# Crawler Technology Based on Scrapy Framework

Wu Hejing

East University of Heilongjiang

Heilongjiang, China

e-mail: 499917928@qq.com

*Abstract—With the development of the times and the popularization of scientific and technological products, the Internet has become inseparable from our lives, and search engines have become a daily necessity of people.In view of the growing needs。* **This topic requires the design of a prototype crawler system based on Scrapy framework. The specific requirements and contents are as follows: analyzing the structure and rules of the target website, looking for data items that need to be crawled. Based on Scrappy framework, a crawler prototype program is implemented by customizing crawling rules.Select the appropriate database for data access and analysis.**

*Keywords-Creeper; Scrapy; framework; Python; Cookie*

## I. INTRODUCTION

With the development of the times and the popularization of scientific and technological products, the Internet has become inseparable from our lives, and search engines have become a daily necessity of people. Users can search information by inputting keywords into search engines to find information related to keywords. But with the explosive growth of network information, it becomes more difficult to find the desired information accurately.

In order to meet the growing needs, this paper chooses Scrapy, an open-source crawler framework based on Python, to crawl the "knowledge", and to learn and analyze the principle and running process of the crawler. On this basis, a prototype program of web crawler is implemented, and data crawling and storage are completed.

Firstly, this paper introduces the development process of the crawler, the working principle and classification of the crawler and the grasping strategy, and focuses on the current popular Cookie and its corresponding Session and Robots protocol.

Secondly, the use of Scrapy framework is introduced in detail. Using Scrapy framework to develop crawlers, the process and implementation details of developing crawlers by Scrapy are introduced in detail.

Finally, the crawler is tested and the results of crawling are shown.

## II. WORKING PRINCIPLE

In a word, a crawler is a script or program that can get information and save it. The first step is to send a request to the target web page or website, and then get the response from the server.

Universal crawler is an important part of search engine. Its main function is to collect web pages on the Internet, then save them and process them.

Focus on crawlers, crawlers for specific needs. When it crawls a web page, it filters the first content and grabs the web page information related to the requirements.
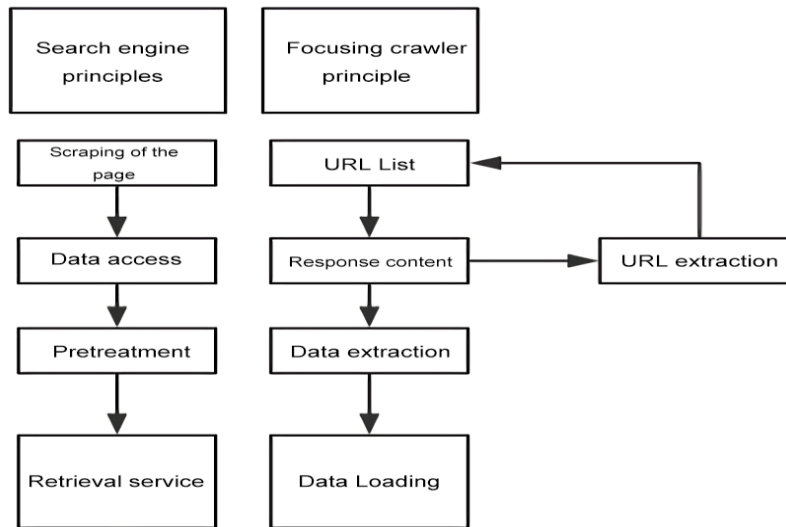
Figure 1.   The Difference between Universal and Focused Reptiles

The crawler workflow is similar to the principle of ordinary users accessing web pages. When a user opens a web page, the browser will send a request to the server visiting the site, and the server will respond to the request and return it to the browser Response. The browser will parse the Response to display the web page[9]. The general crawler framework is shown in Fig. 2below.
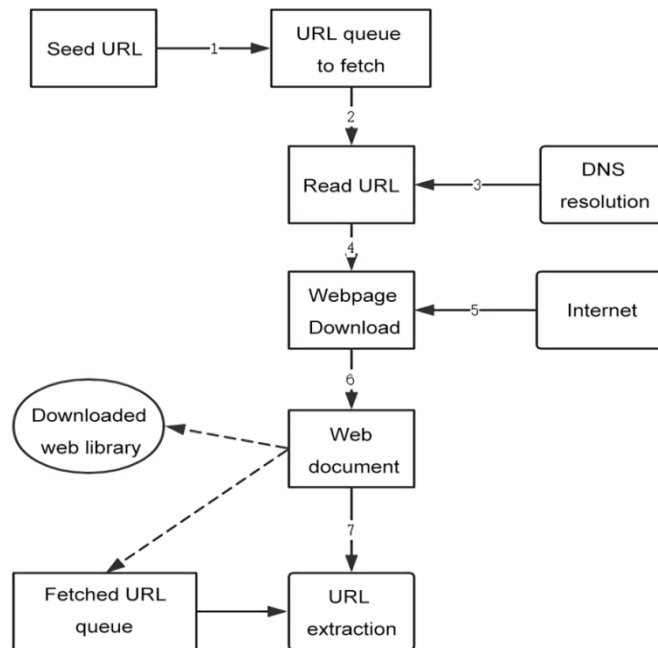


