

课程目标



工作日常



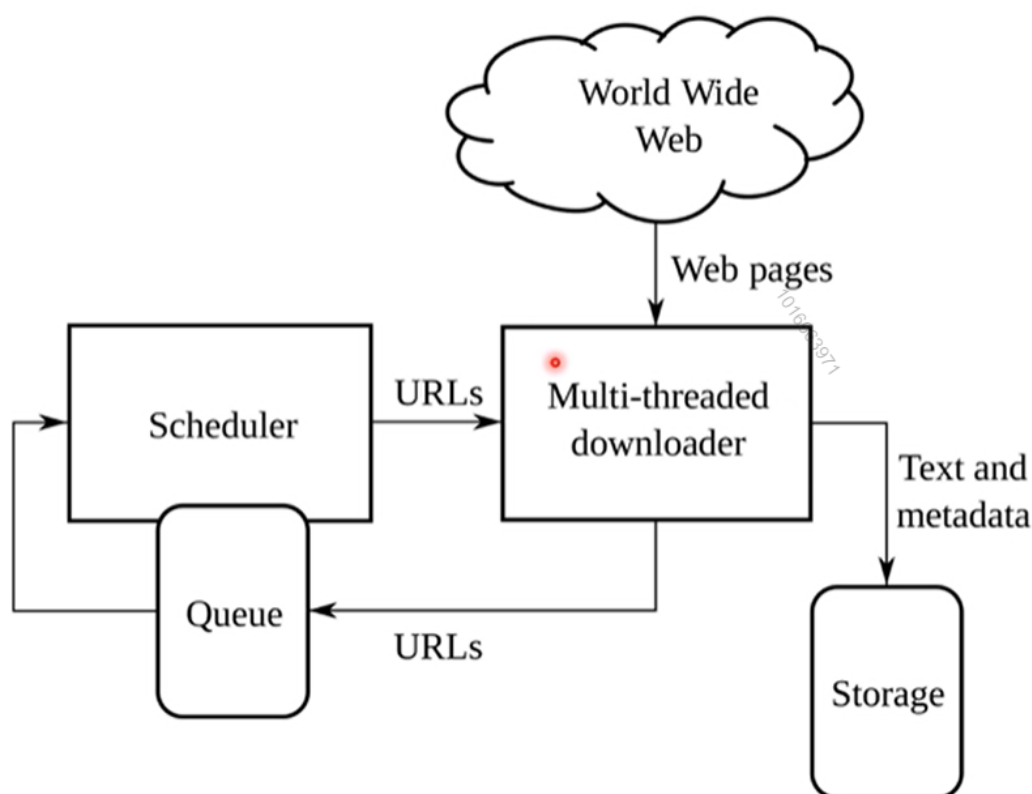
工具链



前端生态

爬取B站数据

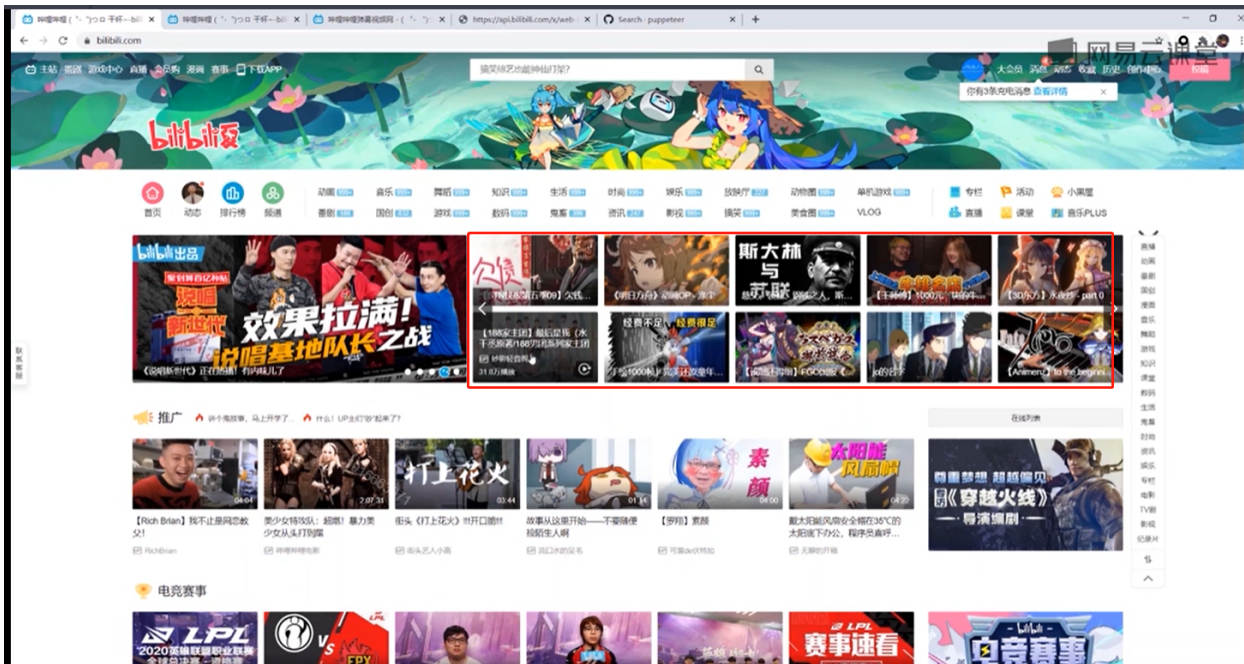
基本概念：



scheduler(一个算法，他的作用是调度)，有的时候我们不只使用一个爬虫，比如你要爬互联网上的数据，互联网上有多少的数据呢？据这种不太专业的数据统计，也有上亿个网站，要把这些网站都爬掉的话都需要调度的，即使你不像百度那样需要爬所有的网站，如果需要爬B站的数据，也是很多的。所有需要一个调度器的，通常爬虫的话，一台机器还不够，可能需要很多机器。爬虫也可以做水军，模拟登录，很多评论都是由爬虫弄进去的，并不是用户这真实的评论。scheduler找到urls，发送给**downloader**下载器，他的主要作用是把

数据下载下来，然后把数据存到Storage里面去。下载器downloader下载后，会进行分析，然后产生新的URLs，比如今天爬B站，我们从B站的首页拿到很多卡片，卡片点进去的需要一个新的连接，跳转。

实战(B站有价值就是相关卡片数据，然后相关详情里面的相关推荐)



