

**Московский авиационный институт
(национальный исследовательский университет)**

**Факультет информационных технологий и прикладной
математики**

Кафедра вычислительной математики и программирования

Лабораторная работа №2 по курсу «Информационный поиск»

Студент: М. М. Сисенов
Преподаватель: А. А. Кухтичев
Группа: М8О-410Б
Дата:
Оценка:
Подпись:

Москва, 2025

Лабораторная работа №2 «Поисковый робот»

Необходимо написать парсер на любом языке программирования:

- Написать поисковый робот — компоненты обкачки документов, используя любой язык программирования.
- Единственным аргументом поисковому роботу подаётся файл конфигурации (формата YAML или JSON), в котором содержатся параметры работы программы.
- База данных должна быть запущена в docker-контейнере, можно использовать docker-compose.
- В качестве хранилища результатов использовать MongoDB или PostgreSQL.
- Робот должен применять нормализацию URL-адресов.
- Робот должен сохранять в БД сырой HTML документа.
- Необходимо сохранять метаинформацию о каждом документе (дата скачивания, источник URL и т.д.).
- При остановке работы робот должен сохранять контрольную точку так, чтобы при повторном запуске он мог продолжить работу с того документа, с которого он остановился.
- Периодически он должен уметь переобкачивать документы, которые уже есть в базе, но только в том случае, если они изменились.

Задание

Требуется реализовать веб-краулер (поисковый робот), который автоматически собирает документы с психологических порталов **b17.ru** и **psychologies.ru**, выбранных в лабораторной работе №1. Робот должен сохранять полный HTML-контент страниц в базу данных MongoDB вместе с метаинформацией, обеспечивать возможность остановки/возобновления работы и отслеживать изменения документов.

Выбор языка программирования и технологий

Для реализации поискового робота был выбран язык программирования **Python**

Стек технологий:

- Язык: Python 3.12.3
- База данных: MongoDB 8.x
- HTTP-клиент: `requests`
- HTML-парсинг: BeautifulSoup4 (библиотека bs4)
- Работа с БД: `pymongo`
- Конфигурация: YAML (PyYAML)

Схема работы

1. **Инициализация:** Чтение конфигурации из YAML-файла, подключение к MongoDB и создание уникальных индексов для URL.
2. **Управление состоянием:** Использование JSON-файла для хранения номеров страниц, на которых остановился сборщик ссылок (Harvester).
3. **Этап 1: Сбор ссылок (Harvesting):**
 - Обход страниц-списков в разделах статей;
 - Извлечение и нормализация URL;
 - Проверка необходимости обкачки (отсутствие в базе или истечение заданного интервала времени);
 - Добавление уникальных ссылок в коллекцию-очередь задач MongoDB.
4. **Этап 2: Обкачка контента (Crawling):**
 - Извлечение задач из очереди;
 - Загрузка HTML-контента с соблюдением задержек;
 - Извлечение очищенного текста и заголовка.
5. **Отслеживание изменений:** Сравнение нового очищенного текста с версией из базы данных. Если текст не изменился, обновляется только метка времени обкачки.
6. **Сохранение:** Запись документов с полями URL, источник, заголовок, «сырой» HTML, очищенный текст и время обкачки в формате Unix timestamp.

Реализация

1. Конфигурационный файл (config.yaml)

```
1 db:
2     host: "localhost"
3     port: 27017
4     name: "search_engine_db"
5     collection_articles: "articles_v5"
6     collection_queue: "task_queue"
7
8 logic:
9     delay_min: 1.0
10    delay_max: 3.0
11    recrawl_interval: 86400
12    state_file: "crawler_state.json"
```

2. Структура документа в MongoDB

```
1 doc = {
2     'url': url,
3     'source': source,
4     'title': title,
5     'raw_html': html,
6     'clean_text': text,
7     'timestamp': time.time()
8 }
```

Пояснение полей:

- `url` — нормализованный URL документа
- `source` — название ресурса
- `title` — Название статьи
- `raw_html` — полный HTML-контент страницы
- `clean_text` — очищенный полезный текст
- `timestamp` — timestamp скачивания (Unix-время)

3. Нормализация URL

```
1 def normalize_url(url):
2     parsed = urlparse(url)
3     clean_url = urlunparse((parsed.scheme.lower(), parsed.netloc.lower(),
4                               (), parsed.path, '', ''))
5     if clean_url.endswith('/'):
6         clean_url = clean_url[:-1]
7     return clean_url
```

При запуске краулер проверяет наличие `checkpoint.json` и продолжает работу с последнего обработанного URL.

