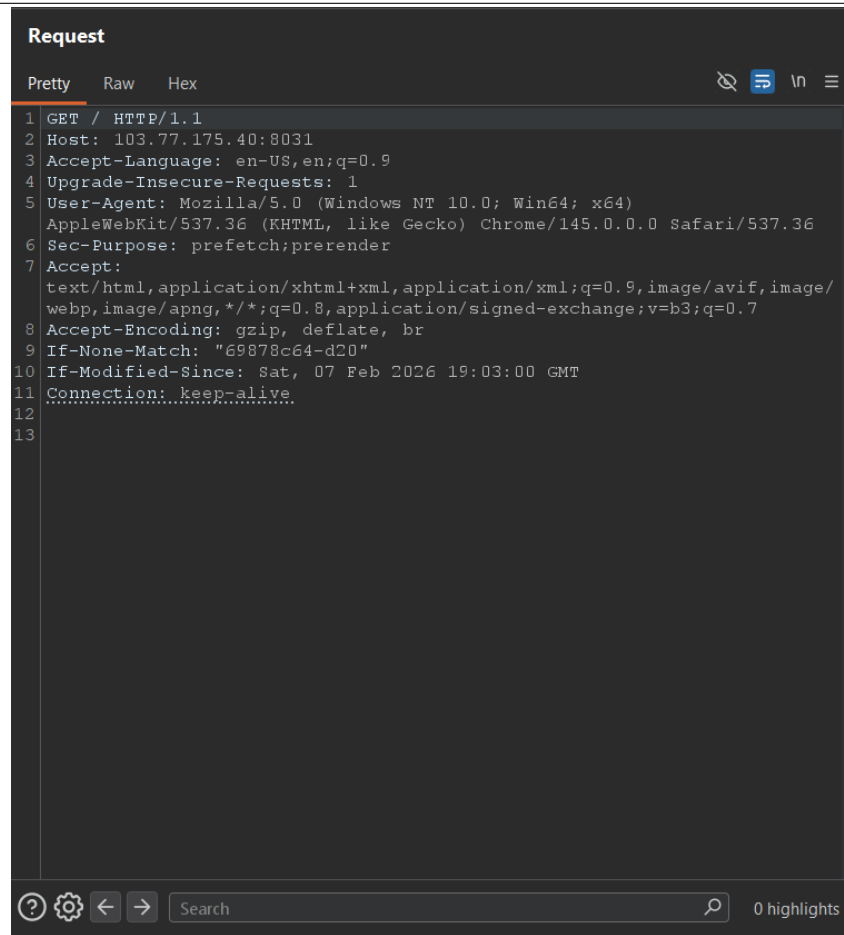
Để làm điều này, **Duc** ném link này vào *Burp Suite* để bắt đầu vọc vạch.

2 Sửa HTTP Header bằng Burp Suite và phát hiện ra hướng đi sai

Duc bật *Intercept* để có thể bắt request dễ hơn, sau đó ném lại đường link vào trình duyệt.

Tiếp theo, **Duc** chọn request *GET* trong mục *HTTP history* và dùng chức năng *Send to Repeater* để chỉnh sửa *HTTP Header* thuận tiện hơn.

Và đây là phần *HTTP Header* của trang web:



```

Request
Pretty Raw Hex
1 GET / HTTP/1.1
2 Host: 103.77.175.40:8031
3 Accept-Language: en-US,en;q=0.9
4 Upgrade-Insecure-Requests: 1
5 User-Agent: Mozilla/5.0 (Windows NT 10.0; Win64; x64)
6 Sec-Purpose: prefetch;prerender
7 Accept:
8 text/html,application/xhtml+xml,application/xml;q=0.9,image/avif,image/
9 webp,image/apng,*/*;q=0.8,application/signed-exchange;v=b3;q=0.7
10 Accept-Encoding: gzip, deflate, br
11 If-None-Match: "69878c64-d20"
12 If-Modified-Since: Sat, 07 Feb 2026 19:03:00 GMT
13 Connection: keep-alive

```

Hình 2: HTTP request header của trang web khi truy cập bằng trình duyệt thông thường.

Duc nhận thấy rằng nếu muốn “giả làm robot”, phần cần ưu tiên chỉnh sửa nhất là trường *User-Agent*.

Vì vậy, **Duc** lần lượt thử thay đổi *User-Agent* thành các giá trị như Robot, rồi đến Mozilla/5.0 (compatible; Googlebot/2.1; +http://www.google.com/bot.html) (*Googlebot*), ...

Tuy nhiên, tất cả các thử nghiệm này đều trả về kết quả **404 Error**.

