

系统设计第三课

小憩一下，马上开课



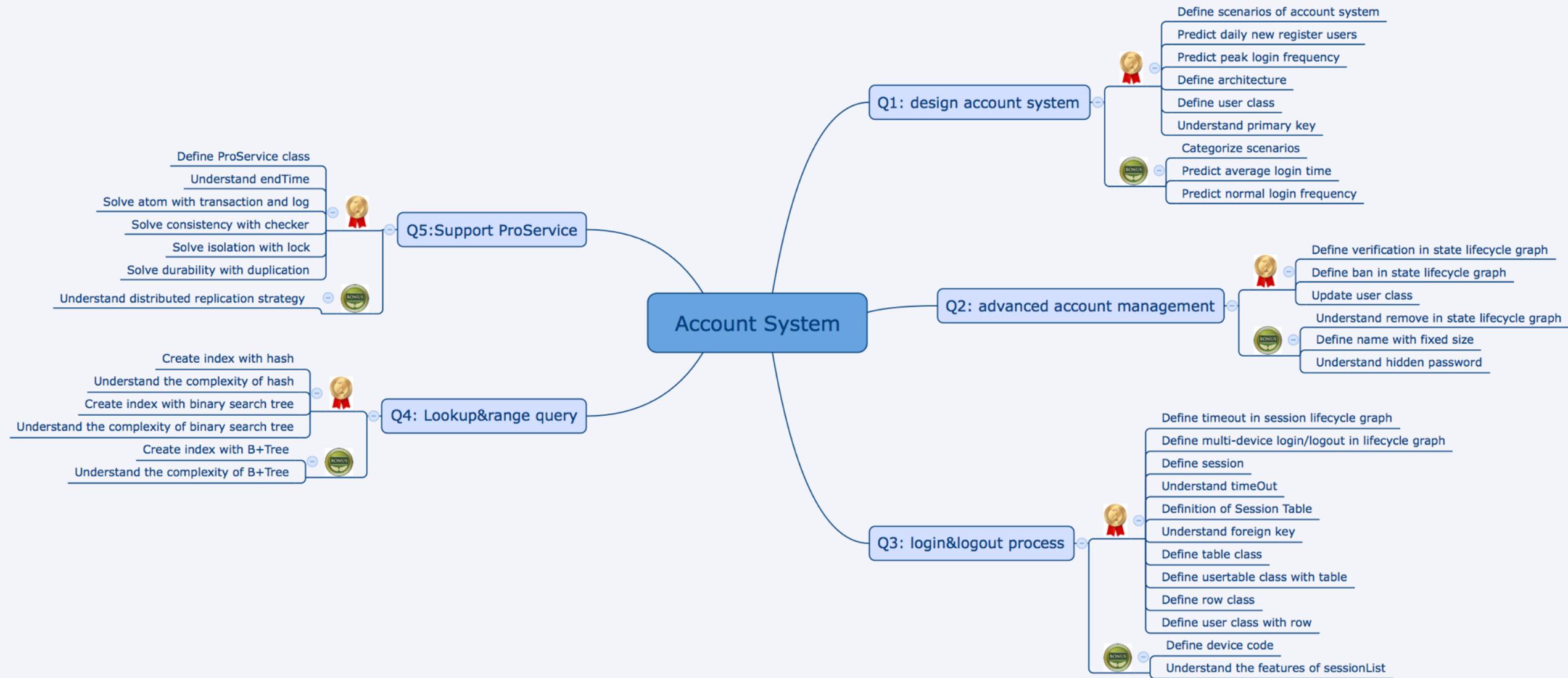
关注微信/微博，获取最新面试题及权威解答

微信: [ninechapter](#)

微博: <http://www.weibo.com/ninechapter>

官网: www.jiuzhang.com

Summary of class





System_Design_3

Web Crawler & Tiny URL

张无忌

2015-10-11

V 3.01

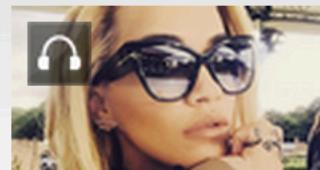
After class_3, you can answer

- Design crawler
 - Dropbox, Google, Turn, Alibaba
- Design thread-safe consumer and producer
 - Google, Amazon, TripAdvisor, Microsoft, Pure Storage, Dropbox, LinkedIn, Palantir, Intel, Bloomberg, ...
- Design TinyURL
 - LinkedIn, Uber, Bloomberg, Hulu

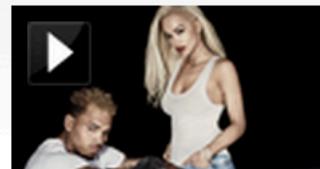
Interviewer

- Let's implement “News on Singers”

LATEST MUSIC NEWS



› EXCLUSIVE: Rita Ora Explains Why It Couldn't Work Between Her And Zayn Malik



› LISTEN: Rita Ora Tells Us What To Expect From Her 'Body On Me' Video Feat Chris Brown



› LISTEN: David Guetta Previews New Song 'Sun Goes Down' With Showtek & Magic!



› WATCH: Miley Cyrus Floats Through Space As She Teases Hosting The 2015 MTV VMA Awards



› Behind The Scenes - Was One Direction's 'Drag Me Down' Vid Shot At NASA Space Station?

[More Music News ➔](#)

