

How to collect data

To get enough book data, I write a crawler by python to collect those data.

<https://www.booktxt.net/xiaoshuodaquan/> In this website, there're around 3k books. By the benefits of multi-threading, the crawler can retrieve hundreds of books in just a few minutes.

Here's the main part of my crawler

```
1  def main():
2      start_time = time.time()
3      books = getBooks()
4      i = 0
5      total_size = len(books)
6      while i < 100: # 只爬前100本书，总共三千本书，实在是太多了，没必要
7          if(threading.activeCount() < 50):
8              thd = threading.Thread(target=getBook,args=(books[i],i))
9              i += 1
10             thd.start()
11             time.sleep(1)
12         for thd in threading.enumerate():
13             if(thd.getName() != threading.currentThread().getName()):
14                 thd.join()
15         stop_time = time.time()
16         diff = stop_time - start_time
17         print(diff,"seconds")
18         print(diff/60,"minutes")
19         print(diff/60/60,"hours")
20
21
22 if __name__ == '__main__':
23     main()
```

How to store the content of books

We didn't store the content of books in database, instead, we store the content in some files, and store the uri location in database using a field called `location`.

All book contents will be stored under a directory named `book`. For specific book, there'e filed called `chapter_no` in our book table, which represents the total chapter number of a book. Books will be

stored in a directory with name of its `location` . And each chapter will be stored in a txt file with chapter number as file name.

Also, every book will have some introduction message stored in a file named `info.txt` and a cover image stored in a file named `cover.jpg` within the same directory.

For instance, suppose there a book with following data

```
1 | {  
2 |     "location": "qwerty",  
3 |     "chapter_no": 10  
4 | }
```

In this example, the first chapter will be stored in `book/qwerty/1.txt` . The second chapter will be stored in `book/qwerty/2.txt` . Book introduction will be stored in `book/qwerty/info.txt` . Cover image will be stored in `book/qwerty/cover.jpg` .

Where's our database

We use a server to deploy our database, which is at www.irran.top. We can login to this database through command like this

```
1 | mysql -h www.irran.top -u yggdrasil -p
```

How to modify book content stored in our server

Client application will issue right formatted post request. I design a few servlet written in JAVA to process those request and modify data accordingly.

Here's one of my servlets

@Override

public void doPost(HttpServletRequest req, HttpServletResponse resp) throws ServletException

try {

String location = req.getParameter("location");

String content = req.getParameter("content");

String path = String.format("%s%s/info.txt", getServletContext().getRealPath(

FileWriter fw = new FileWriter(path);

fw.write(content);

fw.flush();

fw.close();

}catch (Exception e) {

PrintWriter out = resp.getWriter();

out.println(e.getMessage());

}

}