Điều đó cho thấy hướng tiếp cận dựa trên việc chỉnh *User-Agent* có vẻ không đúng với ý đồ của challenge. Vì vậy, **Duc** quyết định chuyển sang một hướng tiếp cận khác để tiếp tục khai thác challenge này.

3 Đổi target resource của HTTP request và thành công tìm được hướng đi đúng

Duc suy luận rằng nếu mình là một *robot* (crawler), thì thay vì truy cập trực tiếp trang chủ, nhiều khả năng sẽ đọc file *robots.txt* trước để lấy các rule điều hướng.

Vì vậy, **Duc** chỉnh dòng đầu của request trong *Repeater* thành:

```
GET /robots.txt HTTP/1.1
```

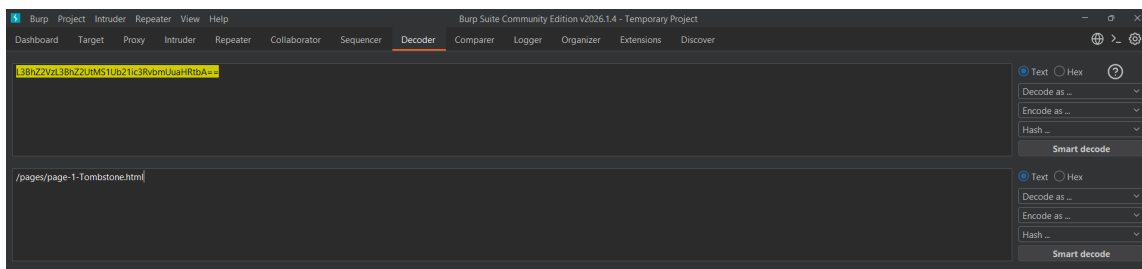
sau đó nhấn *Send*. Kết quả trả về gồm 2 dòng:

```
User-agent: *
Disallow: L3BhZ2VzL3BhZ2UtMS1Ub21ic3RvbmUuaHRtbA==
```

Đây là dấu hiệu cho thấy **Duc** đã đi đúng hướng và tìm được file chứa *PATH* cần thiết. Quan sát kỹ hơn, **Duc** nhận ra giá trị của trường *Disallow* có dạng chuỗi được mã hoá *Base64*. Vì vậy, **Duc** copy chuỗi này, chuyển sang tab *Decoder* trong Burp Suite, paste nội dung vào và chọn *Decode as Base64*.

Sau khi giải mã, **Duc** thu được một *PATH* mới:

```
/pages/page-1-Tombstone.html
```



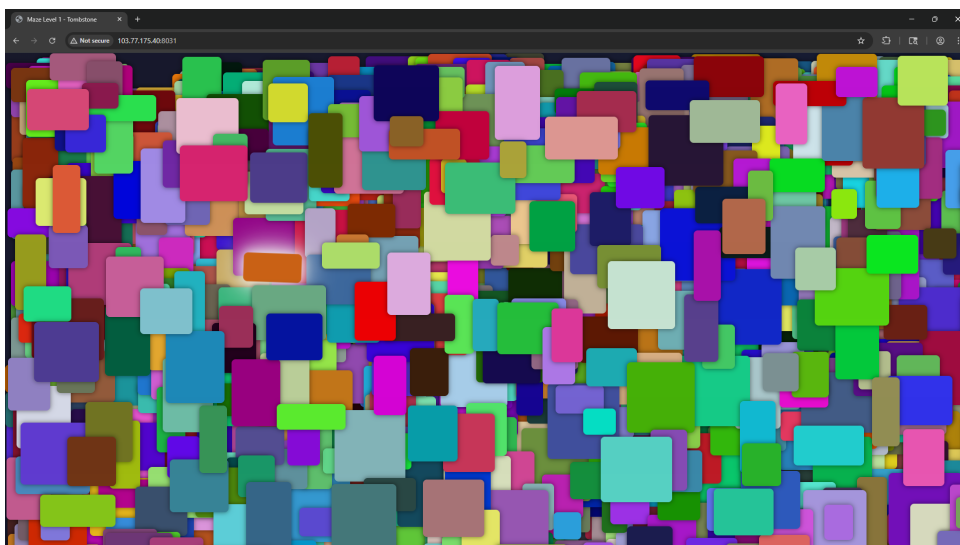
Hình 3: Giải mã giá trị trong trường *Disallow* bằng *Decoder* của Burp Suite để thu được *PATH* kế tiếp.

4 Chạy thử đường *PATH* được Decode và tìm ra cách giải Round 1 của Maze

Sau khi tìm được *PATH* sau khi giải mã, **Duc** đưa lại vào request với dòng đầu:

```
GET /pages/page-1-Tombstone.html HTTP/1.1
```

sau đó nhấn *Send* và thu được giao diện sau:



Hình 4: Giao diện tại */pages/page-1-Tombstone.html* sau khi truy cập bằng *Burp* hoặc trình duyệt.