••○○○ AT&T 17:31 100%

今日头条

本地 视频 社会 音乐 订阅 娱乐

2015.07.20 Monday 今天

日本摇滚女团香港开唱

热 手机新浪网 评论 0



舞蹈女神米拉内地首发唱片 星外星再推新疆流行音乐

热 星外星 评论 0



5分钟前

2015张北草原音乐节24日开唱 带你摇滚带你飞

热 中国青年网 评论 0

8分钟前

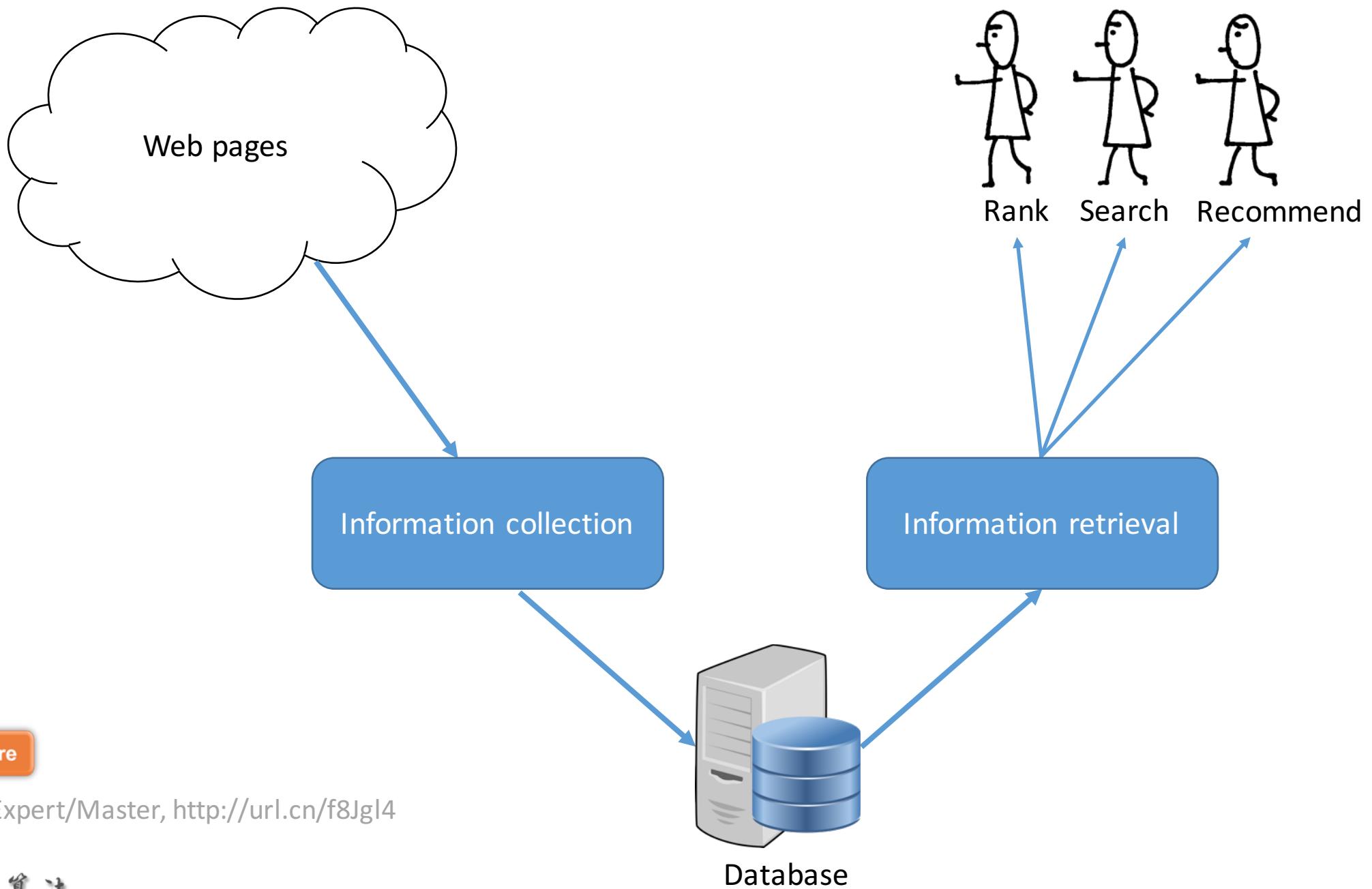
炎亚纶《新声代3》享受音乐 张铭浩被称小王源

热 红网 评论 4

10分钟前



7
Copyright © www.jiuzhang.com



Read More

Novice/Expert/Master, <http://url.cn/f8Jgl4>

Interviewer: crawl a page

音乐家聂耳逝世80周年 曾为国歌作曲

2015年07月17日 10:02

来源：凤凰音乐综合

0人参与 0评论



聂耳

1935年的7月17日，音乐家聂耳溺水逝世。他一生只经历了23个春秋，却创作出包括国歌《义勇军进行曲》、《卖报歌》、《金蛇狂舞》等许多优秀作品。

- Crawl a page from

http://yue.ifeng.com/a/20150717/39747647_0.shtml

Crawl a page in Python



```
import urllib2

# Request source file
url = 'http://yue.ifeng.com/a/20150717/39747647_0.shtml'
request = urllib2.Request(url)          #Write a letter
response = urllib2.urlopen(request)      #Send the letter and get the reply
page = response.read()                  #Read the reply

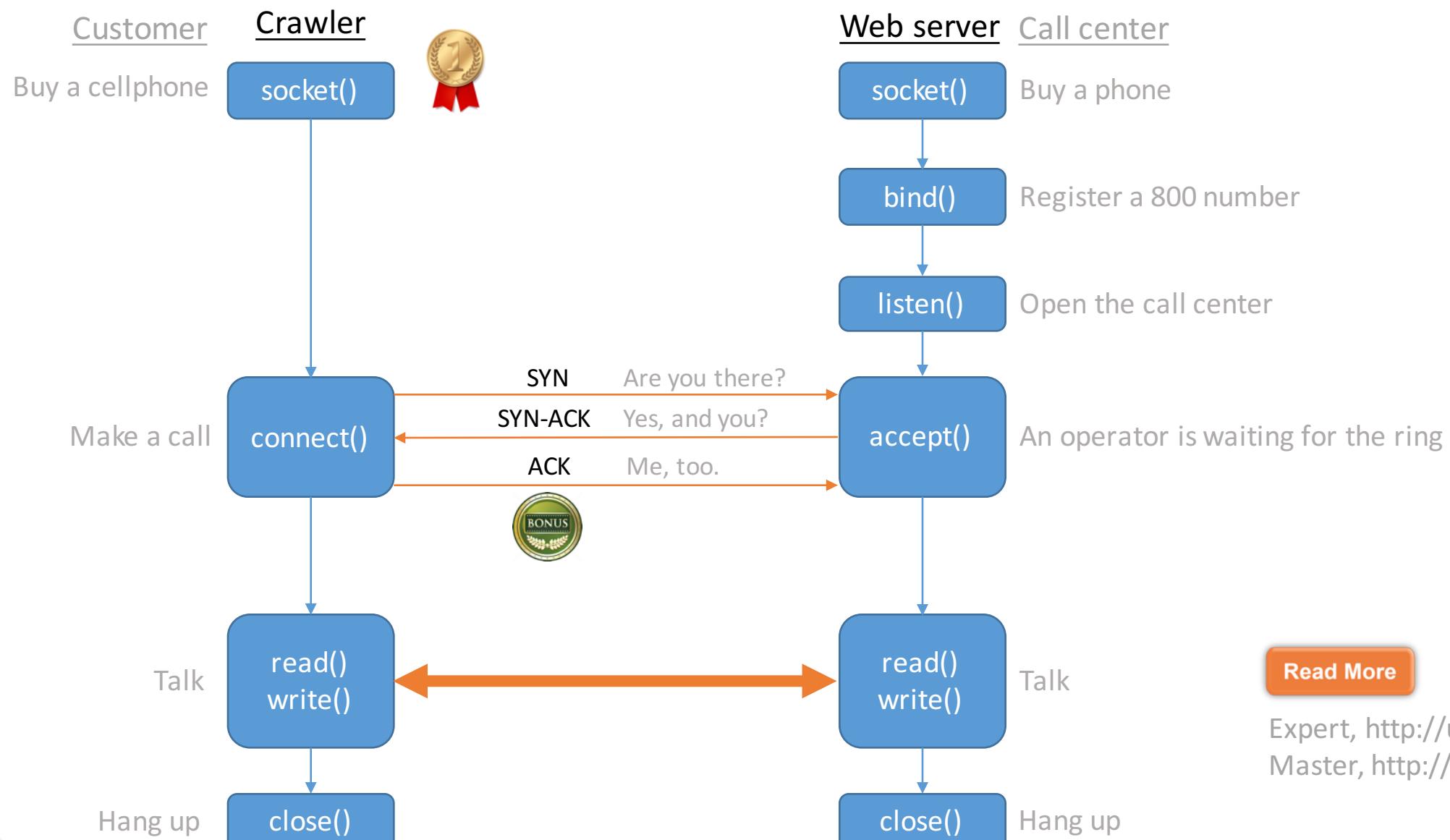
# Save source file
webFile = open('webPage.html', 'wb')
webFile.write(page)
webFile.close()
```

[Read More](#)

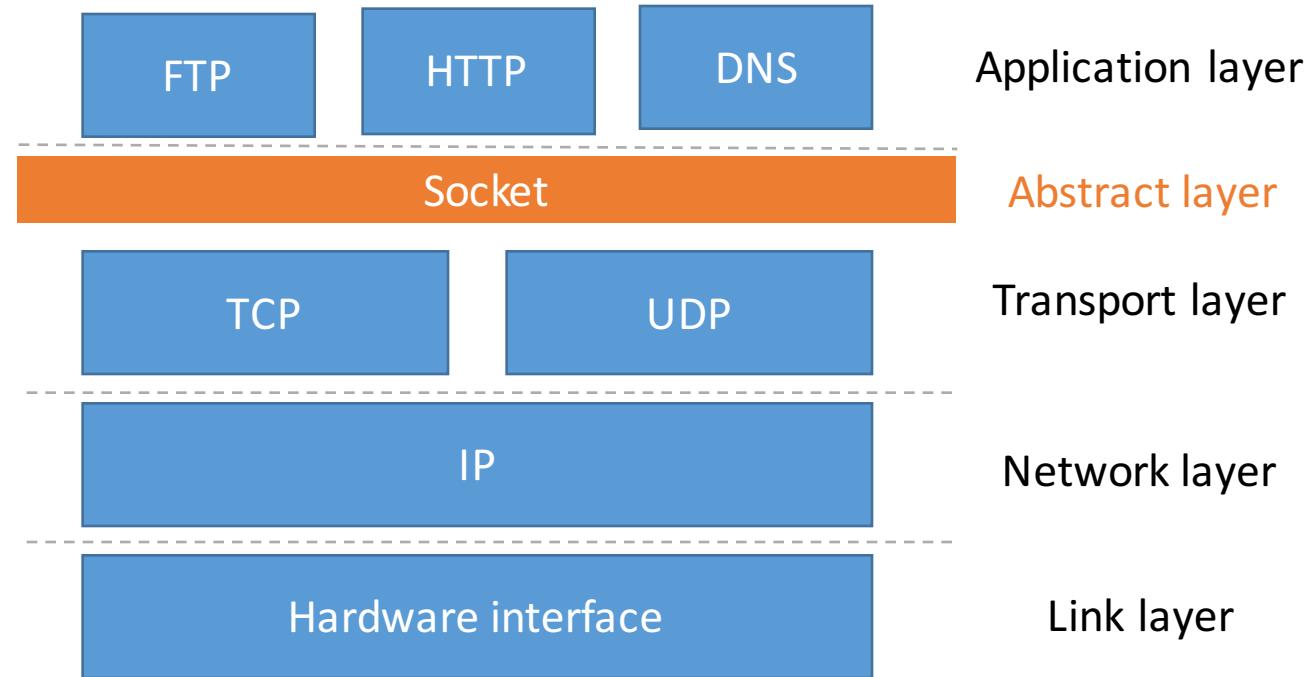
Novice/Expert/Master, <http://url.cn/NEHw7J>
Novice/Expert, <http://url.cn/ZANYur>
Expert, <http://url.cn/agT8mg>
Expert, <http://url.cn/SOBO8x>
Expert/Master, <http://url.cn/1eHONm>

Interviewer: What is the network process
when you are crawling a webpage?

