



白雪爬虫和七个网站的故事

网络信息内容安全作业一

刘迅

2024/04/22



中国科学院大学

University of Chinese Academy of Sciences



目录

1 爬虫实验小结

► 爬虫实验小结

► 创新点与主要贡献

1+X+1 代码框架

可视化流程图分析

► 讨论与展望



实验目标

1 爬虫实验小结

	视频文件	视频简介	播放/点赞量	频道	标题	搜索	备注
ifeng	single	no	yes	no	yes	yes	使用haokan搜索结果
xiaodutv	single	no	yes	no	yes	yes	
thepaper	single	yes	yes	no	yes	yes	
haokan	single	no	yes	no	yes	yes	
ku6	single	no	no	no	yes	no	
cntv	segment	yes	yes	no	yes	yes	
bilibili	single	yes	yes	yes	yes	yes	

图: 目标网站所含功能一览

- 本次实验要求完成 7 个网站的视频内容爬取任务，功能涵盖视频标题、简介、点赞/播放量与频道获取，以及搜索功能，并实现单个/分段视频的下载。
- 由于网站设计差异，不同网站实现的功能有所差异。此处用红框列出与其他网站有较大不同的区别。



实验目标

1 爬虫实验小结

	视频文件	视频简介	播放/点赞量	频道	标题	搜索	备注
ifeng	single	no	yes	no	yes	yes	使用haokan搜索结果
xiaodutv	single	no	yes	no	yes	yes	
thepaper	single	yes	yes	no	yes	yes	
haokan	single	no	yes	no	yes	yes	
ku6	single	no	no	no	yes	no	
cntv	segment	yes	yes	no	yes	yes	
bilibili	single	yes	yes	yes	yes	yes	

图: 目标网站所含功能一览

1. cntv 的视频是分为多个 ts 文件传输, 其他视频网站的视频都是单个媒体文件。
2. ku6 是七个网站中唯一一个没有播放/点赞量的。
3. ku6 还是七个网站中唯一一个没有搜索功能的。
4. xiaodutv 的搜索结果自动跳转到百度, 聚合了 haokan, bilibili, weibo 多家搜索结果内容。为了方便起见, 同时也不影响实验效果, xiaodutv 的搜索使用 haokan 的方法。



实验目标

1 爬虫实验小结

	视频文件	视频简介	播放/点赞量	频道	标题	搜索	备注
ifeng	single	no	yes	no	yes	yes	使用haokan搜索结果
xiaodutv	single	no	yes	no	yes	yes	
thepaper	single	yes	yes	no	yes	yes	
haokan	single	no	yes	no	yes	yes	
ku6	single	no	no	no	yes	no	
cntv	segment	yes	yes	no	yes	yes	
bilibili	single	yes	yes	yes	yes	yes	

图: 目标网站所含功能一览

“频道”项的处理

上课提到了没有频道标签的视频，如果需要分析其频道，应该如何处理，应该是引导大家向侧信道（通过弹幕、评论等信息分析）、聚类算法等角度思考。本次实验受时间所限，仅考虑有显式频道标签的 bilibili 网站，其他网站的频道留空处理。

[illegible]

图: 例子: bilibili 爬取结果整理

¹只要有的话



实验成果

1 爬虫实验小结

$$1+X+1$$

(X=7)

