

# 综合训练-抓取91看剧的视频

首先, 我们先从一个简单的案例入手.

我们想要抓取网上的视频资源就必须要了解我们的视频网站是如何工作的. 这里我用91看剧来做举例. 其他网站的原理是一样的.

## 1.视频网站是如何工作的

假设, 你现在想要做一个视频网站. 也有很多的UP猪帮你上传视频. OK, 作为服务器作者的你. 只需要把视频保存起来. 然后给出一个视频的链接即可. 然后在你的HTML代码中通过 video 标签引入即可.

4_2_多线程.mp4	2021/3/3 下午 4:34	704.7 MB	Monosnap vide
4_1_第四章概述.mp4	2021/3/2 下午 7:35	192.3 MB	Monosnap vide
1_7_补充_关闭resp.mp4	2021/3/2 下午 5:52	169.4 MB	Monosnap vide
3_5_综合训练_抓取网易云音乐评论信息.mp4	2021/3/2 下午 5:43	4.82 GB	Monosnap vide
3_3_防盗链_抓取梨视频.mp4	2021/3/1 下午 9:02	1.89 GB	Monosnap vide
3_4_代理.mp4	2021/3/1 下午 5:56	544.6 MB	Monosnap vide
3_2_处理 cookie_登录小说网.mp4	2021/2/26 下午 6:51	1.23 GB	Monosnap vide
3_1_requests进阶概述.mp4	2021/2/26 下午 6:20	116.7 MB	Monosnap vide
2_9_xpath实战_抓取猪八戒网信息.mp4	2021/2/26 下午 4:19	2.57 GB	Monosnap vide
2_8_xpath入门_02.mp4	2021/2/25 下午 7:23	959.7 MB	Monosnap vide
2_8_xpath入门_01.mp4	2021/2/25 下午 6:07	876.4 MB	Monosnap vide
2_7_bs4解析案例-抓取优美图库图片.mp4	2021/2/25 下午 4:59	2.03 GB	Monosnap vide
2_6_bs4解析入门-搞搞菜价.mp4	2021/2/24 下午 7:24	1.42 GB	Monosnap vide
2_5_bs4解析前戏-html语法规则.mp4	2021/2/24 下午 6:48	549.9 MB	Monosnap vide
2_4_屠戮盗版天堂电影信息.mp4	2021/2/24 下午 5:16	2.94 GB	Monosnap vide
2_4_手刃豆瓣 TOP250 电影排行.mp4	2021/2/24 下午 4:26	1.16 GB	Monosnap vide
2_3_python的re模块使用.mp4	2021/2/23 下午 8:09	1.58 GB	Monosnap vide
2_2_re解析_正则表达式_02.mp4	2021/2/23 下午 6:16	763 MB	Monosnap vide
2_2_re解析_正则表达式_01.mp4	2021/2/23 下午 5:45	967.9 MB	Monosnap vide
2_1_数据解析概述.mp4	2021/2/23 下午 5:19	171.2 MB	Monosnap vide
1_6_3_requests入门.mp4	2021/2/23 下午 4:51	1.1 GB	Monosnap vide
1_6_2_requests入门.mp4	2021/2/23 下午 4:29	373.7 MB	Monosnap vide
1_6_1_requests入门.mp4	2021/2/23 下午 4:12	1.18 GB	Monosnap vide
1_5_http协议.mp4	2021/2/22 下午 4:45	753.6 MB	Monosnap vide
1_4_web请求过程剖析.mp4	2021/2/20 下午 6:45	1.22 GB	Monosnap vide
1_3_手刃一个小爬虫.mp4	2021/2/20 下午 4:03	872.2 MB	Monosnap vide
1_2_本课程使用的软件.mp4	2021/2/19 下午 4:39	610.3 MB	Monosnap vide
1_1_爬虫概述.mp4	2021/2/19 下午 3:36	492.6 MB	Monosnap vide

1 <video src="1\_1\_爬虫概述.mp4"></video>

就可以了. 但是, 如果你这么做. 你的用户和老板一定会把你骂的狗血临头. 为什么呢?