Quan sát giao diện này, **Duc** nhận thấy hướng khai thác nhiều khả năng nằm trong phần *Source* của *Response* thay vì giao diện hiển thị trực tiếp.

Vì vậy, **Duc** thử tìm từ khoá `href` trong *Response* và phát hiện một *PATH* mới, cụ thể là `page-2-Deadwood.html`.

Duc tiếp tục thử request:

```
GET /pages/page-2-Deadwood.html HTTP/1.1
```

Kết quả trả về là một trang mới nhưng có giao diện gần như giống hệt `page-1`. Sau khi lặp lại thao tác này thêm vài lần, kết quả vẫn tương tự: mỗi trang đều dẫn đến một *PATH* tiếp theo được giấu trong *Source*.

Từ đó, **Duc** nghĩ đến việc viết một script *Python* để tự động hoá quá trình này: truy cập một trang, trích xuất URL của trang tiếp theo, rồi tiếp tục truy cập trang kế tiếp và lặp lại như vậy cho đến khi không còn lấy được *PATH* mới thì dừng lại.

5 Tự động hoá việc lần theo các trang bằng Python

Sau khi nhận thấy mỗi trang đều chứa liên kết đến một trang kế tiếp trong phần *Source*, Duc viết một script Python để tự động hoá quá trình crawl: truy cập trang hiện tại, trích xuất URL tiếp theo, rồi lặp lại cho đến khi gặp *flag* hoặc không còn liên kết mới.

```

1 import re # dùng regex để tìm flag hoặc tìm pattern của URL kế tiếp.
2 import time # dùng sleep() để delay giữa các request (tránh spam server).
3 import argparse # đọc tham số dòng lệnh (--start, --proxy, ...) khi chạy script.
4 from urllib.parse import urljoin, urlparse # urlparse: tách URL để lấy netloc (domain:port) nhằm kiểm tra "out of scope".
5 # urljoin: ghép URL tuyệt đối từ URL hiện tại + href tương đối (giống cách browser xử lý link).
6
7 import requests # gửi HTTP request như một "mini browser".
8 from bs4 import BeautifulSoup # parse HTML để trích href từ các thẻ <a>.
9
10 FLAG_PATTERNS = [
11     re.compile(r"BKSEC\{([^\}]+)\}", re.I) # [^}]+ = "mọi ký tự bất kỳ, miễn không phải }", lặp 1+ lần.
12     # re.I = case-insensitive (không phân biệt hoa thường), dù với BKSEC() thường không cần lắm.
13 ]
14
15 NEXT_PAGE_REGEX = re.compile(r"page-\d+[-_A-Za-z0-9]*\.html?&", re.I)
16 # page- literal
17 # \d+ = số trang
18 # [-_A-Za-z0-9]* = phần tên sau số (tùy bài)
19 # \.html? = .htm hoặc .html (vì i là tùy chọn)
20 # $ = kết thúc chuỗi (đảm bảo link kết thúc bằng .htm/.html)
21
22 def find_flag(text: str) -> str | None:
23     for rx in FLAG_PATTERNS:
24         m = rx.search(text)
25         if m:
26             return m.group(0)
27     return None
28 # Nhận vào text (thường là HTML body).
29 # Lặp qua từng regex trong FLAG_PATTERNS.
30 # rx.search(text) tìm match đầu tiên ở bất kỳ vị trí nào.
31 # Nếu match, trả về m.group(0) = toàn bộ chuỗi khớp (ví dụ BKSEC(...)).
32 # Nếu không tìm thấy gì, trả None.
33
34 def extract_next_url(html: str, current_url: str) -> str | None:
35     soup = BeautifulSoup(html, "html.parser")
36     hrefs = []
37     for a in soup.select("a[href]"):
38         href = a.get("href") or ".strip()
39         if not href or href == "#":
40             continue
41         hrefs.append(href)
42
43     if not hrefs:
44         return None
45
46     # ưu tiên chain kiểu page-x
47     for h in hrefs:
48         if NEXT_PAGE_REGEX.search(h):
49             return urljoin(current_url, h)
50
51     # fallback: lấy cái đầu
52     return urljoin(current_url, hrefs[0])
53
54
55 def main():
56     ap = argparse.ArgumentParser(description="CTF Maze crawler (via Burp proxy).")
57     ap.add_argument("--base", default="http://103.77.175.48:8031", help="Base URL")
58     ap.add_argument("--start", default="/pages/page-1-tombstone.html", help="Start path or full URL")
59     ap.add_argument("--max", type=int, default=500, help="Max pages")
60     ap.add_argument("--delay", type=float, default=0.05, help="Delay between requests (seconds)")
61     ap.add_argument("--proxy", default="http://127.0.0.1:8080", help="Burp proxy URL")
62     ap.add_argument("--insecure", action="store_true", help="Disable TLS verification (useful if HTTPS through Burp)")
63     args = ap.parse_args()
64
65     # Burp proxy
66     proxies = {"http": args.proxy, "https": args.proxy}
67
68     session = requests.Session()
69     session.headers.update({
70         "User-Agent": "Mozilla/5.0 (CTF Maze Crawler via Burp)",
71         "Accept": "text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8",
72         "Connection": "close",
73     })
74
75     # Start URL: nếu user đưa full URL thì dùng luôn, không thì join với base
76     if args.start.startswith("http://") or args.start.startswith("https://"):
77         url = args.start
78     else:
79         url = urljoin(args.base, args.start)
80
81     visited = set()
82
83     for step in range(1, args.max + 1):
84         if url in visited:
85             print(f"[!] Loop detected: {url}")
86             return
87         visited.add(url)
88
89         r = session.get(
90             url,
91             timeout=20,
92             allow_redirects=True,
93             proxies=proxies,
94             verify=(not args.insecure),
95         )
96
97         print(f"[step:03] {r.status_code} {url}")
98
99         # scan flag in body
100         flag = find_flag(r.text)
101         if flag:
102             print(f"\n[+] FLAG FOUND: {flag}")
103             print(f"[+] URL: {url}")
104             return
105
106         # keyword hint (không phụ thuộc find_flag)
107         if re.search(r"\bflag\b", r.text, re.I):
108             print("\n[!] Found keyword 'flag' in body. Inspect this page.")
109             print(f"[+] URL: {url}")
110             return
111
112         next_url = extract_next_url(r.text, url)
113         if not next_url:
114             print("[!] No next <a href> found. Stop.")
115             print(f"[!] Last URL: {url}")
116             return
117
118         # stay in-scope
119         if urlparse(next_url).netloc != urlparse(args.base).netloc:
120             print(f"[!] Out of scope next URL: {next_url}")
121             return
122
123         url = next_url
124         time.sleep(args.delay)
125
126     print("[~] Reached max steps without finding flag.")
127
128
129 if __name__ == "__main__":
130     main()

