

404页面识别

前言

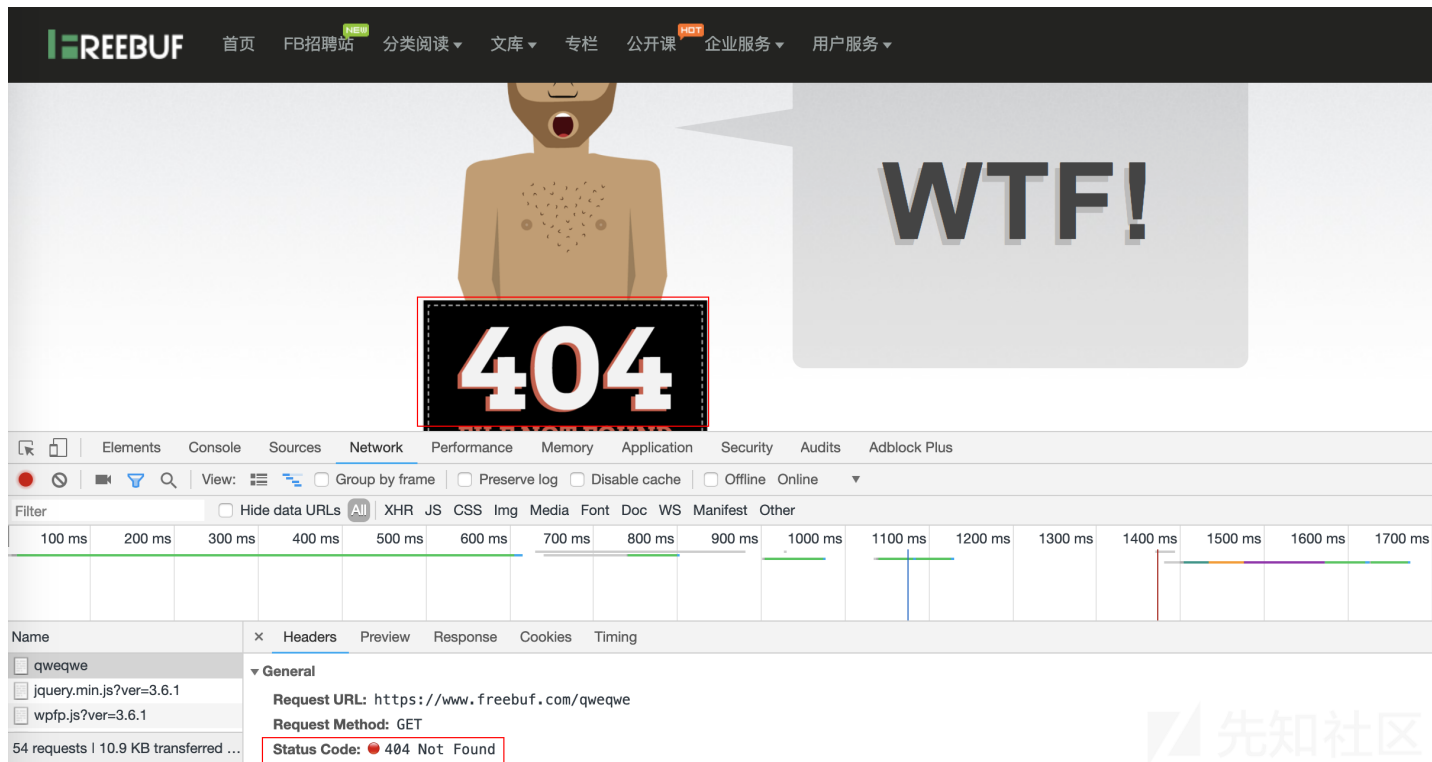
最近在挖洞前做资产收集的时跑了一波子域名，但是目前很多子域名挖掘机挖出来资产还是存在很多水分，准确率一般；于是乎写了个脚本用是否能请求成功作为筛选条件