Call center VS web server



Layers



Read More

Expert, <http://url.cn/ZX9hMs>
Master, <http://url.cn/RTPIgV>

Interviewer: What is HTML?

```

<!DOCTYPE html>
<html>
<head>
<!--cmpp:start-->
<meta charset="utf-8" />
<title>音乐家聂耳逝世80周年 曾为国歌作曲|聂耳_凤凰网</title>
<meta name="keywords" content="聂耳" />
<meta name="description" content="1935年的7月17日，音乐家聂耳溺水逝世。他一生只经历了23个春秋，却创作出包括国歌《义勇军进行曲》、《卖报歌》、《金蛇狂舞》等许多优秀作品。 聂耳资料： 聂耳(1912-1935)，原名聂守信，字" />
<meta name="og: webtype" content="news" />
<meta property="og:url" content="http://yue.ifeng.com/a/20150717/39747647_0.shtml" />
<meta property="og:title" content="音乐家聂耳逝世80周年 曾为国歌作曲" />
<meta property="og:description" content="1935年的7月17日，音乐家聂耳溺水逝世。他一生只经历了23个春秋，却创作出包括国歌《义勇军进行曲》、《卖报歌》、《金蛇狂舞》等许多优秀作品。 聂耳资料： 聂耳(1912-1935)，原名聂守信，字" />
<meta name="og:time" content="2015年07月17日 10:02" />
<meta name="og:category" content="凤凰音乐" />
<meta property="og:image" content="http://y3.ifengimg.com/a/2015_29/dfed5a6d8dd010f_size12_w300_h409.jpg" />
<meta name="og: img_slide" content="" />
<meta name="og: img_video" content="" />
<meta name="viewport" content="width=device-width, initial-scale=1.0, minimum-scale=1.0, maximum-scale=1.0, user-scalable=0" />
<meta http-equiv="mobile-agent" content="format=html5;url=http://i.ifeng.com/ifengurl.f?
vt=5&ch=rj_bd_me&url=http://yue.ifeng.com/a/20150717/39747647_0.shtml" />
<meta http-equiv="mobile-agent" content="format=xhtml;url=http://i.ifeng.com/ifengurl.f?
vt=2&ch=rj_bd_me&url=http://yue.ifeng.com/a/20150717/39747647_0.shtml" />
<!--广告代码--><script src="http://m1.ifengimg.com/ifeng/sources/inice_v1.js"></script>
<script>
window['recommendOpinion'] = function (selector){

    var selector = selector.split('_');
    if(selector.length >1){
        var num = jQuery('#' + selector[1]).html() /1;
        jQuery('#' + selector[1]).html(num +1);
    }
};

```

[Read More](#)



鳳凰 音乐 凤凰网音乐 >内地 >正文

音乐家聂耳逝世80周年曾为国歌作曲

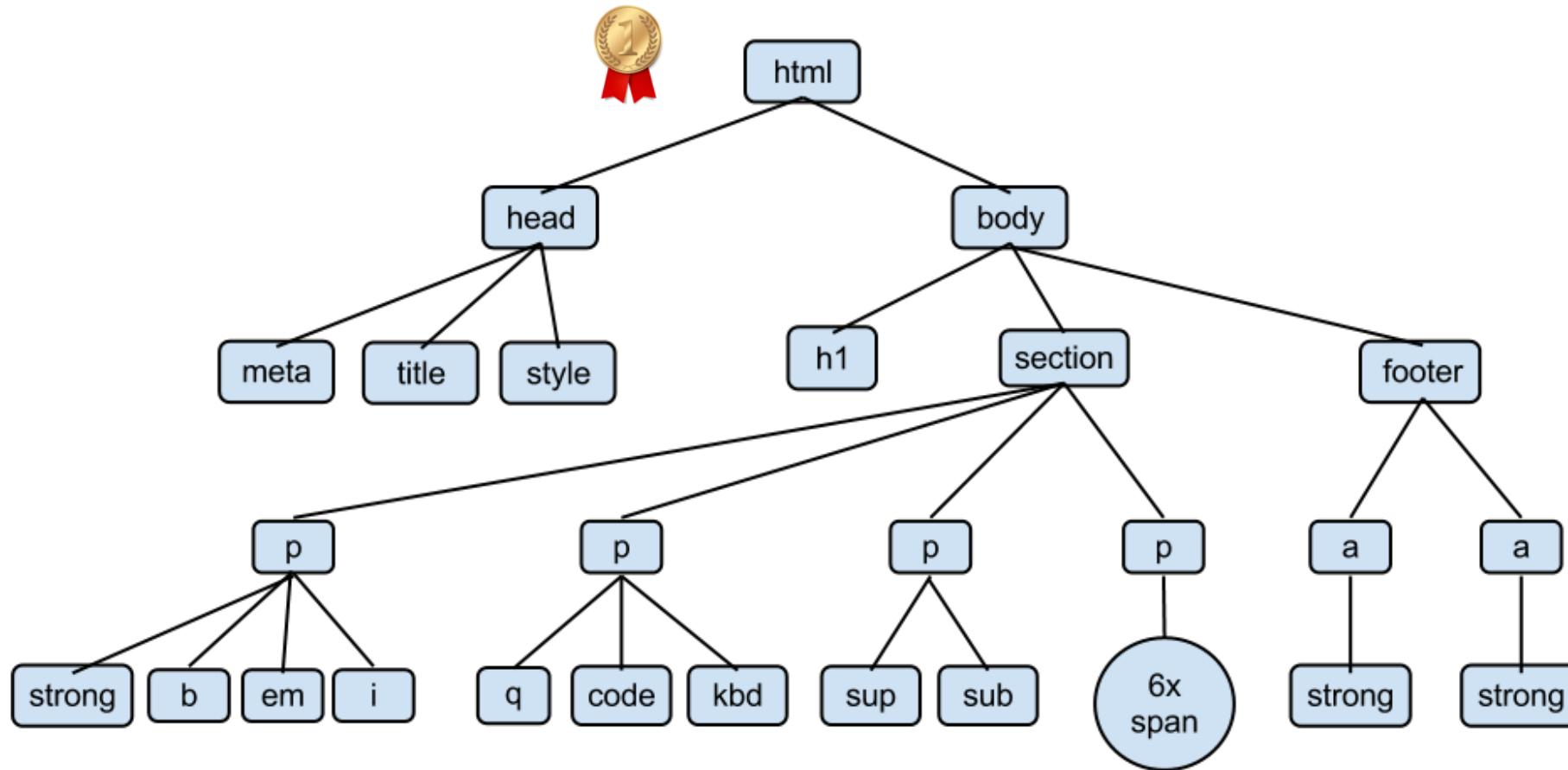
015年07月17日 10:02 来源：凤凰网综合 0人参与 0评论



聂耳

1935年的7月17日，音乐家聂耳溺水逝世。他一生只经历了23个春秋，却创作出包括国歌《义勇军进行曲》、《卖报歌》、《金蛇狂舞》等许多优秀作品。

HTML uses a tree to represent a human



Interviewer: crawl all the news
of a website



阿里成立音乐集团 高晓松宋柯加盟

高晓松：音乐职业经理人生涯最后一站 | 宋柯曾提：音乐已死

周杰伦曾想要儿子 好友赞女儿贴心改其看法

女儿五官遗传昆凌长睫毛 | 粉丝取名“周凌”英文名Jolin

音乐家聂耳逝世80周年 曾为国歌作曲

乐评人质疑姚贝娜假唱称其死不瞑目引不满

朝鲜首次批准外国乐队演出 八月献演两场

爱疯音乐家

张洪量：心如止水鉴常明



我希望我们的音乐要领导别人，而不是被领导；我也希望我们能更贡献好的音乐给全世界人民，做出文化上的外销品。 [详细]

- 黄雅莉：真命天子快出来！
- 张力尹：我不要只当“花瓶”

Crawl all the news from

- <http://yue.ifeng.com>

Identify the list page



- http://yue.ifeng.com/news/list_0/0.shtml

凤凰网 音乐 [凤凰网音乐 > 滚动新闻](#)

[滚动新闻](#)

·专访歌手刘炫颉：畅谈《吉祥经》的创作之路 [01-13 17:30]
·“歌手3”陈洁仪遭淘汰 歌迷力挺求开个唱 [01-10 03:21]
·唱作歌手吴梓涵—新专辑暖心发声 [01-09 40:29]
·罗中旭疑释前嫌献唱旧爱 好歌曲第二期看点全面开扒 [01-09 26:52]
·叶一茜发EP同届超女齐祝福 田亮惊喜现身深情献吻 [01-09 46:42]

·“九公子”创始人荆纯：创业不能只靠粉丝经济(图) [01-08 52:48]
·安以轩王岳伦将出席G客盛典胶泥定格动画《风雪山神庙》备受关注 [01-07 29:39]
·音乐熊猫诗词儿歌歌手选拔赛广州决赛圆满举行 [01-07 20:53]
·第六届香港金紫荆杯普通话朗诵比赛圆满落幕 [01-07 34:42]
·首部舞台芭比真人秀音乐剧《闪亮芭比》发布会召开启动全国巡演 [01-06 03:00]

·何润东焦恩俊联袂助阵“微商春晚” [01-06 01:44]
·我是歌手孙楠全力开唱 热拍爆赛前搞笑花絮 [01-06 58:04]
·The Push推乐队摘第15届华语音乐传媒大奖 [01-06 04:42]
·《有声有色》袁惟仁：那英对我太简单粗暴 [01-06 40:28]
·“大马之光”张祖诚北京发片 “暖男情歌”获盛誉 [01-06 39:05]

·榭霖全新歌曲诠释《旧爱》新欢的感情纠葛 [01-05 24:31]
·专注原创 金立ELIFE携手《中国好歌曲》首亮荧屏 [01-04 33:36]
·肖明超发布明星指数：黄晓明口碑佳 王俊凯粉丝最活跃 [12-31 42:48]
·张信哲牵手顶级人气动漫《秦时明月》被赞烧脑 [12-31 39:04]
·张皓玥全新单曲《哎呀爱呀》MV强势来袭 [12-31 35:17]

·“我的根在草原”齐峰2015草原情怀1月5日上海唱响 [12-30 05:19]

<首页 [5127](#) [5126](#) [5125](#) ... [3](#) [2](#) [1](#) 下页 > 末页 >

Identify the links of news

```
import re
```

BONUS

```
pattern = re.compile('<a href="http://yue.ifeng.com/news/detail_(.*?)" target="_blank">(.*)</a>', re.S)
items = re.findall(pattern, page)
```

```
for item in items:
```

```
    print item[0],item[1]
```

