

Simple Search Engine

简介

本项目实现了一个简单的搜索引擎，下面简单介绍要点，详细技术要点参考 [技术文档](#)

- 功能要点： 查询网页、异步爬虫、添加文档、缓存快速删除文档、搜索文档、其他数据维护功能
- 技术要点：
 1. 项目包含 web、crawler、engine 三大模块，全部由 scala 编写，使用akka提供Actor并发模型支持
 2. web模块：由 Spring Boot 后端和静态 HTML 前端（原生ES6）组成
 3. crawler模块：基于Actor模型的异步爬取和异步添加/更新文档、定期搜集、种子搜索+维护url列表策略
 4. engine模块：Actor并发与Future异步方式、缓存快速删除文档设计、缓存索引表和url库、文件读写锁
 5. 其他模块：定时爬取任务，定时缓存持久化，分词功能
- 搜索引擎过程要点：
 1. web搜集：异步爬虫、定期搜集、种子搜索+维护url列表策略
 2. 预处理：分词、正则提取url、正则提取正文
 3. 维护文件：网页库（文档信息、文档正文）、索引表（文档id、内容hash、网页库文件偏移量）、倒排表（词、文档id、文档位置列表）、文档计数表（总文档数、各文档词数）、url库（url、正文hash）
 4. 查询服务：包含信息列表三要素（标题、URL、摘要）、字符串切份生成摘要
 5. 查询排序：BM25排序算法

用户手册

演示

本项目提供简单的DEMO，以便跳过部署和运行从而直接使用：

<http://39.106.185.26:8080/>

配置

修改 `Config.scala` 相关项目，如：文件存储位置、超时时间、hash范围等

编译

使用Maven管理，springboot插件编译打jar包：`mvn clean package`

部署

建议使用 `systemd`：

1. 编写 `sse.service`，这里提供简单例子：

```
[Unit]
Description=Simple search engine

[Service]
Type=simple
User=root
ExecStart=/usr/bin/java -jar /opt/sse/simple-search-engine.jar --
server.port=8080

[Install]
WantedBy=multi-user.target
```

2. 复制jar包到对应位置
3. `systemctl daemon-reload` 重新加载

运行

```
systemctl start sse
```

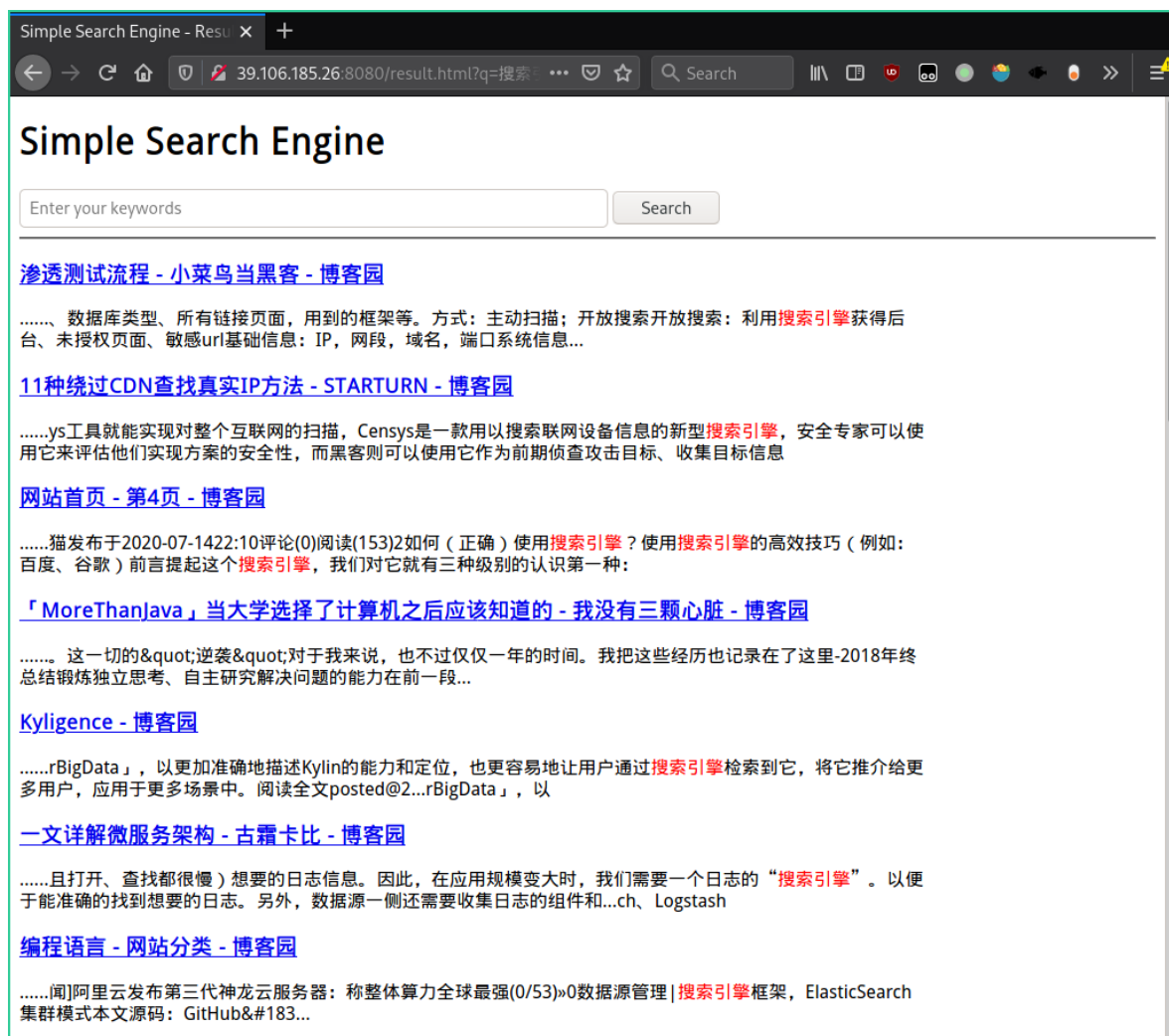
使用

访问8080端口即可使用

使用截图

DEMO配置中设置了只爬取博客园的策略，所以搜索结果都是博客园

PC端



移动端



Simple Search Engine

[待解决高分问题_博问_博客园](#)

.....源于数据库查询的数据集合。suqq小白回答(0)浏览(38)4天前1回答数100vue移动端项目中pdfjs-dist插件报错报错信息有下面这两条Uncaught...<'pd

[博问_程序员问答社区，解决您的IT难题_博客园](#)

.....据，导致爬虫爬不到，渲染的页面就不对了，我看了下network，也没看到发出过请vue axios遗失的美好,,,回答(0)浏览(24)1周前3回答数5webuplo...cd~user1进入个

[所有评论 - 博客园](#)

.....经很开心了(●~▽~●)倾叶子伶评论于2020-07-1522:41Re:记录下vue中引用echarts出现"TypeError:Cannotrea...步加载完可以不h祝h评

[精华区 - 第4页 - 博客园](#)

.....时下热门的任何一种技术选型都有成熟的

方案，比如：vue+vue-i18nangular+angular-tran...

[伍华聪 - 博客园](#)

.....见的图标，不过数量不是很多，应该是30

