# 爬虫简介

# 大纲

- 爬虫是什么
- 基础知识
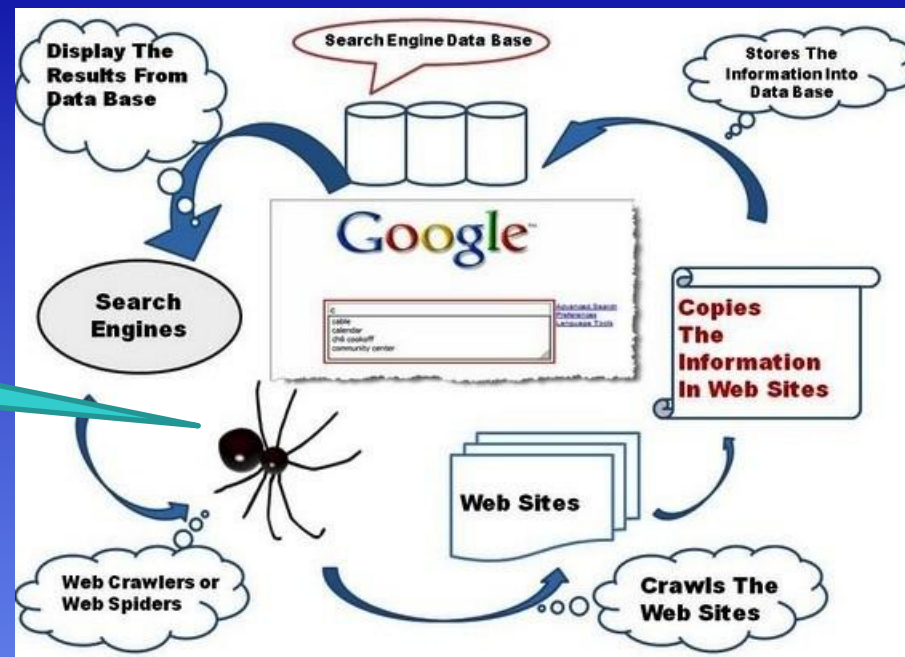- 爬虫分类
- 爬虫架构
- 反抓取
- 网页解析
- 工具介绍