```

1) `find_flag(text)` Hàm này dùng để tìm *flag* trong nội dung phản hồi (thường là HTML response body). Cụ thể, hàm sẽ duyệt qua danh sách các biểu thức chính quy trong biến `FLAG_PATTERNS`, sau đó dùng `search()` để tìm chuỗi khớp đầu tiên. Nếu tìm thấy, hàm trả về toàn bộ chuỗi *flag* (ví dụ `BKSEC{...}`); nếu không tìm thấy thì trả về `None`.

2) `extract_next_url(html, current_url)` Hàm này dùng để trích xuất URL của trang kế tiếp từ mã HTML hiện tại. Đầu tiên, hàm parse HTML bằng `BeautifulSoup`, sau đó lấy toàn bộ các thẻ `<a>` có thuộc tính `href`. Tiếp theo:

- Bỏ qua các `href` rỗng hoặc chỉ là `#`.
- Ưu tiên chọn liên kết khớp mẫu `page-{số}-*.html` bằng regex `NEXT_PAGE_REGEX`.
- Nếu không có link nào khớp mẫu ưu tiên, hàm dùng *fallback* là lấy liên kết đầu tiên.

Cuối cùng, hàm dùng `urljoin()` để chuyển liên kết tương đối thành URL đầy đủ dựa trên `current_url`.

3) `main()` Đây là hàm điều phối toàn bộ quá trình crawl. Các nhiệm vụ chính:

- Khởi tạo các tham số dòng lệnh bằng `argparse` (base URL, start path, số bước tối đa, delay, proxy Burp, v.v.).
- Tạo `requests.Session()` và thiết lập các *HTTP Header* cơ bản (đặc biệt là *User-Agent*).
- Xác định URL bắt đầu: nếu `-start` là URL đầy đủ thì dùng luôn, nếu không thì ghép với `-base`.
- Dùng tập `visited` để tránh lặp vô hạn nếu gặp vòng lặp URL.
- Trong mỗi vòng lặp:
 - Gửi request đến trang hiện tại (qua Burp proxy nếu có).
 - In ra trạng thái HTTP `status code` và URL.
 - Gọi `find_flag()` để tìm flag trong response body.
 - Nếu chưa có flag, gọi `extract_next_url()` để lấy URL kế tiếp.
 - Kiểm tra *scope* bằng `urlparse()` để đảm bảo crawler không đi ra ngoài domain challenge.
 - Chờ một khoảng `delay` rồi tiếp tục vòng lặp.
- Nếu chạy đến giới hạn `-max` mà chưa thấy flag, script sẽ dừng và in thông báo tương ứng.