代码框架，更具可拓展性的贡献之处



目录

2 创新点与主要贡献

► 爬虫实验小结

► 创新点与主要贡献
1+X+1 代码框架
可视化流程图分析

► 讨论与展望



目录

2 创新点与主要贡献

- ▶ 爬虫实验小结
- ▶ 创新点与主要贡献
 - 1+X+1 代码框架
 - 可视化流程图分析
- ▶ 讨论与展望



1+X+1 代码框架

2 创新点与主要贡献

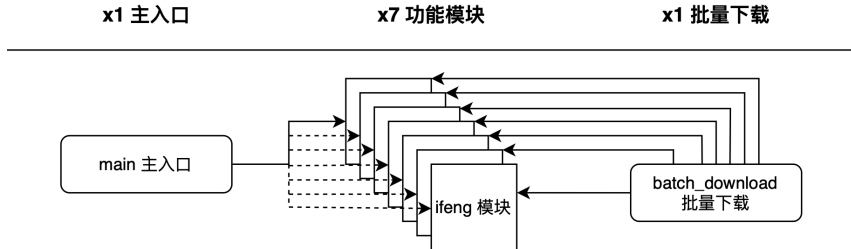


图: 1+X+1 代码框架

项目构建了一套统一的代码框架，用户只需要在 main 中选择目标视频网站，后续操作即可屏蔽不同网站的差异，通过功能一致的模块对外交互。



1+X+1 代码框架：入口

2 创新点与主要贡献

```
src > main.py ...
1 if __name__ == '__main__':
2     print('=====欢迎来到聚合视频下载系统! =====')
3     website = input("请输入视频网站: ")
4     if website == 'haokan':
5         import haokan as website
6     elif website == 'thepaper':
7         import thepaper as website
8     elif website == 'ifeng':
9         import ifeng as website
10    elif website == 'xiaodutv':
11        import xiaodutv as website
12    elif website == 'ku6':
13        import ku6 as website
14    elif website == 'cntv':
15        import cntv as website
16    elif website == 'bilibili':
17        import bilibili as website
18    id = input("请输入视频ID: ")
19    keyword = input("请输入搜索词: ")
20    if id:
21        website.get_video_info(id)
22    elif keyword:
23        website.search_video(keyword)
24    else:
25        print("请输入视频ID或搜索词...")
```

```
=====欢迎来到聚合视频下载系统! =====
请输入视频网站: ifeng
请输入视频ID:
请输入搜索词: 小米
视频URL: https://video19.ifeng.com/video09/2024/04/21/p7187813613816517180-102-222416.mp4
正在保存 ifeng 视频: 吉利、智己、极越包围小米，围剿还是蹭流量? .mp4
成功保存 ifeng 视频: 吉利、智己、极越包围小米，围剿还是蹭流量? .mp4
-----根据 关键词 获取视频信息-----
视频标题是: 吉利、智己、极越包围小米，围剿还是蹭流量?
ifeng 没有简介
88 点赞
ifeng 没有频道
视频存储路径: data/ifeng/吉利、智己、极越包围小米，围剿还是蹭流量? .mp4
ifeng 视频日志存储: logs/ifeng/吉利、智己、极越包围小米，围剿还是蹭流量? .txt
视频URL: https://video19.ifeng.com/video09/2024/04/20/p7187229787327308461-102-080538.mp4
正在保存 ifeng 视频: 苹果全球销量降10%，小米传音却猛增为何? .mp4
成功保存 ifeng 视频: 苹果全球销量降10%，小米传音却猛增为何? .mp4
-----根据 关键词 获取视频信息-----
视频标题是: 苹果全球销量降10%，小米传音却猛增为何?
ifeng 没有简介
1102 点赞
ifeng 没有频道
视频存储路径: data/ifeng/苹果全球销量降10%，小米传音却猛增为何? .mp4
ifeng 视频日志存储: logs/ifeng/苹果全球销量降10%，小米传音却猛增为何? .txt
视频URL: https://video19.ifeng.com/video09/2024/04/16/p7185921135341478844-102-170536.mp4
正在保存 ifeng 视频: 小米SU7细节，对比才有差距#小米su7 #新能源汽车 #小米汽车.mp4
成功保存 ifeng 视频: 小米SU7细节，对比才有差距#小米su7 #新能源汽车 #小米汽车.mp4
-----根据 关键词 获取视频信息-----
视频标题是: 小米SU7细节，对比才有差距#小米su7 #新能源汽车 #小米汽车
ifeng 没有简介
```