# 爬虫是什么

- 名字
  - Crawler、Spider、Bot
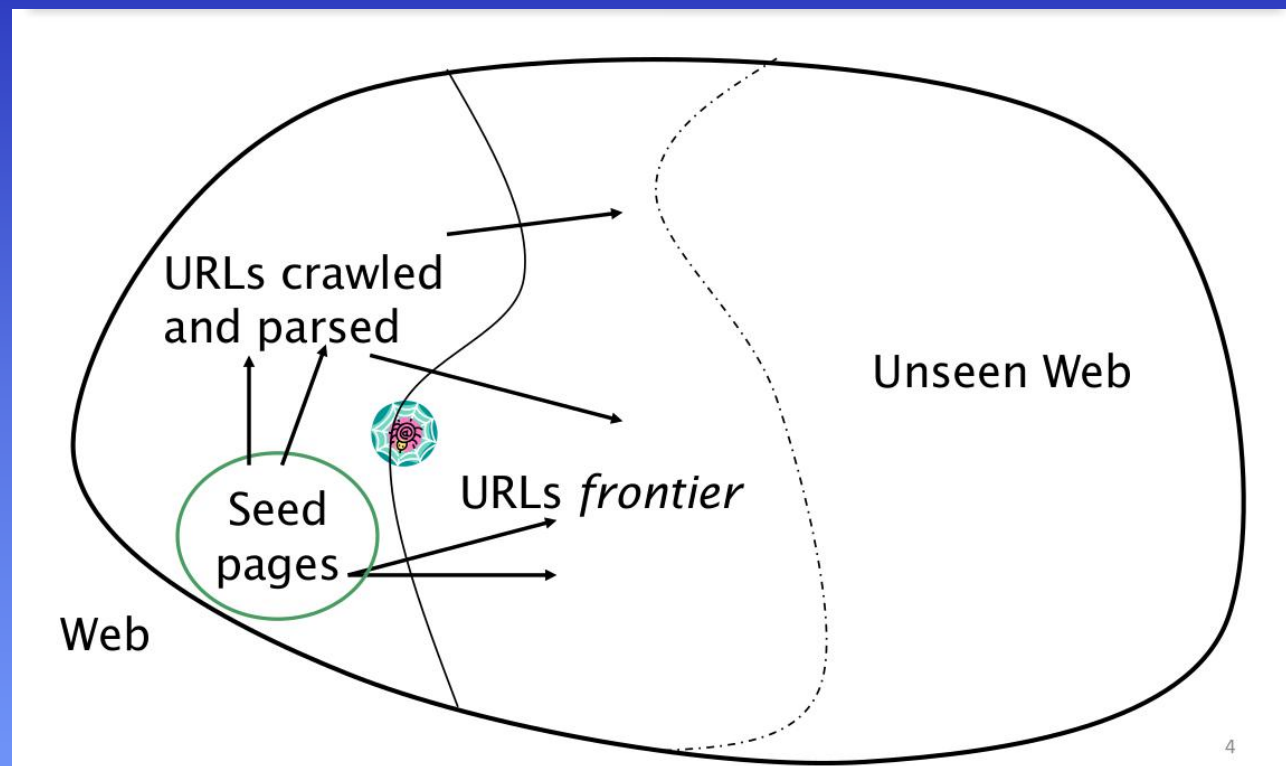  - Web Scraping(verb)
- 定义(wiki)
  - A Web crawler, sometimes called a spider, is an Internet bot that systematically browses the World Wide Web, typically for the purpose of Web indexing

Spider

# 爬虫的抓取过程

- 从"种子"链接开始
  - 把它们加到待抓取"队列"

- While True:
  - 从队列中取任务
  - 抓取内容并且抽取链接
  - 把链接去重后加入队列

URLs crawled and parsed
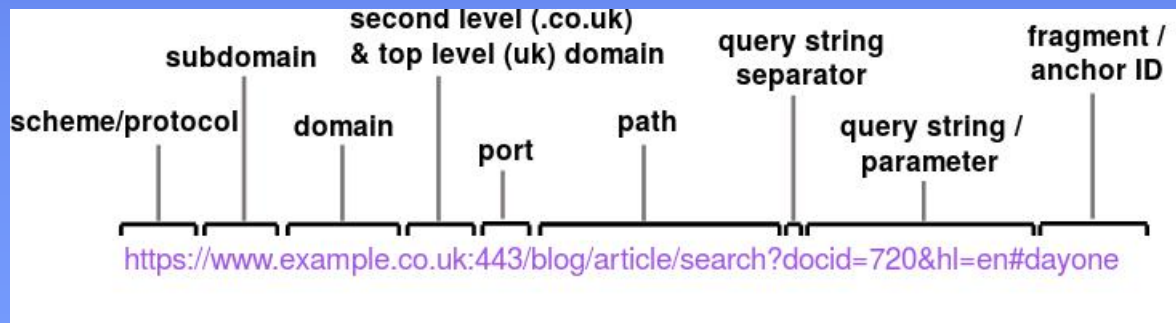
Unseen Web

Seed pages

URLs *frontier*

Web

# 基础知识

- URL
- HTML/CSS/JS
- HTTP
- DNS

# URL

- Uniform Resource Locator
  - colloquially termed a web address, is a reference to a web resource that specifies its location on a computer network and a mechanism for retrieving it
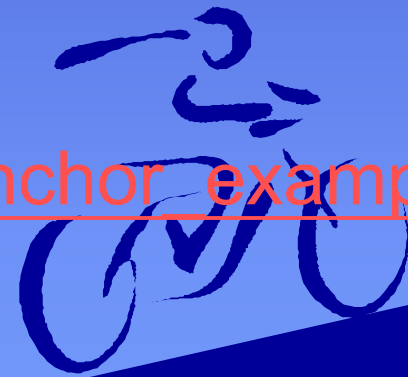- scheme:[//[user:password@]host[:port]][/]path[?query][#fragment]