1

```
262 <li><a href="http://yue.ifeng.com/news/detail_2015_01/10/40163401_0.shtml"
263 target="_blank">·歌手3“陈洁仪遭淘汰 歌迷力挺求开个唱</a><span>[01-10 03:21]</span></li>
264 <li><a href="http://yue.ifeng.com/news/detail_2015_01/09/40155354_0.shtml"
265 target="_blank">·唱作歌手吴梓涵-新专辑暖心发声</a><span>[01-09 40:29]</span></li>
266 <li><a href="http://yue.ifeng.com/news/detail_2015_01/09/40155188_0.shtml"
target="_blank">·罗中旭疑释前嫌献唱旧爱 好歌曲第二期看点全面开扒</a><span>[01-09 26:52]</span></li>
267 <li><a href="http://yue.ifeng.com/news/detail_2015_01/09/40154764_0.shtml"
target="_blank">·叶一茜发EP同届超女齐祝福 田亮惊喜现身深情献吻</a><span>[01-09 46:42]</span></li>
268 <div class="space10"></div>
269 <li><a href="http://yue.ifeng.com/news/detail_2015_01/08/40141951_0.shtml"
target="_blank">·“九公子”创始人荆纯：创业不能只靠粉丝经济(图)</a><span>[01-08 52:48]</span></li>
```

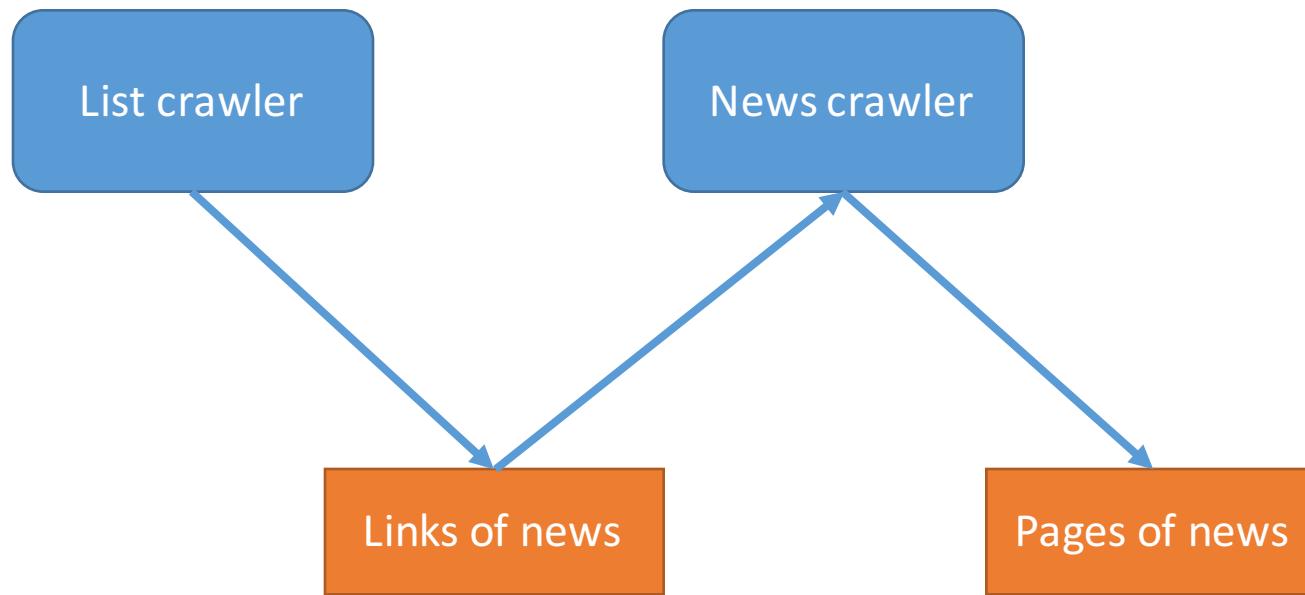
Read More

Novice/Expert, http://url.cn/TWdTAp

Novice/Expert, http://url.cn/a8U9VH

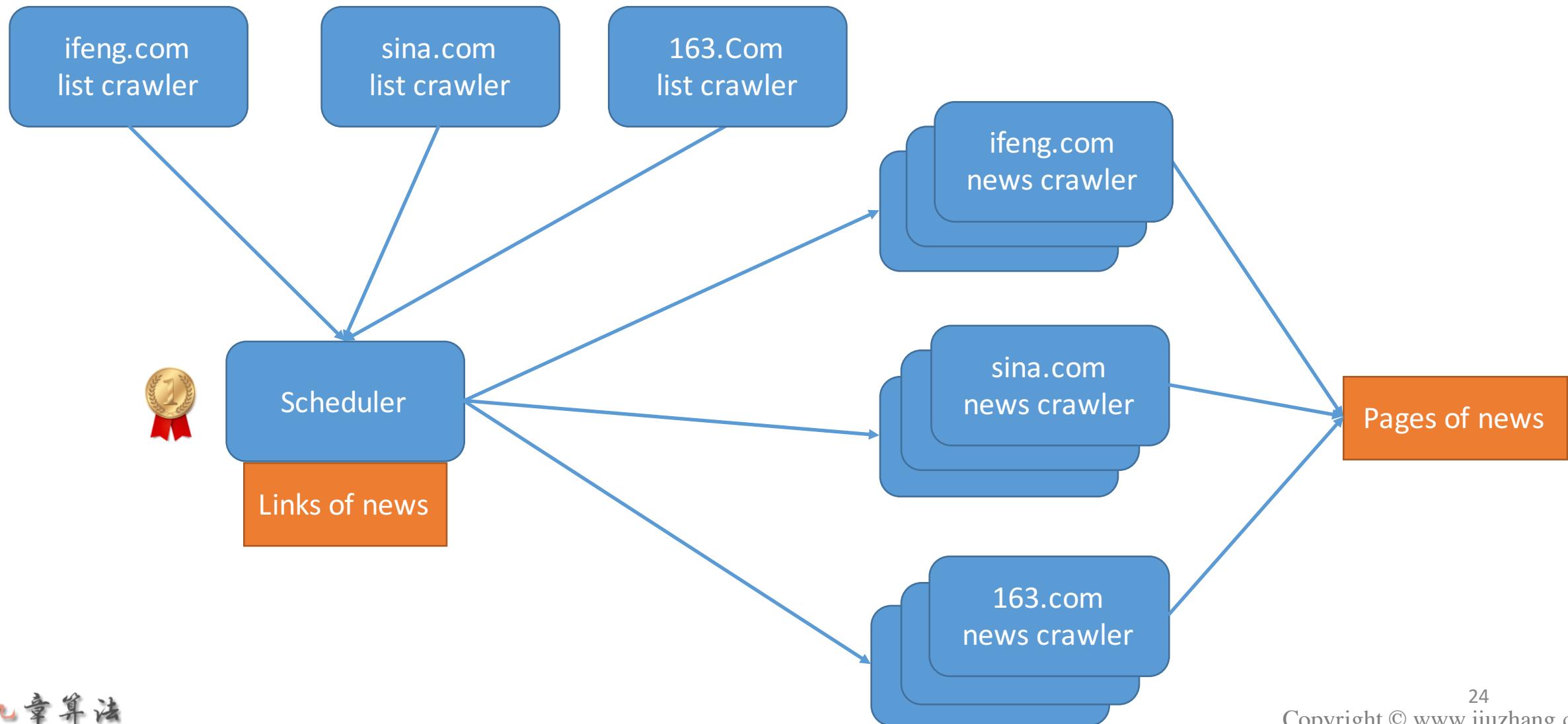
Expert/Master, http://url.cn/dl93as

Architecture (v1)