假设你的视频是10个G的高清无码大资源. 那么此时, 你的用户和你老板将面临如下困境

1. 用户: 这个视频怎么加载的这么慢. 点击快进也快进不了. 太慢了. 塔喵的烦死了.
2. 老板: 怎么这个月的流量费又这么高啊. 要死的拉好不~

为什么会这样？聪明的我告诉你答案。你的视频那么大。每次用户打开的时候。可能只是差了最后几分钟没看呢。那此时它必须把整个视频都传输完毕。才能看到他想看的那里。等待时间肯定超长的好不。而每次都要把10G的文件进行网络传输。流量费~你懂的。三大运营商最喜欢的就是你这种朴实无华的送钱行为。

OK~ 不扯了。但凡有点儿经验的程序员肯定会想办法把用户上传好的视频进行转码(不同清晰度)做切片(ts)处理。这样既方便用户进行大跨度的调整进度条(最小延迟)。也能为公司节省大量的流量费。

既然要把视频切成非常多个小碎片。那就需要有个文件来记录这些小碎片的路径。该文件一般为M3U文件。M3U文件中的内容经过UTF-8的编码后, 就是M3U8文件。今天, 我们看到的各大视频网站平台使用的几乎都是M3U8文件。

如何解读M3U8文件。

```
#EXTM3U
#EXT-X-VERSION:3
#EXT-X-TARGETDURATION:13
#EXT-X-MEDIA-SEQUENCE:0
#EXT-X-KEY:METHOD=AES-128,URI="key.key"
#EXTINF:12.600000,
cFN8o3436000.ts
#EXTINF:10.000000,
cFN8o3436001.ts
#EXTINF:10.000000,
cFN8o3436002.ts
#EXTINF:10.000000,
cFN8o3436003.ts
#EXTINF:10.000000,
cFN8o3436004.ts
#EXTINF:10.000000,
cFN8o3436005.ts
#EXTINF:6.880000,
cFN8o3436006.ts
```

每个视频切片最大时长。

切片文件的加密方式以及加密的密钥地址  
如果有加密, 需要先解密才能播放

持续时间

这里面不带#开头的就是每个ts文件的地址

基本知道这些就够了.

## 2. 先来个简单的试试

目标: 哲仁皇后. 走起

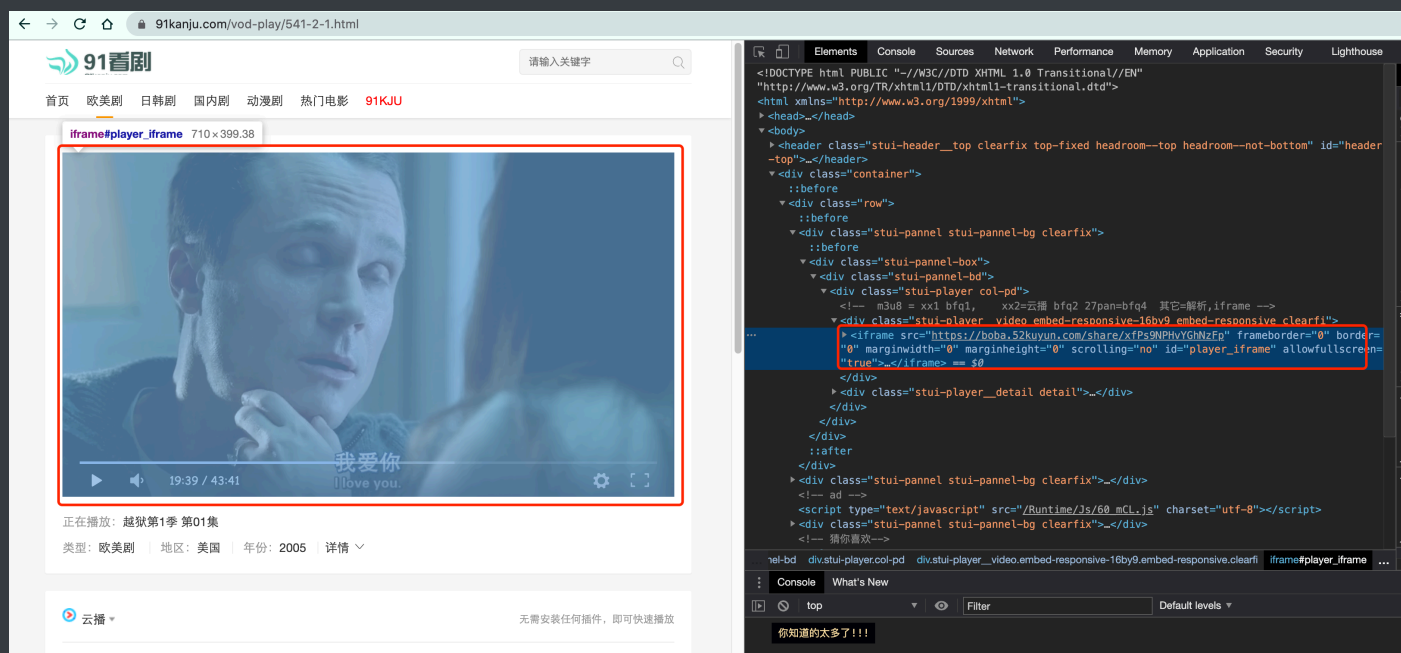
```
1 import requests
2 import re
3
4 headers = {
5     "user-agent": "Mozilla/5.0 (Macintosh; Intel
    Mac OS X 10_15_4) AppleWebKit/537.36 (KHTML, like
    Gecko) Chrome/88.0.4324.192 Safari/537.36"
6 }
7
8 url = "https://www.91kanju.com/vod-play/54812-1-
    1.html"
9 resp = requests.get(url, headers=headers)
10
11 obj = re.compile(r"url: '(?P<url>.*?)'", re.S)
12 result = obj.search(resp.text)
13 resp.close()
14
15 url = result.group('url')
```

```
16 resp2 = requests.get(url, headers=headers)
17 f = open("哲仁王后.m3u8", mode="wb")
18 f.write(resp2.content)
19 resp2.close()
20 f.close()
21
22 # 读取m3u8视频文件
23 with open("哲仁皇后.m3u8", mode="r",
    encoding="utf-8") as f:
24     count = 1
25     for line in f:
26         line = line.strip()
27         if line.startswith("#"):
28             continue
29
30         resp3 = requests.get(line,
    headers=headers)
31         with open(f"video2/{count}.ts",
    mode="wb") as ff:
32             ff.write(resp3.content)
33             resp3.close()
34             count += 1
35             print("完事儿1个")
36             if count >= 10: # 测试. 别整太过分
37                 break
38
```