# URL

- Encoding/Decoding
  - https://www.baidu.com/s?ie=utf-8&f=8&rsv_bp=0&rsv_idx=1&tn=baidu&wd=%E5%A4%A9%E6%B0%94

- Fragment
  - http://www.tagindex.net/html/link/anchor_example2.html#a003
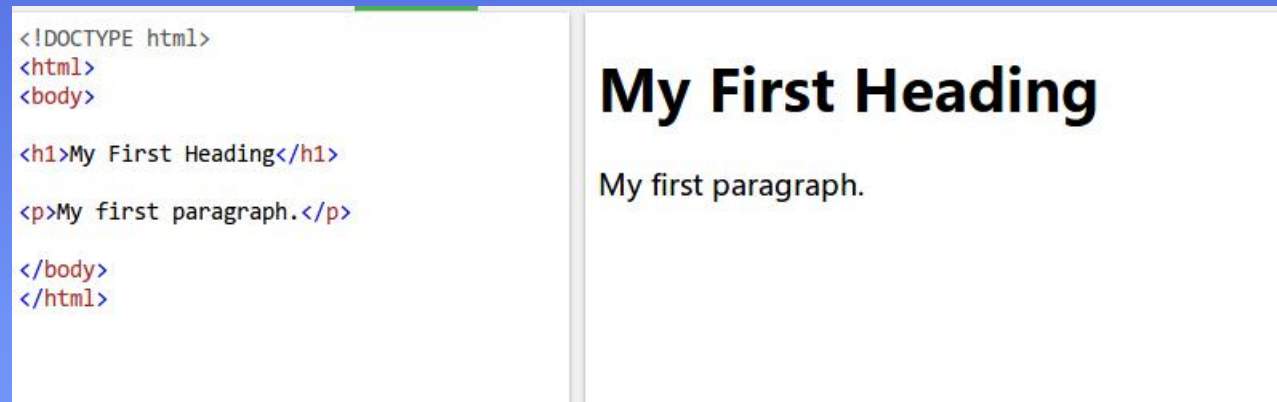
# URL Normalization

- HTTP://www.Example.com/ → http://www.example.com/
- http://www.example.com/a%c2%b1b → http://www.example.com/a%C2%B1b
- http://www.example.com:80/bar.html → http://www.example.com/bar.html

- Google Bot remove fragment
- https://en.wikipedia.org/wiki/URL_normalization

# HTML

- Hypertext Markup Language (HTML) is the standard markup language for creating web pages and web applications. With Cascading Style Sheets (CSS) and JavaScript, it forms a triad of cornerstone technologies for the World Wide Web

```
<!DOCTYPE html>
<html>
<body>

<h1>My First Heading</h1>

<p>My first paragraph.</p>

</body>
</html>
```

**My First Heading**

My first paragraph.

- https://www.w3schools.com/html/tryit.asp?filename=tryhtml_basic_document

# HTML

- [http://news.sina.com.cn/c/2018-07-03/doc-ihevauxi3422393.shtml](http://news.sina.com.cn/c/2018-07-03/doc-ihevauxi3422393.shtml)

# 不只是一个HTML！

# 浏览器渲染过程

# 浏览器渲染过程

- https://www.html5rocks.com/en/tutorials/internals/howbrowserswork/

# Ajax

- Ajax ( "Asynchronous JavaScript And XML") is a set of Web development techniques using many Web technologies on the client side to create asynchronous Web applications