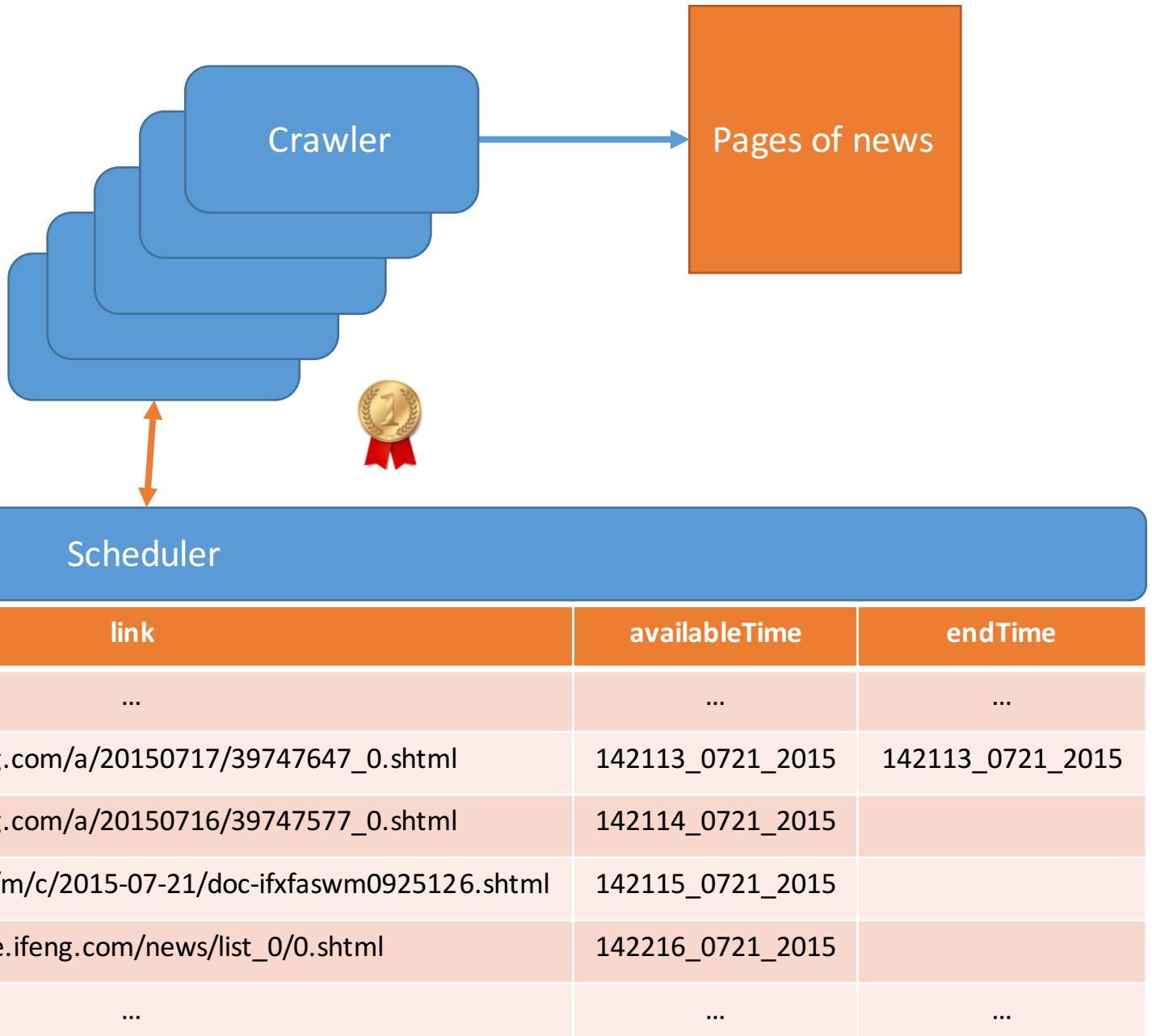
Interviewer: crawl from more websites

Architecture (v2)



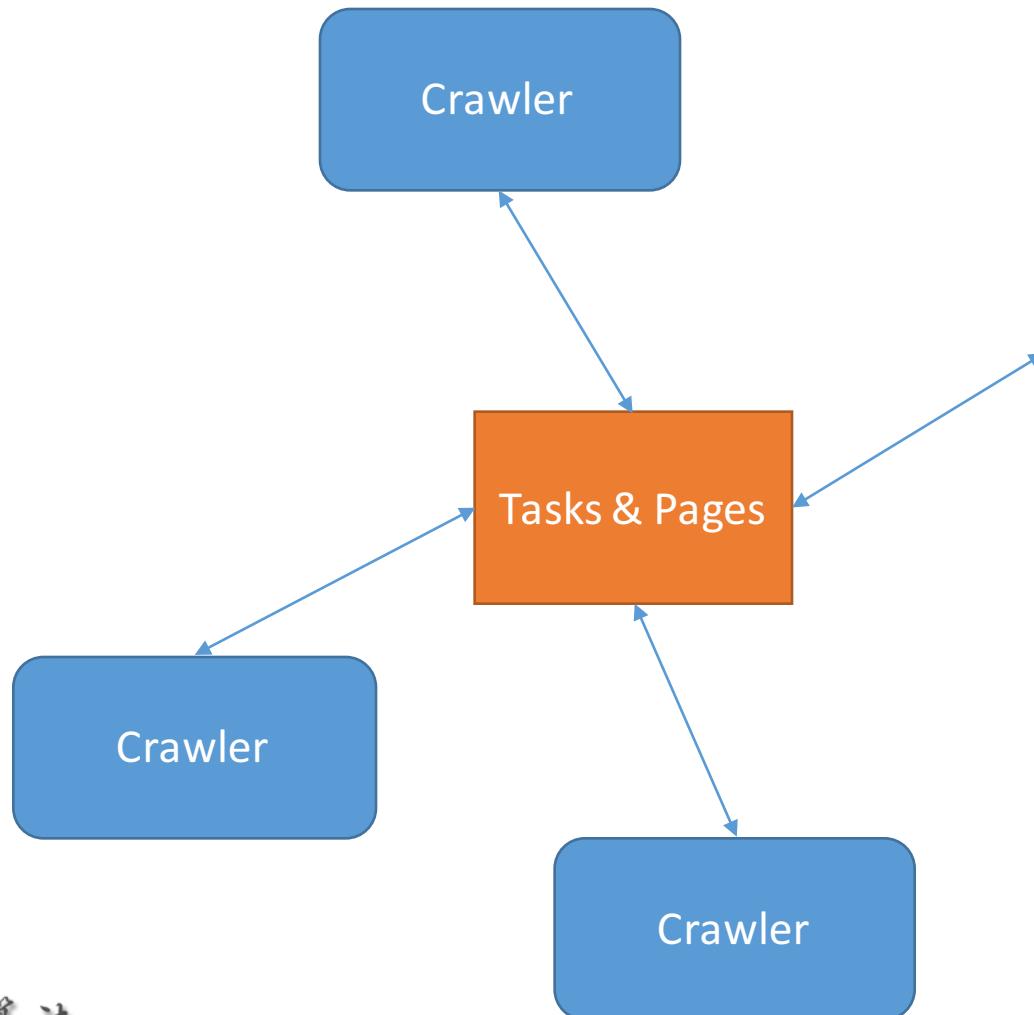
Interviewer: Your crawlers are wasted!

Architecture (v3)



Interviewer: design scheduler?

Solution with sleep



```
Crawler
While(true)
Lock(taskTable)
While taskTable.Find("state=new") == NULL
Release(taskTable)
Sleep(1s)
Lock(taskTable)

task = taskTable.FindOne("state=new")
task.state = "working"
Release(taskTable)

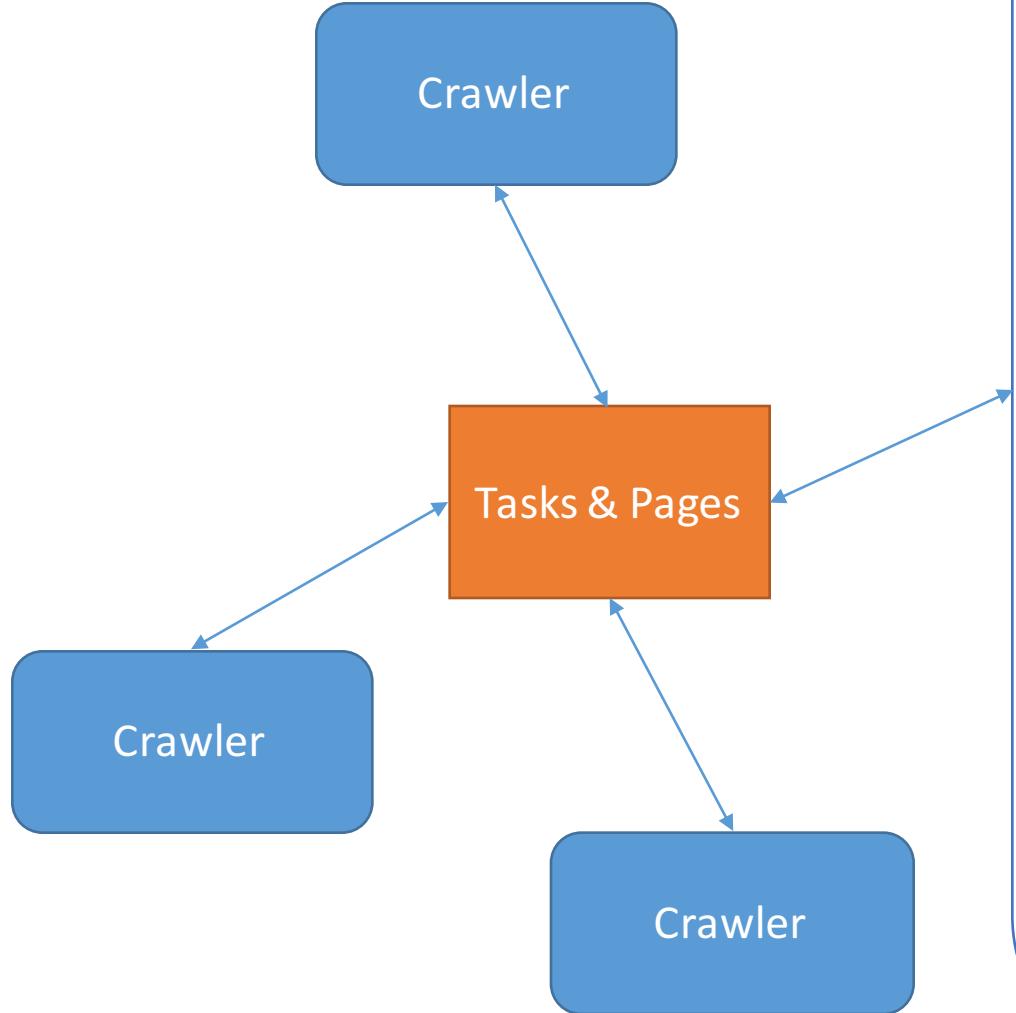
page = Crawl(task.url)

If task.type == "list"
Lock(taskTable)
For newTask In page:
taskTable.Add(newTask)
task.state = "done"
Release(taskTable)
Else
Lock(pageTable)
pageTable.Add(page)
Release(pageTable)
Lock(taskTable)
task.state = "done"
Release(taskTable)
```