404页面识别思路

在编写之前先看一下目前网站404的呈现有哪些方式

1. web容器设置404错误页面，服务端返回404状态码

例如freebuf，当我们随便访问一个[不存在的页面](#)，返回的404页面，此时状态码返回了404。



这种情况下，可以根据返回response的状态码来直接判断是不是404页面。

1. 将404错误页面指向一个新的页面，页面上显示404信息，但是此时状态码并不为404，一般返回状态码有301、302，或者直接返回状态码为200的错误页面。

比如Baidu，当我们访问一个[不存在的页面](#)，会返回<https://www.baidu.com/search/error>。

很抱歉，您要访问的页面不存在！

温馨提示：

请检查您访问的网址是否正确

如果您不能确认访问的网址，请浏览[百度更多](#)页面查看更多网址。

[回到顶部](#)重新发起搜索

如有任何意见或建议，请及时[反馈给我们](#)。



根据以上这2种情况，大致得出404页面的识别思路。

- 从状态码是否为404判断
- 获取域名的404页面，然后判断请求的页面和404页面是否相似，相似则可以判断为404页面。

404页面过滤

页面相似度

根据上面提到的识别思路，首先要解决的第一个问题是■■■■■的判断。如何判断两个页面的相似度呢？

这里使用hashes.simhash，对两个页面的body计算hash值，再调用similarity获取两个页面的相似值，自定义一个阈值作为标准判断是否相似，radio可以根据具体情

```
from hashes.simhash import simhash

def is_similar_page(res1, res2, radio=0.85):
    if res1 is None or res2 is None:
        return False

    body1 = res1.body
    body2 = res2.body

    url1 = res1.get_url()
    url2 = res2.get_url()

    simhash1 = simhash(body1.decode('utf-8'))
    simhash2 = simhash(body2.decode('utf-8'))

    calc_radio = simhash1.similarity(simhash2)
    # print("[%s]■[%s]■■■■■■■■■■: %s" % (url1, url2, calc_radio))
    if calc_radio >= radio:
        return True
    else:
        return False
```

构造404页面

这里很简单，访问一个随机生成字符串为后缀的页面。

```
def generate_404_kb(self, url):
    # ■■URL■■■■■
    domain = url.get_domain()          #www.freebuf.com
    domain_path = url.get_domain_path() #https://www.freebuf.com
    rand_file = rand_letters(8) + '.html'
```

整体思路

荒废的域名

比如直接访问<http://e.apptaxi.com.cn>，会发生302跳转到定义好的404页面https://img-ys011.didistatic.com/static/dfc_default_page/index.html，但是直接访问https://img-ys011.didistatic.com/static/dfc_default_page/index.html

1. 6 条回复



[hpdoger](#) 2019-03-18 16:06:13

这个hashes模块Pip下不能安装么？

0 回复Ta



[Bojack](#) 2019-03-18 16:10:59

[@hpd****](#) 包名python-hashes，如果实在下载不了我也可以发你，你留个邮箱

0 回复Ta



[Bojack](#) 2019-03-18 16:30:21

页面相似度的计算还能参考 <https://thief.one/2018/04/12/1/>

0 回复Ta



[zhangzhongnan](#) 2019-03-19 15:07:10

非常巧合的是这篇帖子的ID是以404结尾的

2 回复Ta



[张喵喵是小仙女](#) 2019-04-11 15:35:26

小哥哥，源码可以发一下嘛~

0 回复Ta



[Passer6y](#) 2019-05-05 00:10:28

[@hpdoger](#) 踩了一遍坑，这个库在16年爆了个issue，到现在也没修... 然后有个人fork了他的仓库，帮他把洞修好了2333，安装 `pip install changanya`
导入：`from changanya.simhash import Simhash`

0 回复Ta

[登录](#) 后跟帖

先知社区

[现在登录](#)

热门节点

[技术文章](#)

[社区小黑板](#)

目录

[RSS](#) [关于社区](#) [友情链接](#) [社区小黑板](#)