Figure 2.   Universal crawler framework process

First of all, select some sites in the Internet, and take it as a starting point.Put these starting points into the queue to be grabbed, perform the queue out operation, and read out the queue elements.Resolve the URL of the target site through DNS.DNS will convert the domain name to the corresponding IP. The

Downloader Downloads the target page through the server.The URL of the download page will be extracted.Reduplicate the crawl URL queue.The URL of the crawl URL queue continues to loop until the waiting URL queue is empty.

Figure 3.　Network crawler flow chart

## III.　DETAILS OF CRAWLER IMPLEMENT

Data items are obtained by debugging web pages that know the user interface.The　fieldin Python can accept almost any data type.

Follower_info_parse has two functions: first, it can initiate requests for user information through the attention list; second, it has the function of turning over pages. By parsing the response, it can obtain all users of the current target attention list and obtain detailed information of users. There is also the function of page

scheduling to get the list of concerns on the next page. Further requests are then retrieved recursively for circular crawling. Followee_info_parse, which can request user's detailed information through fan list, also has the function of turning pages. Its implementation logic is exactly the same as follower_info_parse, except that the object of the request is different. One is to request detailed information from the person concerned, the other is to request detailed information from the person concerned with the current user.

```
class UserItem(Item):
    name = Field()
    avatar_url = Field()
    headline = Field()
    description = Field()
    url = Field()
    url_token = Field()
    cover_url = Field()
    answer_count = Field()
    articles_count = Field()
    commercial_question_count = Field()
    favorited_count = Field()
    follower_count = Field()
    following_columns_count = Field()
    following_count = Field()
    pins_count = Field()
    question_count = Field()
    thanked_count = Field()
    vote_from_count = Field()
    following_question_count = Field()
    following_topic_count = Field()
    hosted_live_count = Field()
```

Figure 4.   Field definitions in item.py file

```
def followee_info_parse(self, response):
    result = json.loads(response.text)
    if 'data' in result.keys():
        for user in result.get('data'):
            yield                    Request(self.user_info_url.format(user=user.get('url_token'),
include=self.followee_query),
                        callback=self.user_info_parse)
    if 'paging' in result.keys() and result.get('paging').get('is_end') == False:
        next_url = result.get('paging').get('next')
        yield Request(next_url, callback=self.followee_info_parse)
```

Figure 5.   Followee_info_parse method

Spiders.py is the core of web crawling module and an important part of the whole project. It defines the core business logic. Followee_info_parse, which can request user's detailed information through fan list, also has the function of turning pages. Its implementation logic is exactly the same as follower_info_parse, except that the object of the request is different. One is to request detailed information from the person concerned, the other is to request detailed information from the person concerned with the current user.

```
class ZhihuSpider(Spider):
    name = 'zhihu'
    allowed_domains = ['www.zhihu.com']
    start_urls = ['http://www.zhihu.com/']
    start_user = 'india'  |
    user_info_url = 'https://www.zhihu.com/api/v4/members/{user}?include={include}'

    user_query = :...
    follower_url = :...
    follower_query = :...
    followee_url = :...
    followee_query = :...
```