4) Khối `if __name__ == "__main__":` Đây là điểm khởi chạy của chương trình Python. Khi file được chạy trực tiếp, Python sẽ gọi `main()`. Cách viết này giúp tách phần “định nghĩa hàm” và phần “thực thi”, đồng thời thuận tiện nếu sau này muốn import file này như một module.

5.1 Chạy file code Python và thành công bước vào round 2

Duc mở *Terminal* tại thư mục chứa file `maze_crawler_burp.py` và chạy lệnh sau:

```
python .\maze_crawler_burp.py -base http://103.77.175.40:8031 -start /pages/page-1-Tombstone.html -max 200 -delay 0.02 -proxy "http://127.0.0.1:8080"
```

Sau khi chạy, Duc thu được kết quả như sau:


```
Windows PowerShell
[076] 200 http://103.77.175.40:8031/pages/page-76-Mowry_Station.html
[077] 200 http://103.77.175.40:8031/pages/page-77-Washington_Camp_Junction.html
[078] 200 http://103.77.175.40:8031/pages/page-78-Duquesne_Junction.html
[079] 200 http://103.77.175.40:8031/pages/page-79-Lochiel_Station.html
[080] 200 http://103.77.175.40:8031/pages/page-80-Sunnyside_Junction.html
[081] 200 http://103.77.175.40:8031/pages/page-81-Patagonia_Mountains.html
[082] 200 http://103.77.175.40:8031/pages/page-82-Santa_Rita_Mountains.html
[083] 200 http://103.77.175.40:8031/pages/page-83-Huachuca_Mountains.html
[084] 200 http://103.77.175.40:8031/pages/page-84-Dragoon_Mountains.html
[085] 200 http://103.77.175.40:8031/pages/page-85-Whetstone_Mountains.html
[086] 200 http://103.77.175.40:8031/pages/page-86-Mustang_Mountains.html
[087] 200 http://103.77.175.40:8031/pages/page-87-Empire_Mountains.html
[088] 200 http://103.77.175.40:8031/pages/page-88-Santa_Catalina_Mountains.html
[089] 200 http://103.77.175.40:8031/pages/page-89-Rincon_Mountains.html
[090] 200 http://103.77.175.40:8031/pages/page-90-Tucson_Mountains.html
[091] 200 http://103.77.175.40:8031/pages/page-91-Tortolita_Mountains.html
[092] 200 http://103.77.175.40:8031/pages/page-92-Silver_Bell_Mountains.html
[093] 200 http://103.77.175.40:8031/pages/page-93-Waterman_Mountains.html
[094] 200 http://103.77.175.40:8031/pages/page-94-Roskrige_Mountains.html
[095] 200 http://103.77.175.40:8031/pages/page-95-Tucson_Basin.html
[096] 200 http://103.77.175.40:8031/pages/page-96-Avra_Valley.html
[097] 200 http://103.77.175.40:8031/pages/page-97-Altar_Valley.html
[098] 200 http://103.77.175.40:8031/pages/page-98-San_Rafael_Valley.html
[099] 200 http://103.77.175.40:8031/pages/page-99-Sonoita_Valley.html
[100] 200 http://103.77.175.40:8031/pages/page-100-Sulphur_Springs_Valley.html
[101] 200 http://103.77.175.40:8031/pages/page-101-qncdl1248dbsl.html
[102] 200 http://103.77.175.40:8031/round2_qncdl1248dbsl/index.html
[103] 200 http://103.77.175.40:8031/round2_qncdl1248dbsl/giftsgssul/index.html
[!] Loop detected: http://103.77.175.40:8031/round2_qncdl1248dbsl/index.html
PS D:\CTF Writeup\Web Ex\Maze Maze>
```