图: 运行结果示例

图: 主入口示例



1+X+1 代码框架：功能模块

2 创新点与主要贡献

```
src > structure.py > ...
94 # 获取视频简介
95 def get_video_intro(id):
96     # 输入: 视频 ID string
97     # 输出: 视频简介 string
98     return f"({WEBSITE_NAME}) 没有简介"
99
100 # 获取视频播放量
101 def get_video_play(id):
102     # 输入: 视频 ID string
103     # 输出: 视频播放量和点赞量 string
104     pass
105
106 # 获取视频频道
107 def get_video_channel(id):
108     # 输入: 视频 ID string
109     # 输出: 视频频道 string
110     return f"({WEBSITE_NAME}) 没有频道"
111
112 # 根据当前视频网站决定下载方式
113 def download_video(id):
114     # 输入: 视频 URL string
115     # 输出: 视频文件存储路径 string
116     pass
117
118 # 收集视频基本信息
119 def get_video_info(id, title='ID'):
120     title = f'-----视频 {title} 获取视频信息-----'
121     video_title = get_video_title(id)
122     video_intro = get_video_intro(id)
123     video_play = get_video_play(id)
124     video_chan = get_video_channel(id)
125     video_path = download_video(id)
126     print(title)
127     print(f"视频标题是: {video_title}")
128     print(video_intro)
129     print(video_play)
```

图: 功能模块框架示例-1

```
src > structure.py > ...
36 WEBSITE_NAME = ""
37 VIDEO_URL_PREFIX = ""
38 SEARCH_URL_PREFIX = ""
39 SEARCH_API_URL = ""
40 SEARCH_NUM = 10
41
42 # 根据视频id得到视频url
43 def get_url(id):
44     return VIDEO_URL_PREFIX.format(id)
45
46 # 获取视频文件
47 def download_video_by_url(url, filename=None):
48     # 输入: 视频 URL string; 视频文件名 string (可选)
49     # 输出: 视频文件存储路径 string
50
51     if filename is None:
52         filename = url.split('/')[-1]
53
54     # 如果WEBSITE_NAME目录不存在那么创建WEBSITE_NAME目录
55     if not os.path.exists("data/" + WEBSITE_NAME):
56         os.makedirs("data/" + WEBSITE_NAME)
57     file_path = "data/" + WEBSITE_NAME + "/" + filename
58     with open(file_path, mode="wb") as f:
59         print(f"正在保存({WEBSITE_NAME})视频: {filename}")
60         video_content = requests.get(url=url).content
61         f.write(video_content)
62         print(f"已成功保存({WEBSITE_NAME})视频: {filename}")
63
64     return file_path
65
66 # 分段下载视频
67 def download_video_by_segment(url, filename=None):
68     # 输入: 视频 URL string; 视频文件名 string (可选)
69     # 输出: 视频文件存储路径 string
70
71     if filename is None:
```

图: 功能模块框架示例-2



目录

2 创新点与主要贡献

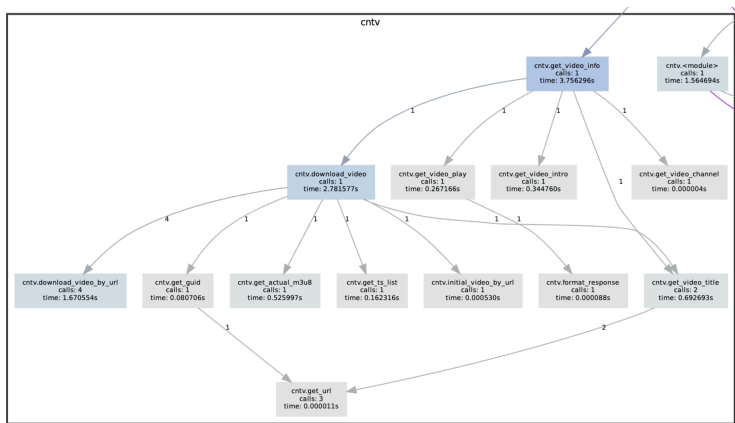
- ▶ 爬虫实验小结
- ▶ 创新点与主要贡献
 - 1+X+1 代码框架
 - 可视化流程图分析
- ▶ 讨论与展望