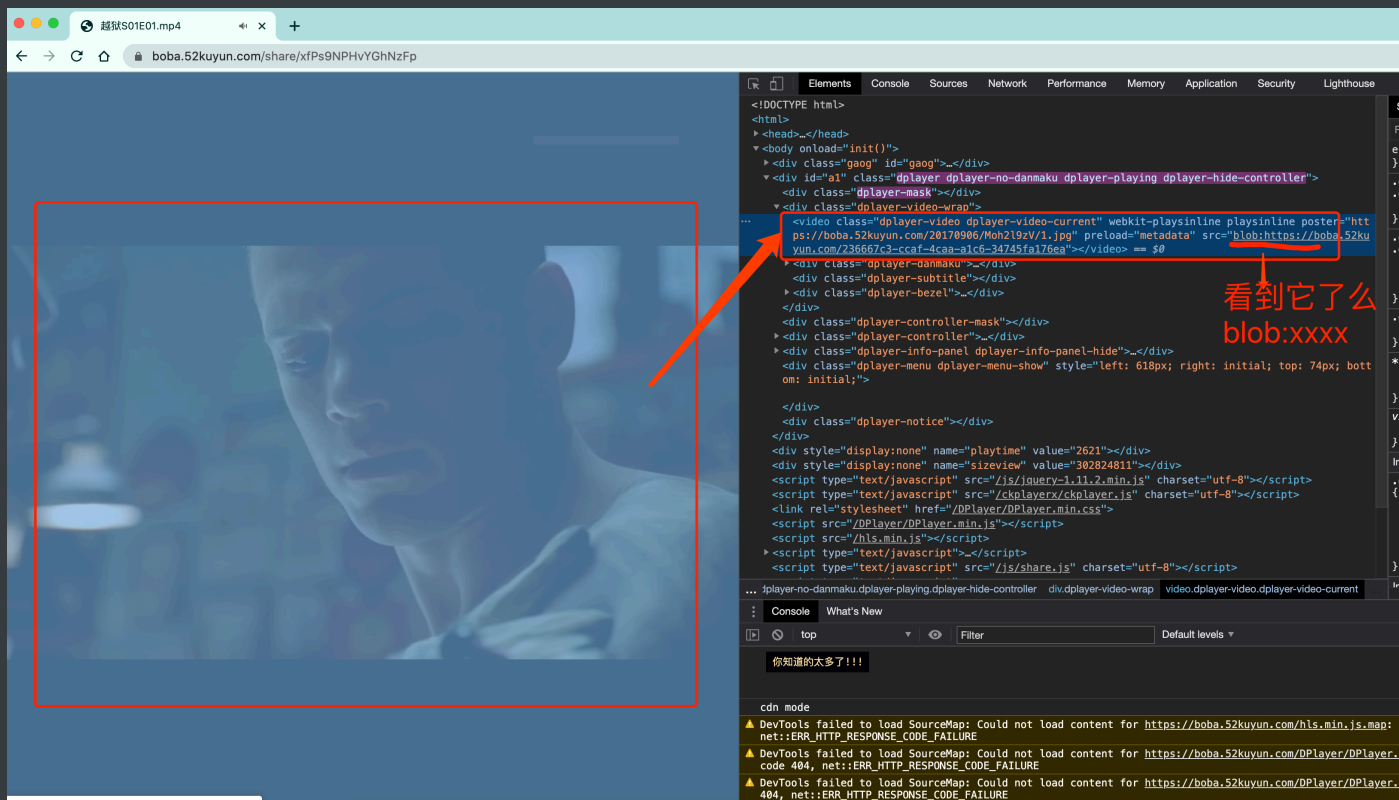
- 1 # 最后，需要用quick time把视频拼接起来。
- 2 # 看到了么。上面那个居然没有用协程，也没有多线程。为什么呢？因为太简单了。我们来个复杂的。