Interviewer: design scheduler with
conditional variable

Solution with conditional variable



Crawler

```
While(true)
  Lock(taskTable)
  While( taskTable.Find("state=new") == NULL )
    Cond_Wait (cond, taskTable)
    task = taskTable.FineOne("state=new")
    task.state = "working"
    Release(taskTable)

    page = Crawl(task.url)

    If task.type == "list"
      Lock(taskTable)
      For newTask In page:
        taskTable.Add(newTask)
      Cond_Signal(cond)
      task.state = "done"
      Release(taskTable)
    Else
      Lock(pageTable)
      pageTable.Add(page)
      Release(pageTable)
      Lock(taskTable)
      task.state = "done"
      Release(taskTable)
```



Conditional variable

Cond_Wait(cond, mutex)

```
Lock(cond.threadWaitList);
cond.threadWaitList.Add(this.thread);
Release(cond.threadWaitList);
Release(mutex);
Block(this.thread);
Lock(mutex);
```

Cond_Signal(cond)

```
Lock(cond.threadWaitList);
If( cond.threadWaitList.size()>0 );
thread = cond.threadWaitList.Pop();
Wakeup(thread);
```

```
Release(cond.threadWaitList);
```

Read More

Novice, <http://url.cn/44qusn>

Novice, <http://url.cn/cYsLZ3>

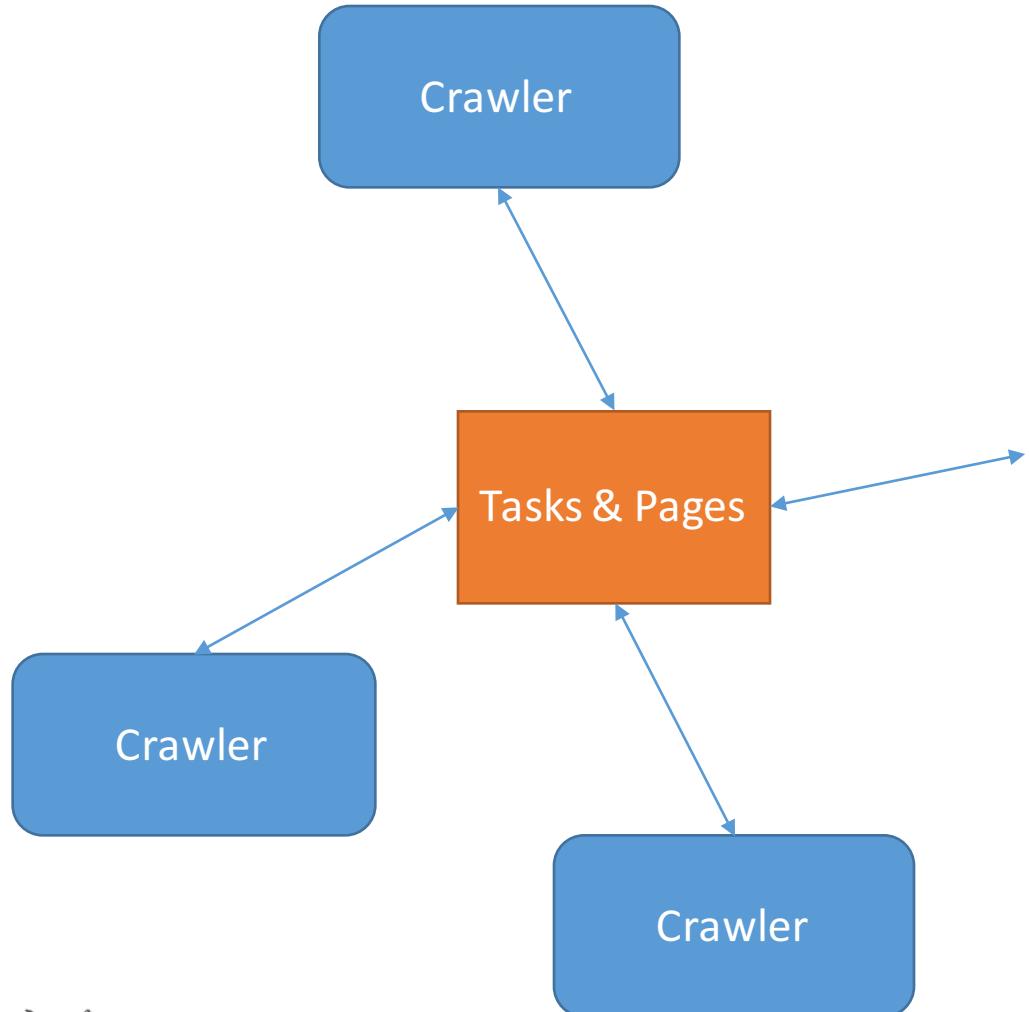
Expert, <http://url.cn/coAAeW>

Expert, <http://url.cn/TNVxRu>

Expert, <http://url.cn/edly1C>

Interviewer: design scheduler with semaphore

Solution with semaphore



Crawler

```
While(true)
    Wait(numberOfNewTask)
    Lock(taskTable)
    task = taskTable.FindOne("state=new")
    task.state = "working"
    Release(taskTable)

    page = Crawl(task.url)

    If task.type == "list"
        Lock(taskTable)
        For newTask In page:
            taskTable.Add(newTask)
            Signal(numberOfNewTask)
            task.state = "done"
            Release(taskTable)
    Else
        Lock(pageTable)
        pageTable.Add(page)
        Release(pageTable)
        Lock(taskTable)
        task.state = "done"
        Release(taskTable)
```



Semaphore

Wait(semaphore)

Lock(semaphore);

 semaphore.value--;

If(semaphore.value<0)

 semaphore.processWaitList.**Add(this.process);**

Release(semaphore);

Block(this.process);

Else

Release(semaphore);

Signal(semaphore)

Lock(semaphore);

 semaphore.value++;

If(semaphore.value<=0)

 process = semaphore.processWaitList.**Pop();**

Wakeup(process);

Release(semaphore)

Read More

Novice, <http://url.cn/XWQJWo>

Novice, <http://url.cn/g7xO6G>

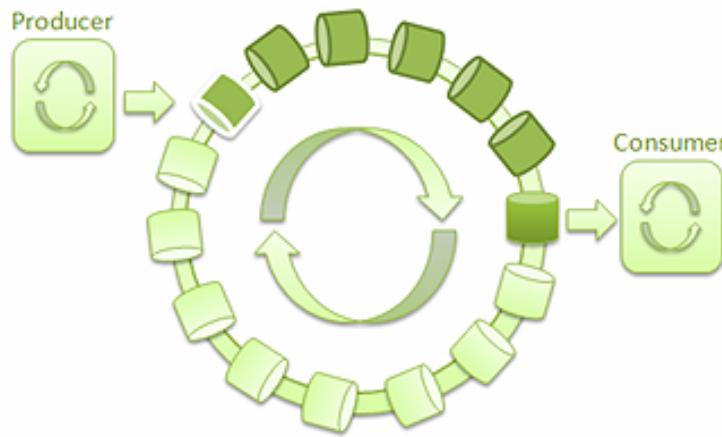
Expert, <http://url.cn/daQTsi>

Expert, <http://url.cn/YUG6Nd>

Expert/Master, <http://url.cn/a9NI3i>

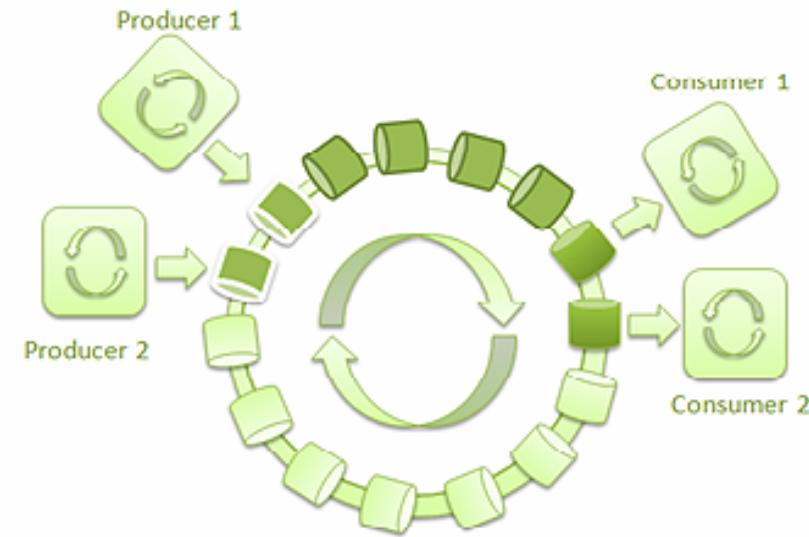
Interviewer: design the fastest consumers and producers