Исходный код

```
1 import requests
2 from bs4 import BeautifulSoup
3 from pymongo import MongoClient
4 import time
5 import random
6 import json
7 import os
8 import sys
9 import yaml
10 import argparse
11 from urllib.parse import urlparse, urlunparse
12
13
14 def load_config(path):
15     with open(path, 'r') as f:
16         return yaml.safe_load(f)
17
18
19 parser = argparse.ArgumentParser(description='SearchRobot')
20 parser.add_argument('config', type=str, help='Path to config.yaml')
21 args = parser.parse_args()
22
23 cfg = load_config(args.config)
24
25 client = MongoClient(cfg['db']['host'], cfg['db']['port'])
26 db = client[cfg['db']['name']]
27 a_coll = db[cfg['db']['collection_articles']]
28 q_coll = db[cfg['db']['collection_queue']]
29
30 a_coll.create_index("url", unique=True)
31 q_coll.create_index("url", unique=True)
32
33 HEADERS = {
34     'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/124.0.0.0 Safari/537.36',
35     'Accept-Language': 'ru-RU,ru;q=0.9,en-US;q=0.8,en;q=0.7'
36 }
37
38
39 def normalize_url(url):
40     parsed = urlparse(url)
41     clean_url = urlunparse((parsed.scheme.lower(), parsed.netloc.lower(),
42                             parsed.path, '', '', ''))
43     if clean_url.endswith('/'):
44         clean_url = clean_url[:-1]
45     return clean_url
46
47 def load_state():
48     if os.path.exists(cfg['logic']['state_file']):
```

```

49         with open(cfg['logic']['state_file'], 'r') as f:
50             return json.load(f)
51     return {"b17_page": 1, "psych_page": 1}
52
53
54 def save_state(b17_p, psych_p):
55     with open(cfg['logic']['state_file'], 'w') as f:
56         json.dump({"b17_page": b17_p, "psych_page": psych_p}, f)
57
58
59 def add_to_queue(url, source):
60     url = normalize_url(url)
61     existing_doc = a_coll.find_one({'url': url})
62
63     if existing_doc:
64         last_scan = existing_doc.get('timestamp', 0)
65         if time.time() - last_scan < cfg['logic']['recrawl_interval']:
66             return False
67         else:
68             print(f"Re-crawl:{url}")
69
70     try:
71         q_coll.update_one(
72             {'url': url},
73             {'$set': {'url': url, 'source': source, 'added_at': time.
74                         time()}},
75             upsert=True
76         )
77         return True
78     except:
79         return False
80
81
82 def harvest_lists():
83     state = load_state()
84     b17_page = state['b17_page']
85     psych_page = state['psych_page']
86
87     print("Starting harvester...")
88
89     try:
90         while True:
91             # B17
92             try:
93                 r = requests.get(f"https://www.b17.ru/article/?page={
94                     b17_page}", headers=HEADERS)
95                 if r.status_code == 200:
96                     soup = BeautifulSoup(r.text, 'lxml')
97                     links = soup.find_all('a', href=True)
98                     if not links:
99                         break
99
99                     count = 0

```

```

100         for a in links:
101             href = a['href'].split('#')[0]
102             if href.startswith('/article/') and href.count(
103                 '/') == 3:
104                 if add_to_queue("https://www.b17.ru" + href,
105                     'b17'):
106                     count += 1
107             print(f"B17_page:{b17_page}:{count} added")
108             b17_page += 1
109     except Exception as e:
110         print(f"Error B17: {e}")
111
112     time.sleep(1)
113
114     # Psychologies
115     try:
116         r = requests.get(f"https://www.psychologies.ru/articles
117                         /{psych_page}/", headers=HEADERS)
118         if r.status_code == 200:
119             soup = BeautifulSoup(r.text, 'lxml')
120             count = 0
121             for a in soup.select('a.rubric-anons_title, .rubric
122             -anons_list a.link'):
123                 href = a.get('href')
124                 if href:
125                     full = href if href.startswith('http') else
126                         "https://www.psychologies.ru" + href
127                     if add_to_queue(full, 'psychologies'):
128                         count += 1
129                     print(f"Psych_page:{psych_page}:{count} added")
130                     psych_page += 1
131     except Exception as e:
132         print(f"Error Psych: {e}")
133
134     save_state(b17_page, psych_page)
135     time.sleep(1)
136
137     except KeyboardInterrupt:
138         print("\nStopped.")
139
140
141 def process_queue():
142     print(f"Crawler started (Delay: {cfg['logic']['delay_min']}-{cfg['
143         logic']['delay_max']}s)")
144
145     while True:
146         task = q_coll.find_one_and_delete({})
147         if not task:
148             print("Queue empty.")
149             break
150
151         url = task['url']
152         source = task['source']