node服务端爬部分公开api

```
1 const fetch = require('node-fetch')
2 const app = require('express')()
3
4 function fetchData(){
5   return new Promise((resolve,reject)=>{
6     fetch('https://api.bilibili.com/x/web-interface/ranking/region?rid=1&day=3&original=0')
7     .then(res=>{
8       return res.json()
9     }).then(res=>{
10       resolve(res)
11     })
12   })
13 }
14 app.get('/',(req,res)=>{
15   async function BiData(){
16     let data = await fetchData()
17     res.json({stats:200,data,message:'成功'})
18   }
19 })
```

```
18 }

19 BiData()

20 })

21

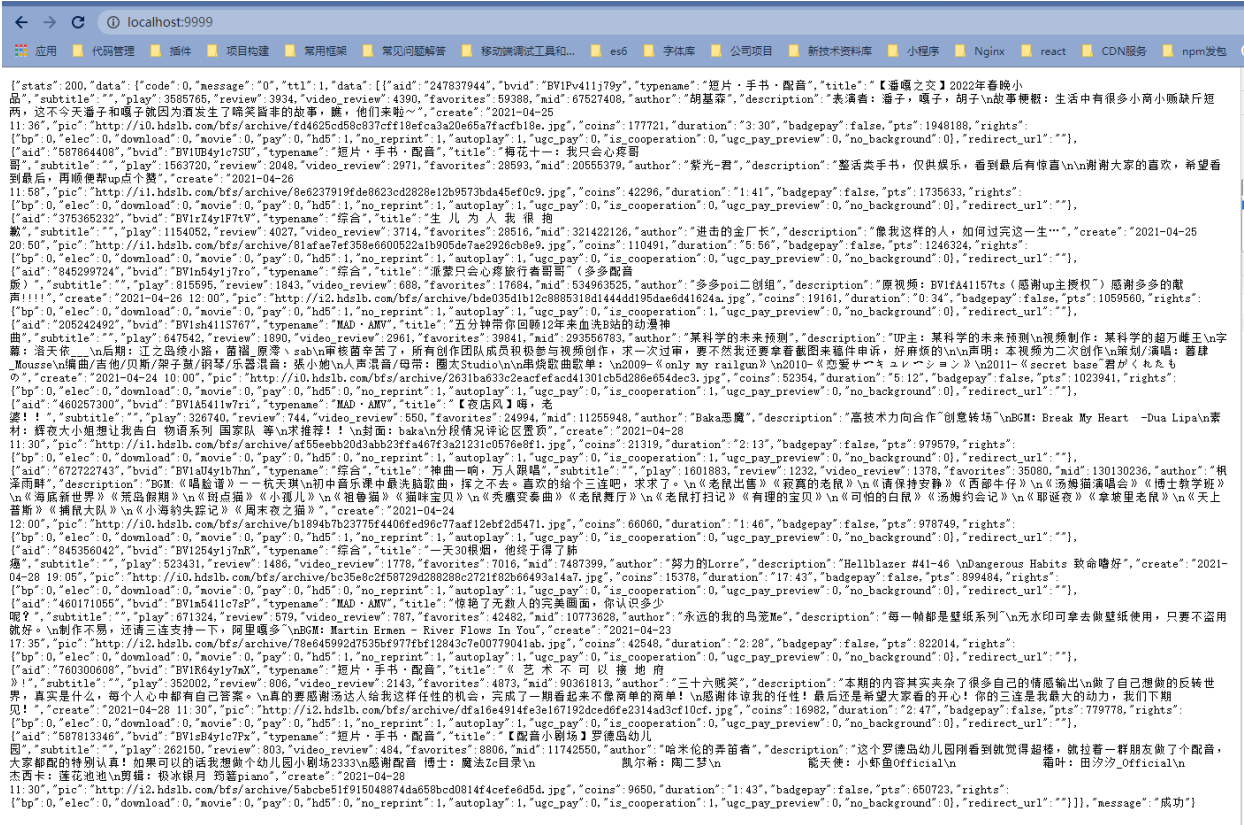
22

23 app.listen(9999,()=>{

24   console.log('9999服务');

25 })
```

爬虫结果：



但是有些需要登录的api就需要模拟登录了

```
{ "stats": 200, "data": { "code": -101, "message": "账号未登录", "ttl": 1 }, "message": "成功" }
```

模拟登录

手动添加上token再去请求数据。

```
1 const fetch = require('node-fetch')
```

```

2  const app = require('express')()
3
4  function fetchData(){
5    return new Promise((resolve,reject)=>{
6      fetch('https://api.bilibili.com/x/member/web/account',{
7        headers:{'cookie': '_uuid=00B3D359-2310-1DA1-5D34-E2161E2D27B670006infoc;
        buvid3=987581A0-8085-4350-9271-FCF22039A337185014infoc; bsource=search_baid
        u; buvid_fp=987581A0-8085-4350-9271-FCF22039A337185014infoc; SESSDATA=c365f
        90b%2C1635237967%2Ce5e94%2A41; bili_jct=bcf74cb7f953cbdcdb80b25eb5d6a6f9; D
        edeUserID=17868069; DedeUserID__ckMd5=cb00c0c926e3d66e; sid=a9g9b6qs; bfe_i
        d=6f2695e1895fb89897286b11ddc486b0; bp_video_offset_17868069=51781851868354
        1265; fingerprint3=a318682c183fc382c274fb22e5789b9f; fingerprint=336eb44fea
        45f179a053172b33f62e24; fingerprint_s=590ece7250adf7287ffc8227dda15e48; buv
        id_fp_plain=987581A0-8085-4350-9271-FCF22039A337185014infoc; PVID=2'}
8      })
9      .then(res=>{
10       return res.json()
11     }).then(res=>{
12       resolve(res)
13     })
14   })
15 }
16 app.get('/',(req,res)=>{
17   async function BiData(){
18     let data = await fetchData()
19     res.json({stats:200,data,message:'成功'})
20   }
21   BiData()
22 })
23
24
25 app.listen(9999,()=>{
26   console.log('9999服务');
27 })

```

上面说的这些都是固态话，要更加浏览器化，得需要重新像一个浏览器一样。需要一个浏览器内核这时候可以再github找到一个库Puppeteer谷歌内核一个库