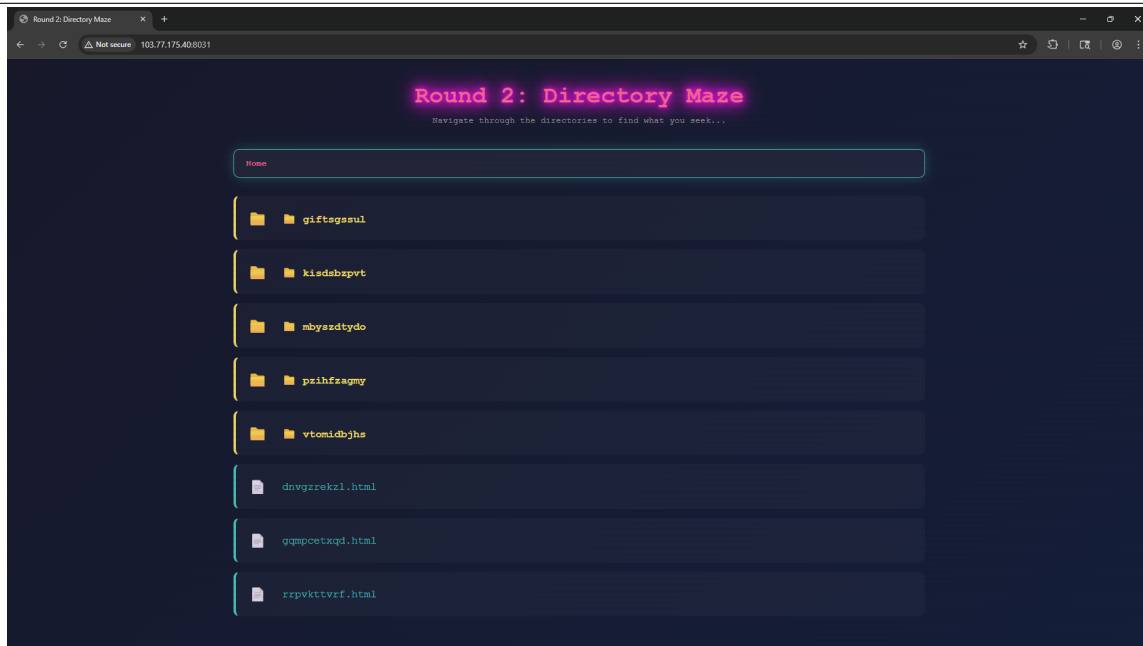
Hình 6: Kết quả chạy script maze_crawler_burp.py trong Terminal.

Dòng cuối của kết quả cho biết chương trình đã bị *loop detected*, đồng thời trả ra một đường link mới:

```
http://103.77.175.40:8031/round2_qncdl1248dbsl/index.html
```

Điều này cho thấy script đã lần theo maze thành công và chạm tới một nhánh mới của challenge.

Dục tiếp tục ném đường link này vào *Burp Suite* (hoặc mở trực tiếp trên trình duyệt) và thu được giao diện của **Round 2** như sau:



Hình 7: Giao diện Round 2 – Directory Maze sau khi lần theo maze thành công.

6 Giải mã Round 2

Theo như trang web gợi ý thì round này Duc cần phải đi vào từng URL có trong Round 2, sau đó kiểm tra phần *Response* để tìm *flag*. Nếu chưa thấy *flag* thì tiếp tục đi vào URL tiếp theo, lặp lại cho đến khi đã duyệt hết toàn bộ các URL có thể đi được.

Duc quyết định dùng **ChatGPT** để viết một file Python tự động thực hiện công việc này. Ở lần đầu, Duc chọn hướng duyệt **DFS** (*Depth-First Search*) để có thể lần theo cây thư mục và đi sâu vào từng nhánh.