```

```

147
148     time.sleep(random.uniform(cfg['logic']['delay_min'], cfg['logic']
149                   ]['delay_max']))
150
151     try:
152         r = requests.get(url, headers=HEADERS, timeout=10)
153         if r.status_code == 200:
154             html = r.text
155             soup = BeautifulSoup(html, 'lxml')
156
157             if source == 'b17':
158                 title, text = extract_b17_data(soup)
159             else:
160                 title, text = extract_psych_data(soup)
161
162             if len(text) > 200:
163                 old_doc = a_coll.find_one({'url': url})
164                 if old_doc and old_doc.get('clean_text') == text:
165                     a_coll.update_one({'url': url}, {'$set': {'
166                         timestamp': time.time()}})
167                     print(f"Unchanged:{url}")
168                 else:
169                     doc = {
170                         'url': url,
171                         'source': source,
172                         'title': title,
173                         'raw_html': html,
174                         'clean_text': text,
175                         'timestamp': time.time()
176                     }
177                     a_coll.replace_one({'url': url}, doc, upsert=
178                         True)
179                     action = "Updated" if old_doc else "Saved"
180                     print(f"{action}:{url}")
181                 else:
182                     print(f"Short content skip:{url}")
183             else:
184                 print(f"Status:{r.status_code}:{url}")
185
186     except Exception as e:
187         print(f"Network error:{url}{e}")
188
189
190     def extract_b17_data(soup):
191         title_tag = soup.find('h1', class_='from_bb_h1') or soup.find('h1')
192         title = title_tag.get_text(strip=True) if title_tag else "No Title"
193         content_div = soup.find('div', attrs={'itemprop': 'articleBody'}) or \
194             \
195             soup.find('div', attrs={'itemprop': 'articleBody'}) \
196             or \
197             soup.find('div', id='article_body')
198
199         if content_div:
200             for j in content_div.find_all('div', class_='art_start'): j.
201                 decompose()

```

```

194         for t in content_div(['script', 'style']): t.decompose()
195     return title, content_div.get_text(separator=' ', strip=True)
196     return title, ""
197
198
199 def extract_psych_data(soup):
200     title_tag = soup.find('h1', class_='article__title')
201     title = title_tag.get_text(strip=True) if title_tag else "No Title"
202     parts = []
203     lead = soup.find('p', class_='article__lead-paragraph')
204     if lead: parts.append(lead.get_text(strip=True))
205     body = soup.find('section', attrs={'itemprop': 'articleBody'})
206     if body:
207         for bl in body.find_all('div', class_='article__block-type-text'):
208             parts.append(bl.get_text(separator=' ', strip=True))
209     return title, ''.join(parts)
210
211
212 if __name__ == "__main__":
213     print("1. Harvest")
214     print("2. Crawl")
215     choice = input("Choice: ")
216     if choice == '1':
217         harvest_lists()
218     elif choice == '2':
219         process_queue()

```

Вывод

В ходе лабораторной работы был разработан полнофункциональный поисковый робот на языке Python, интегрированный с СУБД MongoDB. Реализованная двухэтапная архитектура (разделение на сбор ссылок и загрузку контента) позволила эффективно управлять очередью задач и гибко настраивать интенсивность запросов к каждому источнику.

Использование контрольных точек (checkpoint) обеспечило отказоустойчивость: робот успешно возобновляет работу после прерывания. Была внедрена система отслеживания изменений, исключающая дублирование данных и позволяющая актуализировать корпус документов без лишних затрат ресурсов. Итоговый робот полностью соответствует требованиям ТЗ и сформировал надежную базу данных для последующих этапов построения поискового движка.

Литература

- [1] Маннинг, Рагхаван, Шютце *Введение в информационный поиск* — Издательский дом «Вильямс», 2011. — 528 с.