### 3. 真正的挑战. 抓取<越狱>视频第一集

首先, 盯紧我们的目标, 你会发现, 我们想要的视频被放在了一个iframe里面



我们打开这个iframe的链接, 发现这个链接里嵌套了一个播放器。



而这个播放器中的视频地址是一个blob:xxxx, 这是个什么鬼.  
blob:https并不是一种协议, 而是html5中blob对象在赋给video标签后生成的一串标记. 说简单点儿. 就是防止爬虫能直接拿到视频下载地址, 中间倒了一手而已.

此路不通, 我们再去观察一下这个页面的页面源代码.

```

19     </div>
20 </div>
21
22     <div id="a1"></div>
23
24     <div style="display:none" name="playtime" value="2621"></div>
25     <div style="display:none" name="sizeview" value="302824811"></div>
26     <script type="text/javascript" src="/js/jquery-1.11.2.min.js" charset="utf-8"></script>
27     <script type="text/javascript" src="/ckplayerx/ckplayer.js" charset="utf-8"></script>
28     <link rel="stylesheet" href="/DPlayer/DPlayer.min.css">
29     <script src="/DPlayer/DPlayer.min.js"></script>
30     <script src="/hls.min.js"></script>
31     <script type="text/javascript">
32         var video_player= 'dplayer'
33         var tracker_url = ''
34         var signaler_url = ''
35         var hosts = '';
36         var redirecturl = "http://vip.okzybo.com";
37         var videoid = "xfPs9NPHvYGhNzFp";
38
39         var id = 'xfPs9NPHvYGhNzFp'
40         var l = ''
41         var r = ''
42         var t= '15'
43         var d= ''
44         var u= ''
45
46         var main = "/20170906/Moh2l9zV/index.m3u8?sign=548ae366a075f0f9e7c76af215aa18e1";
47         var xml = "/20170906/Moh2l9zV/index.xml?sign=548ae366a075f0f9e7c76af215aa18e1";
48         var pic = "/20170906/Moh2l9zV/1.jpg";
49         var thumbnails = "/20170906/Moh2l9zV/thumbnails.jpg";
50     </script>
51     <script type="text/javascript" src="/js/share.js" charset="utf-8"></script>
52     <script type="text/javascript">
53
54     </script>
55
56     <!--广告添加区域-->
57 </body>
58
59 </html>

```

抓取这个视频的整体思路:

- 1 拿到主页的页面源代码. 获取到iframe的链接
- 2 拿到iframe链接的页面源代码. 获取到m3u8链接
- 3 下载m3u8文件(2层), 稍显繁琐.
- 4 对文件进行检索, 拿到小视频的链接, 进行异步下载
- 5 解密合并小的切片文件还原为大的MP4文件进行播放

并且, 我们在浏览器的network中也能非常清晰的看到整个浏览器的请求过程, 几乎和我们刚才总结的思路是一致的 接下来, 上代码吧. 我依然会分步骤的去描述如何抓取到最后的视频文件

```

1 import requests
2 from bs4 import BeautifulSoup
3 import re