项目代码流程图分析

2 创新点与主要贡献

使用 pycallgraph 工具绘制代码调用关系图，更好地分析逻辑关系和性能瓶颈。





目录

3 讨论与展望

- ▶ 爬虫实验小结
- ▶ 创新点与主要贡献
 - 1+X+1 代码框架
 - 可视化流程图分析
- ▶ 讨论与展望



亮点与可改进之处

3 讨论与展望

本项目的一大特点在于

- **全程使用 requests 库构建与网站的交互**，相较于使用 selenium 或 playwright，效率上有所优势。

显而易见的，本次项目只是一个初步的尝试，在性能、鲁棒性等角度都还有很大的提升空间，具体来说有

1. 实践过程中，偶尔会因为个别请求超时或中断而影响整个程序运行，可以加入 try-except 提升整体的鲁棒性。

未解之谜

好看视频 cookie 易过期。假如一段时间爬取过于频繁，即便设置 user-agent 与 cookie，也会被拒绝爬取。放着一段时间又好了！



亮点与可改进之处

3 讨论与展望

本项目的一大特点在于

- **全程使用 requests 库构建与网站的交互**，相较于使用 selenium 或 playwright，效率上有所优势。

显而易见的，本次项目只是一个初步的尝试，在性能、鲁棒性等角度都还有很大的提升空间，具体来说有

1. 实践过程中，偶尔会因为个别请求超时或中断而影响整个程序运行，可以加入 try-except 提升整体的鲁棒性。
2. 可以使用异步、高并发的方式提升下载性能。

未解之谜

好看视频 cookie 易过期。假如一段时间爬取过于频繁，即便设置 user-agent 与 cookie，也会被拒绝爬取。放着一段时间又好了！



思考：从爬虫对抗一窥攻防博弈

3 讨论与展望

对各个网站的爬取过程是对攻防实践非常生动的诠释，这是一个螺旋上升的过程。

- 攻击：从静态网页硬编码的 mp4 URL 中提取。
- 防御：动态网页加载 mp4URL；静态网页显示 blob:xxx。
- 攻击：selenium（本次实验未采用）提取渲染后的 HTML 文件；或者观察对 mp4 文件的 js 请求。
- 防御：分段传输视频文件。
- 攻击：通过全局表示 guid 寻找视频片段（凤凰网）。
- 攻击：将 User-Agent 切换到手机查看（B 站）。
- 防御：检查时间戳 timestamp（好看视频）。



思考：屡禁不绝 v.s. 有意为之？

3 讨论与展望

在这些攻防对抗中，看到了很多防御/混淆措施，但是爬虫之所以不能够根治的原因是什么呢？

个人猜测有技术和非技术两方面的因素：

- 技术上：对于访问媒体文件的真实用户，用户本地的浏览器必须拿到媒体资源的 URL，才能够播放。无论传输过程如何加密，由于视频文件的播放这一步骤必须在本地进行，因此总会暴露出媒体文件的地址。



思考：屡禁不绝 v.s. 有意为之？

3 讨论与展望

在这些攻防对抗中，看到了很多防御/混淆措施，但是爬虫之所以不能够根治的原因是什么呢？

个人猜测有技术和非技术两方面的因素：

- 技术上：对于访问媒体文件的真实用户，用户本地的浏览器必须拿到媒体资源的 URL，才能够播放。无论传输过程如何加密，由于视频文件的播放这一步骤必须在本地进行，因此总会暴露出媒体文件的地址。
- 非技术：视频网站需要被搜索引擎收录、被互联网用户使用，倘若反爬虫做到极致，也就与内网无异了。



致谢

3 讨论与展望

感谢杨杼鑫、秦子茗同学的交流、讨论与指导！

感谢朱老师不那么 push 的 DDL



白雪爬虫和七个网站的故事

大作业完结撒花！

任何问题？