# Ajax

- https://music.163.com/#/song?id=168091

# HTTP

- ## Hypertext Transfer Protocol

  – The Hypertext Transfer Protocol (HTTP) is an application protocol for distributed, collaborative, and hypermedia information systems.[1] HTTP is the foundation of data communication for the World Wide Web.

```
GET /index.html HTTP/1.1
Host: www.example.com
```

```
HTTP/1.1 200 OK
Date: Mon, 23 May 2005 22:38:34 GMT
Content-Type: text/html; charset=UTF-8
Content-Length: 138
Last-Modified: Wed, 08 Jan 2003 23:11:55 GMT
Server: Apache/1.3.3.7 (Unix) (Red-Hat/Linux)
ETag: "3f80f-1b6-3e1cb03b"
Accept-Ranges: bytes
Connection: close

<html>
<head>
  <title>An Example Page</title>
</head>
<body>
  Hello World, this is a very simple HTML document.
</body>
</html>
```

# HTTP Headers

- 请求
  - Referer
  - Cookie
  - User-Agent
  - Accept-Encoding
- 响应
  - Cache-Control/Expires
  - Last-Modified/If-Modified-Since
  - ETag/If-None-Match

# HTTP Header Example

```
(py3-env) lili@lili-Precision-7720:~$ curl -I https://www.baidu.com/img/bd_logo1.png
HTTP/1.1 200 OK
Accept-Ranges: bytes
Cache-Control: max-age=315360000
Connection: Keep-Alive
Content-Length: 7877
Content-Type: image/png
Date: Tue, 03 Jul 2018 09:54:44 GMT
Etag: "1ec5-502264e2ae4c0"
Expires: Fri, 30 Jun 2028 09:54:44 GMT
Last-Modified: Wed, 03 Sep 2014 10:00:27 GMT
P3p: CP=" OTI DSP COR IVA OUR IND COM "
Server: Apache
Set-Cookie: BAIDUID=84B64DAB4F10D451E01F3BE199B69237:FG=1; expires=Wed, 03-Jul-19 09:54:44 GMT; max-age=31536000;
path=/; domain=.baidu.com; version=1
```

```
(py3-env) lili@lili-Precision-7720:~$ curl -I -H "If-Modified-Since: Wed, 03 Sep 2014 10:00:27 GMT" https://www.ba
idu.com/img/bd_logo1.png
HTTP/1.1 304 Not Modified
Cache-Control: max-age=315360000
Connection: Keep-Alive
Date: Tue, 03 Jul 2018 09:55:08 GMT
Etag: "1ec5-502264e2ae4c0"
Expires: Fri, 30 Jun 2028 09:55:08 GMT
Server: Apache
Set-Cookie: BAIDUID=10CAD5741EC0F3097A5E46690BC19D6D:FG=1; expires=Wed, 03-Jul-19 09:55:08 GMT; max-age=31536000;
path=/; domain=.baidu.com; version=1
```