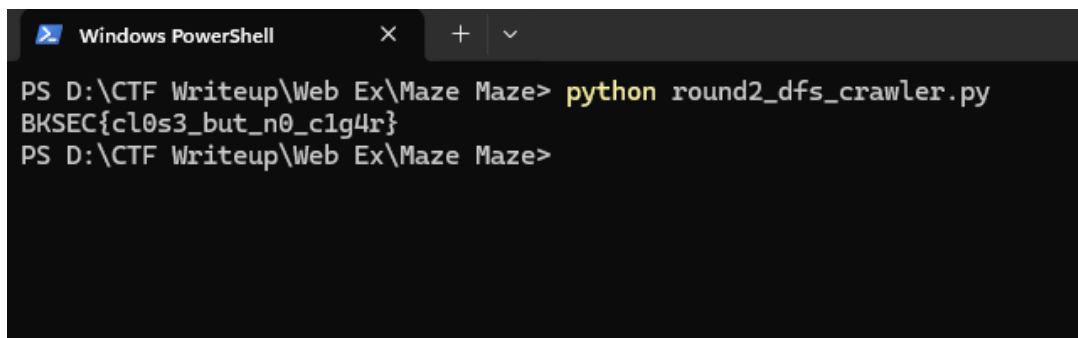
Và đây là đoạn code *Python* mà ChatGPT đã viết:

```

1 import re
2 import sys
3 import time
4 from urllib.parse import urljoin, urlparse
5
6 import requests
7 from bs4 import BeautifulSoup
8
9 START_URL = "http://103.77.175.40:8031/round2_qncd11248dbs1/"
10 TIMEOUT_SEC = 20
11 DELAY_SEC = 0.0 # increase (e.g., 0.02) if you want to be gentler
12 MAX_PAGES = 200000 # safety cap
13
14 FLAG_RE = re.compile(r"BKSEC\[([^\]]+)\]")
15
16 SKIP_SCHEMES = ("javascript:", "mailto:", "tel:")
17 SKIP_HREFS = {"", "#", "http://", "https://", "mailto:", "tel:"}
18
19
20 def extract_links(html: str, base_url: str) -> list[str]:
21     """
22     Extract absolute URLs from <a href="..."> in HTML.
23     Handles directory listings too.
24     """
25     soup = BeautifulSoup(html, "html.parser")
26     links: list[str] = []
27
28     for a in soup.select("a[href]"):
29         href = (a.get("href") or "").strip()
30         if not href or href in SKIP_HREFS:
31             continue
32         if href.lower().startswith(SKIP_SCHEMES):
33             continue
34
35         abs_url = urljoin(base_url, href)
36         # Drop fragment
37         abs_url = abs_url.split("#", 1)[0]
38         links.append(abs_url)
39
40     # De-duplicate while preserving order
41     seen = set()
42     uniq = []
43     for u in links:
44         if u not in seen:
45             seen.add(u)
46             uniq.append(u)
47     return uniq
48
49
50 def in_scope(url: str, start_netloc: str) -> bool:
51     p = urlparse(url)
52     return p.scheme in ("http", "https") and p.netloc == start_netloc
53
54
55 def main() -> int:
56     start_netloc = urlparse(START_URL).netloc
57     session = requests.Session()
58     session.headers.update(
59         {
60             "User-Agent": "Mozilla/5.0 (CTF DFS Crawler)",
61             "Accept": "/*/*",
62             "Connection": "close",
63         }
64     )
65
66     visited: set[str] = set()
67     stack: list[str] = [START_URL]
68     pages_fetched = 0
69
70     while stack:
71         url = stack.pop() # DFS
72         if url in visited:
73             continue
74         if not in_scope(url, start_netloc):
75             continue
76         visited.add(url)
77
78         try:
79             r = session.get(url, timeout=TIMEOUT_SEC, allow_redirects=True)
80         except requests.RequestException:
81             continue
82
83         pages_fetched += 1
84         if pages_fetched > MAX_PAGES:
85             print("Not found flag")
86             return 1
87
88         text = r.text if isinstance(r.text, str) else ""
89         m = FLAG_RE.search(text)
90         if m:
91             print(m.group(0))
92             return 0

```

Duc chạy file `round2_dfs_crawler.py` và thu được kết quả như sau:



```
Windows PowerShell
PS D:\CTF Writeup\Web Ex\Maze Maze> python round2_dfs_crawler.py
BKSEC{cl0s3_but_n0_clg4r}
PS D:\CTF Writeup\Web Ex\Maze Maze>
```

Hình 9: Kết quả chạy lần đầu: script tìm được một flag và in ra màn hình.

Tuy nhiên, khi Duc nhập flag mà chương trình trả về thì hệ thống báo *Incorrect*. Điều này cho thấy challenge này có **nhiều flag gây nhiễu**, chứ không phải chỉ có một flag duy nhất.

Vì vậy, Duc tiếp tục nhờ **ChatGPT** viết lại file Python. Điểm khác biệt là ở phiên bản mới, chương trình sẽ **không dừng ngay khi gặp flag đầu tiên**, mà sẽ tiếp tục duyệt hết các URL và **thu thập tất cả flag** xuất hiện, sau đó in ra *output* theo từng dòng.