Thread-safe consumer and producer



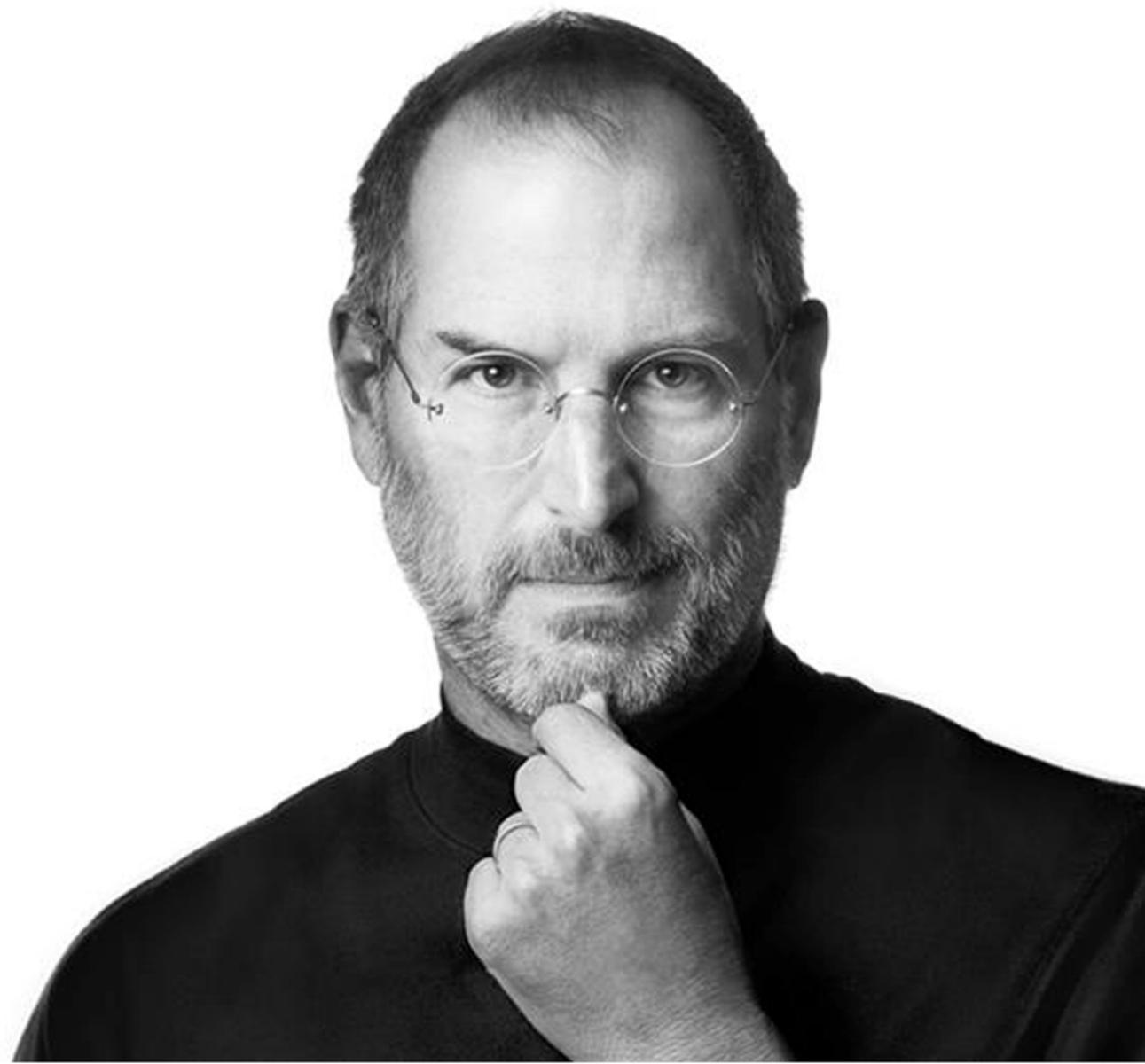
[Read More](#)

Novice, <http://url.cn/cNG8WK>
Expert, <http://url.cn/cFPpaN>
Expert, <http://url.cn/C1llpq>
Expert/Master, <http://url.cn/6CCw4i>
Master, <http://url.cn/9kdFF2>
Master, <http://url.cn/ORvrlX>
Master, <http://url.cn/dYRtFj>



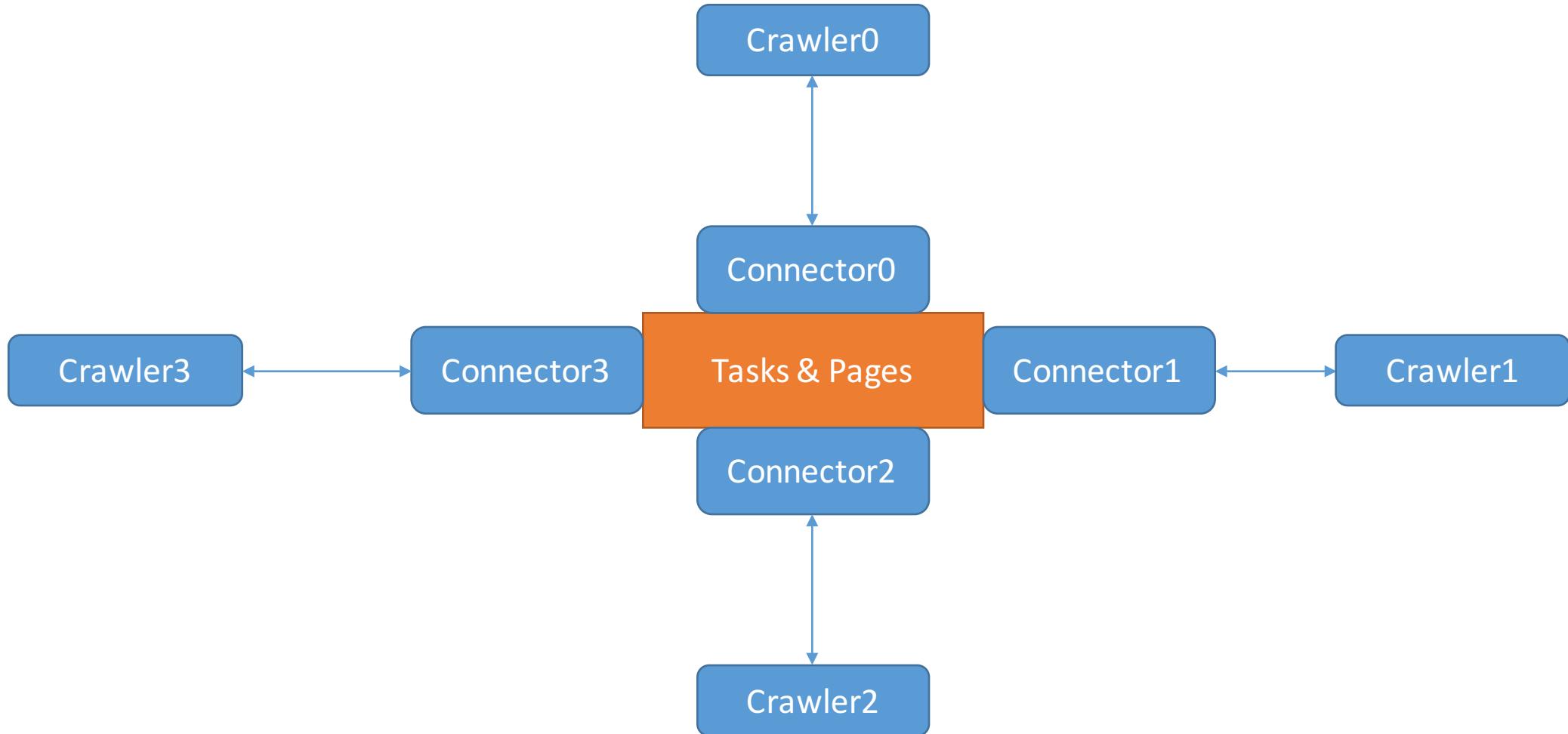
LMAX Disruptor

Stay Hungry,
Stay Foolish.



Interviewer: distribute crawlers
in multiple machines

Solution with one-to-one





Connector0

```
crawler = WaitConnection()
```

```
While(true)  
  Wait(numberOfNewTask)  
  Lock(taskTable)  
  task = taskTable.FindOne("state=new")  
  task.state = "working"  
  Release(taskTable)  
  
  Send(crawler0, task)  
  page = Receive(crawler0)
```

...

...

...