# DNS

- The Domain Name System (DNS) is a hierarchical decentralized naming system for computers, services, or other resources connected to the Internet or a private network.

```
(py3-env) lili@lili-Precision-7720:~$ nslookup www.baidu.com
Server:         127.0.1.1
Address:        127.0.1.1#53

Non-authoritative answer:
Name:   www.baidu.com
Address: 119.75.216.20
Name:   www.baidu.com
Address: 119.75.213.61
```

```
[easemob@vip3-ali-hangzhou-ai-ecs-sdb-poc-mesos1 1.16.4.ZHONGYUAN.FINAL]$ nslookup www.baidu.com
Server:         10.143.22.116
Address:        10.143.22.116#53

Non-authoritative answer:
www.baidu.com   canonical name = www.a.shifen.com.
Name:   www.a.shifen.com
Address: 220.181.111.188
Name:   www.a.shifen.com
Address: 220.181.112.244
```
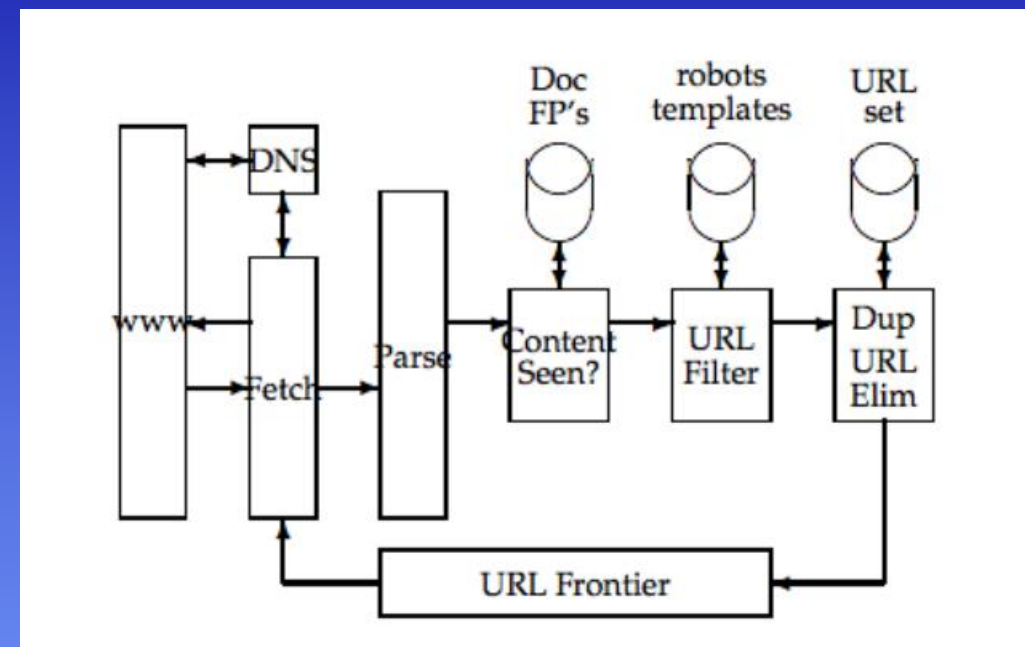
# 爬虫分类

- 通用爬虫
  - 通用搜索引擎
- 垂直(Vertical)爬虫
  - 垂直搜索引擎
- 定向爬虫
  - Web Scraping

# 通用爬虫架构

- URL frontier
- DNS resolution
- Fetching module
- Parsing module
- Duplicate elimination

# 通用爬虫的挑战

- 海量数据
  - Google索引了上千亿网页
- 实时
  - 新的文章能在几分钟搜到
- 深度网络
  - 社交媒体、自媒体
- 反作弊
  - 爬虫陷阱

定向爬虫

# 反爬虫

- robots.txt https://www.baidu.com/robots.txt
- 封ip/验证码
- headers检验
- 登录
- js "加密"

# 网页解析

- 抽取网页的重要数据
  - 标题
  - 正文
  - 发表时间
- 方法
  - 手写模板/规则
  - 通用方法

# XPath

- https://www.w3schools.com/xml/xpath_intro.asp

- https://doc.scrapy.org/en/xpath-tutorial/topics/xpath-tutorial.html

# 工具介绍

- Firebug
- Tcpflow
- Xpath插件

# 练习

- 在终端打印百度首页的DOM树
  - 递归遍历DOM树
- 抓取一个网站
  - 使用介绍的框架
- 根据关键词抓取搜狗微信文章
  - Referer、速度控制
- 抓取https://music.163.com/#/song?id=168091 的所有评论
  - url前端加密

# 练习

- 通用的标题和正文抽取器
  - anchor text ratio
- 网页截图，超过一屏的截屏，某个元素(验证码)的截屏
  - AShot
- 抓取qunar明天北京到上海所有航班的票价
  - 前端视觉混淆