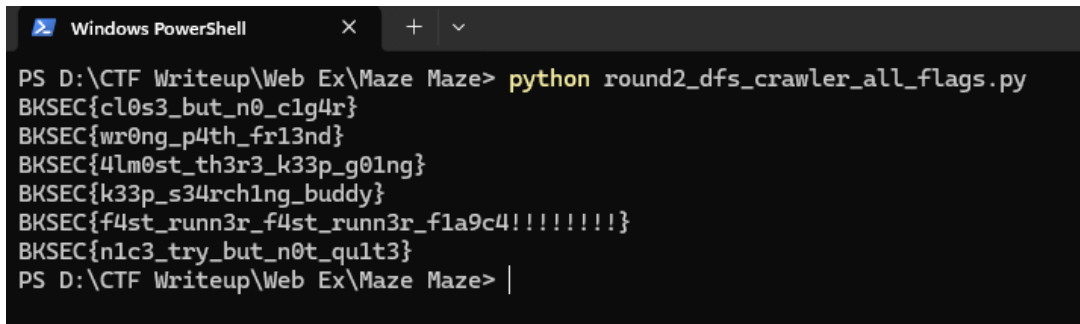
Và đây là file Python mới:

```

1 import re
2 import sys
3 import time
4 from urllib.parse import urljoin, urlparse
5
6 import requests
7 from bs4 import BeautifulSoup
8
9 START_URL = "http://103.77.175.40:8031/round2_qncdl1248dbs1/"
10 TIMEOUT_SEC = 20
11 DELAY_SEC = 0.0 # set e.g. 0.02 to be gentler
12 MAX_PAGES = 300000 # safety cap
13
14 FLAG_RE = re.compile(r"BKSEC\[([^\]]+)\]")
15
16 SKIP_SCHEMES = ("javascript:", "mailto:", "tel:")
17 SKIP_HREFS = {"", "#", "http://", "https://", "mailto:", "tel:"}
18
19
20 def extract_links(html: str, base_url: str) -> list[str]:
21     """Extract absolute URLs from <a href=...> in HTML (directory listings included)."""
22     soup = BeautifulSoup(html, "html.parser")
23     links: list[str] = []
24
25     for a in soup.select("a[href]"):
26         href = (a.get("href") or "").strip()
27         if not href or href in SKIP_HREFS:
28             continue
29         if href.lower().startswith(SKIP_SCHEMES):
30             continue
31
32         abs_url = urljoin(base_url, href).split("#", 1)[0]
33         links.append(abs_url)
34
35     # De-duplicate while preserving order
36     seen = set()
37     uniq = []
38     for u in links:
39         if u not in seen:
40             seen.add(u)
41             uniq.append(u)
42     return uniq
43
44
45 def in_scope(url: str, start_netloc: str) -> bool:
46     p = urlparse(url)
47     return p.scheme in ("http", "https") and p.netloc == start_netloc
48
49
50 def main() -> int:
51     start_netloc = urlparse(START_URL).netloc
52
53     session = requests.Session()
54     session.headers.update(
55         {
56             "User-Agent": "Mozilla/5.0 (CTF DFS Crawler - collect flags)",
57             "Accept": "*/*",
58             "Connection": "close",
59         }
60     )
61
62     visited: set[str] = set()
63     stack: list[str] = [START_URL] # DFS stack
64     pages_fetched = 0
65
66     found_flags: set[str] = set()
67     found_order: list[str] = [] # keep discovery order for nicer output
68
69     while stack:
70         url = stack.pop()
71         if url in visited:
72             continue
73         if not in_scope(url, start_netloc):
74             continue
75         visited.add(url)
76
77         try:
78             r = session.get(url, timeout=TIMEOUT_SEC, allow_redirects=True)
79         except requests.RequestException:
80             continue
81
82         pages_fetched += 1
83         if pages_fetched > MAX_PAGES:
84             break
85
86         text = r.text if isinstance(r.text, str) else ""
87
88         # Collect ALL flags in this page (not just first)
89         for m in FLAG_RE.finditer(text):
90             flag = m.group(0)
91             if flag not in found_flags:
92                 found_flags.add(flag)
93                 found_order.append(flag)
94
95         # Only parse links if it looks like HTML
96         ct = (r.headers.get("Content-Type") or "").lower()
97         if ("text/html" in ct) or ("html" in text.lower()):
98             links = extract_links(text, url)
99
100         # Reverse to keep a more natural DFS order
101         for nxt in reversed(links):
102             if in_scope(nxt, start_netloc) and nxt not in visited:

```

Sau khi chạy file Python mới, Duc thu được kết quả như sau:



```
Windows PowerShell
PS D:\CTF Writeup\Web Ex\Maze Maze> python round2_dfs_crawler_all_flags.py
BKSEC{c10s3_but_n0_c1g4r}
BKSEC{wr0ng_p4th_fr13nd}
BKSEC{4lm0st_th3r3_k33p_g01ng}
BKSEC{k33p_s34rch1ng_buddy}
BKSEC{f4st_runn3r_f4st_runn3r_f1a9c4!!!!!!!!}
BKSEC{n1c3_try_but_n0t_qu1t3}
PS D:\CTF Writeup\Web Ex\Maze Maze> |
```

Hình 11: Kết quả chạy phiên bản DFS thu thập toàn bộ flag trong Round 2.

Duc nhập thử từng flag được in ra, và cuối cùng xác định được **flag đúng** của bài này là:

BKSEC{f4st_runn3r_f4st_runn3r_f1a9c4!!!!!!!!}

Problem solved!