```



```
4 import asyncio
5 import aiohttp
6 import aiofiles
7 from Crypto.Cipher import AES
8 import os
9
10 def get_iframe_url(url):
11     resp = requests.get(url)
12     main_page = BeautifulSoup(resp.text,
13                               "html.parser")
14     iframe = main_page.find("iframe")
15     resp.close()
16     return iframe.get("src")
17
18 def get_first_m3u8_url(iframe_url):
19     iframe_resp = requests.get(iframe_url)
20     obj = re.compile(r'var main = "(?
P<first_m3u8>.*?)"', re.S)
21     result = obj.search(iframe_resp.text)
22     iframe_resp.close()
23     return result.group('first_m3u8').strip()
24
25
26 def download_m3u8_file(url, name):
27     resp = requests.get(url)
```

```
28     resp.encoding='utf-8'
29     with open(name, mode="w", encoding="utf-8")
    as f:
30         f.write(resp.text)
31     resp.close()
32
33
34 async def download_ts(session, url, name):
35     async with session.get(url) as resp:
36         async with
aiofiles.open(f"video3/{name}", mode="wb") as f:
37             await f.write(await
resp.content.read())
38     print(f"{name}~OK!")
39
40
41 async def aio_download(up_url):
42     tasks = []
43     async with aiofiles.open("第二层
_tmp_m3u8.txt", mode="r", encoding="utf-8") as
f:
44         async with aiohttp.ClientSession() as
session:
45             async for line in f:
46                 if line.startswith("#"):
47                     continue
```



```
69
70     bs = await f1.read()
71     await f2.write(aes.decrypt(bs))
72     print(f"{name}~解密OK~")
73
74
75 async def aio_dec(key):
76     tasks = []
77     async with aiofiles.open("第二层
78     _tmp_m3u8.txt", mode="r", encoding="utf-8") as
79     f:
80         async for line in f:
81             if line.startswith("#"):
82                 continue
83             line = line.strip()
84             tasks.append(dec_ts(line, key))
85             # 异步解密
86             await asyncio.wait(tasks)
87
88 def merge_ts():
89     s = []
90     with open("第二层_tmp_m3u8.txt") as f:
91         for line in f:
92             if line.startswith("#"):
93                 continue
```

```
93         line = line.strip()
94         s.append("./video4/tmp_"+line)
95     names = " ".join(s)
96     os.system(f"cat {names} > movie.mp4")
97
98 # 主程序
99 def main(url):
100     iframe_url = get_iframe_url(url)
101     # 1. iframe的域名
102     iframe_domain = iframe_url.split("/share/")
103     [0]
104     # 2. 拿到m3u8地址(顶级)
105     first_m3u8_url =
106     get_first_m3u8_url(iframe_url)
107     # 3.1 下载first_m3u8
108     first_m3u8_url = iframe_domain +
109     first_m3u8_url
110     download_m3u8_file(first_m3u8_url, "顶层
111     _tmp_m3u8.txt")
112     # 3.2 拿到m3u8的上层url, 目的是拼接第二层m3u8的
113     url以及后面的key.key等等
114     first_m3u8_url_up =
115     first_m3u8_url.rsplit("/", 1)[0]
116     # 3.3 下载第二层m3u8文件
117     with open("顶层_tmp_m3u8.txt", mode="r",
118             encoding='utf-8') as f:
```

```
112         for line in f:
113             if line.startswith("#"):
114                 continue
115             else:
116                 second_m3u8_url =
first_m3u8_url_up + "/" + line.strip()
117
118         download_m3u8_file(second_m3u8_url, "第二层
_tmp_m3u8.txt")
119         # 4. 读取第二层m3u8文件. 开始下载ts
120         second_m3u8_url_up =
second_m3u8_url.rsplit("/", 1)[0]
121         # 异步下载
122
123         asyncio.run(aio_download(second_m3u8_url_up))
124         # 测试的时候慎重啊
125         # 5. 解密与合并
126         # 5.1. 拿到key
127         key = get_key(second_m3u8_url_up +
"/key.key") # 这里就不去m3u8里找了. 偷个懒
128         # print(key)
129         # key = "c5878c26baaaac8c"
130         # 5.2. 解密
131         asyncio.run(aio_dec(key))
132         # 5.3. 合并
133         merge_ts()
```

```
132
133
134 if __name__ == '__main__':
135     url = "https://www.91kanju.com/vod-play/541-
        2-1.html"
136     main(url)
```

备注: 不是所有的视频网站都是这样的. 例如, bilibili, 爱奇艺等, 必须得单独想其他办法. 原因我就不讲了. 稍显复杂.