Figure 6.   ZhihuSpider

IV.   RUNNING STATUS AND TESTING

After testing, the crawler capture data of a single host can reach 400,000 users per day. The crawling speed can be artificially controlled by setting it in the code.

```
'description': '深自缄默，如云漂泊。<br/><br/>公号：梁悦<br/>微博：梁悦同学。<br/>约稿转载请私信。<br/><br/>个人简介这
种东西<br/>并不是我讲了你就能搞明白的',
'educations': [],
'employments': [{'company': {'avatar_url': 'https://pic4.zhimg.com/e82bab09c_is.jpg',
                             'id': '',
                             'name': '微信: (约稿/读者)',
                             'type': 'topic',
                             'url': ''},
                 'job': {'avatar_url': 'https://pic4.zhimg.com/e82bab09c_is.jpg',
                         'id': '',
                         'name': 'wx60105991',
                         'type': 'topic',
                         'url': ''}},
                {'company': {'avatar_url': 'https://pic4.zhimg.com/e82bab09c_is.jpg',
                             'id': '20506777',
                             'name': '深圳图书馆',
                             'type': 'topic',
                             'url': 'https://www.zhihu.com/topics/20506777'},
                 'job': {'avatar_url': 'https://pic2.zhimg.com/ec1267951_is.jpg',
                         'id': '19709116',
                         'name': '图书管理员',
                         'type': 'topic',
                         'url': 'https://www.zhihu.com/topics/19709116'}}],
'favorite_count': 3,
'favorited_count': 332298,
'follower_count': 135383,
'following_columns_count': 28,
'following_count': 953,
'following_favlists_count': 16,
'following_question_count': 917,
'following_topic_count': 101,
'gender': 1,
'headline': '公众号：梁悦/微博：梁悦同学',
'hosted_live_count': 0,
'id': '6949ebdcf63dbbcc4b2d6b198927b489',
'locations': [{'avatar_url': 'https://pic2.zhimg.com/v2-d1e9b2b9f276a1e6f21d1179aa356baa_is.jpg',
               'id': '19560551',
               'name': '深圳市',
               'type': 'topic',
               'url': 'https://www.zhihu.com/topics/19560551'}],
'marked_answers_count': 0,
'mutual_followees_count': 0,
'name': '梁悦',
```

Figure 7.   Screenshots of crawling 1

Figure 8.   Screenshots of crawling 2

After the program runs, it gets a database named "zhihu" and stores all the information in the user table.



Figure 9.   Database screenshots

## V.  CONCLUSION

When the crawler development is completed, it should be tested. Testing is a very important step. First of all, we need to know the performance of the crawler through testing, check whether the crawler has problems, and whether it can crawl the required data. Secondly, we should explore the anti-crawler strategy of the target website and improve the crawler. Finally, check the data that has been crawled to see if it achieves the expected goal of the project.The crawler system can also be extended, there are many technologies not added to it, and then added to it is the requirements of the enterprise level. In the process of writing this system, I consulted a lot of information about Scrapy. Scrapy framework is a new thing for me. New APIs and libraries. Fortunately, I have done some crawler projects before, which is not particularly difficult for me.

## ACKNOWLEDGMENT

## REFERENCE

[1]  Jing Wang,YuchunGuo. Scrapy-Based Crawling and User-Behavior Characteristics Analysison Taobao[P]. Cyber-Enabled Distributed Computing and Knowledge Discovery (CyberC), 2012 International Conference on,2012.

[2]  James W. McGuffee. Non-profit geographically constrained locator[J]. ACM SIGCASComputers and Society,2015,45(2).

[3]  Yuhao Fan. Design and Implementation of Distributed Crawler System Based on Scrapy[J].IOP Conference Series: Earth and Environmental Science,2018,108(4).

[4]  Shen Jie, Li Yifan. Application of Web crawler system in cloud media [J]. China Cable Television, 2018 (05): 595-597.

[5]  Zhang Jin. Research on Web crawler technology based on Hadoop platform [D].Nanjing University of Posts and Telecommunications, 2017.

[6]  Zhao Fen, Lei Zhenzhen, Yang Xiaoyun, Su Pengju and Wang Shunye.Based on Baidu Tieba College Students'Network Public Opinion Analysis [J]. Computer Knowledge and Technology, 2018, 14 (28): 227-229.

[7]  Ding Zhongxiang, Yang Yanhong, DuYanming. Design and implementation of video information crawling based on Scrappy framework [J]. Journal of Beijing Printing Institute, 2018, 26 (09): 92-97.

[8]  Xie Zhu. Emotional Tendency Analysis for Chinese Short Texts [D].Hunan University, 2018.

[9]  Wei Chengcheng. Data Information Crawler Technology Based on Python [J]. Electronic World, 2018 (11): 208-209.

[10] Ye Xiqiezhong. Research and Implementation of Tibetan Text Automatic Classification Based on Web [D]. Qinghai University for Nationalities, 2014.

[11] ZhongJiajun. Research on Copyright Infringement Recognition of News Aggregation Platform [D].Lanzhou University, 2018.