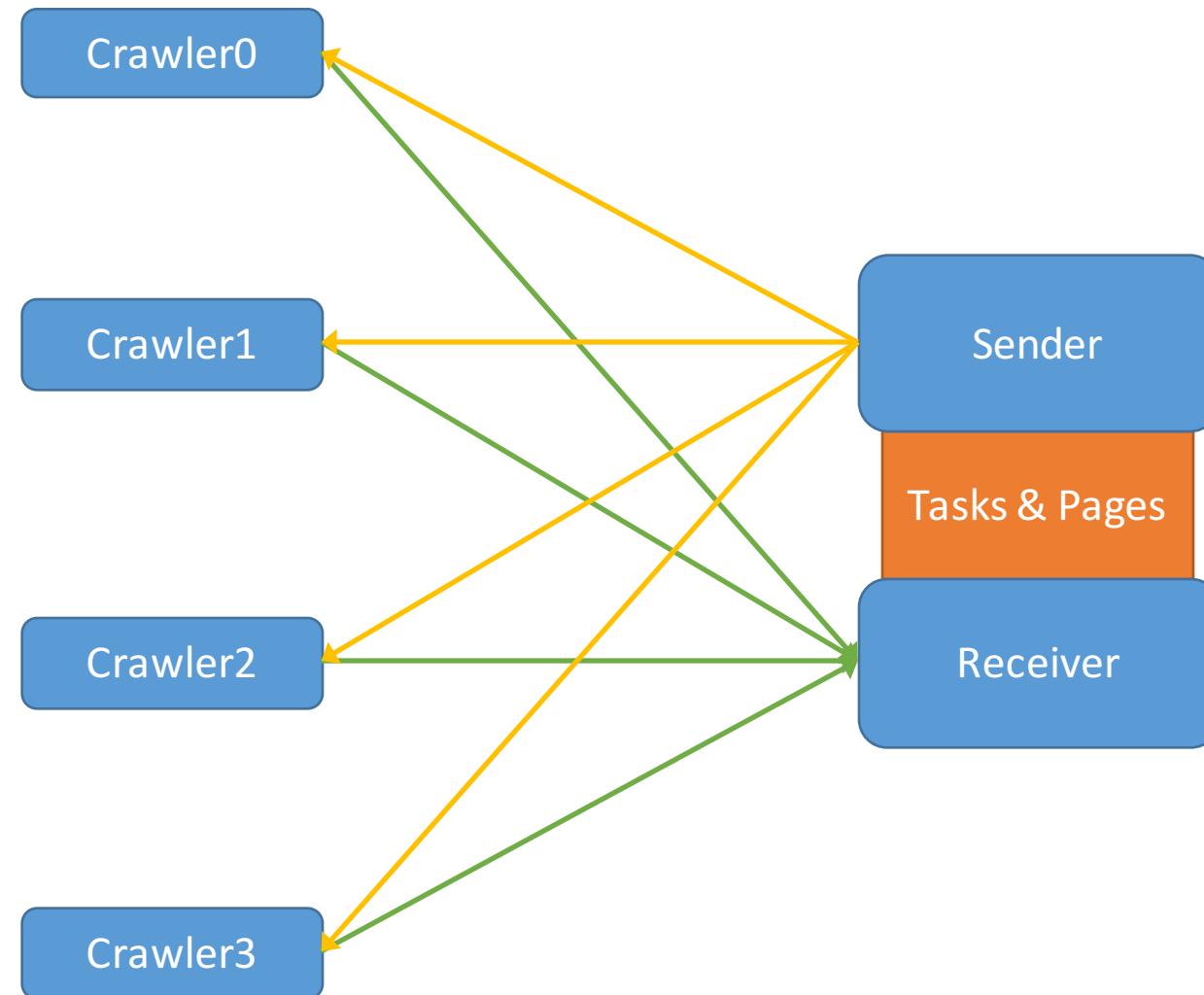
Tasks & Pages

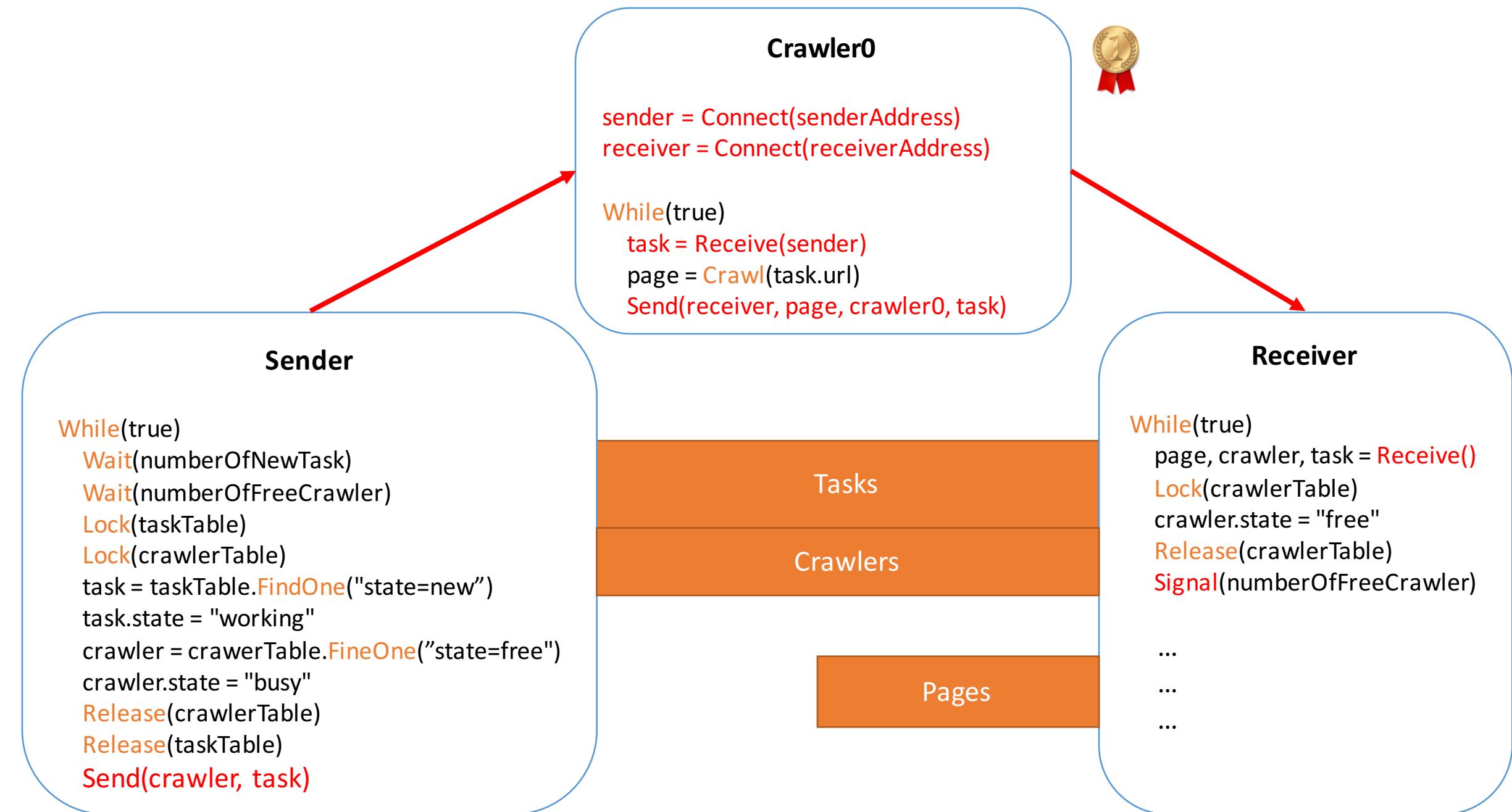
Crawler0

```
connector = Connect(connector0Address)  
  
While(true)  
  task = Receive(connector)  
  page = Crawl(task.url)  
  Send(connector, page)
```

Interviewer: Can you reduce the number of connectors?

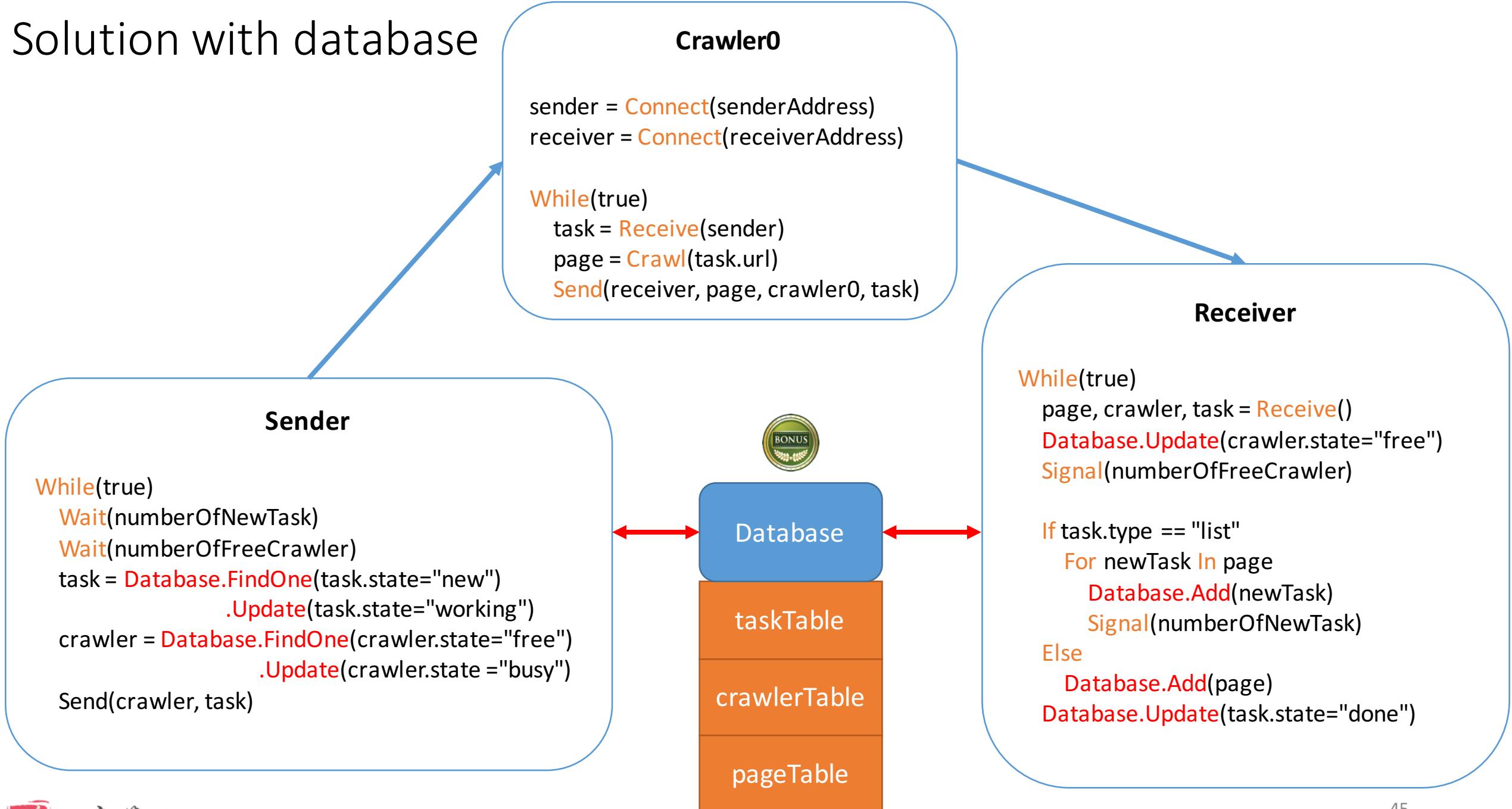
Solution with one-to-many



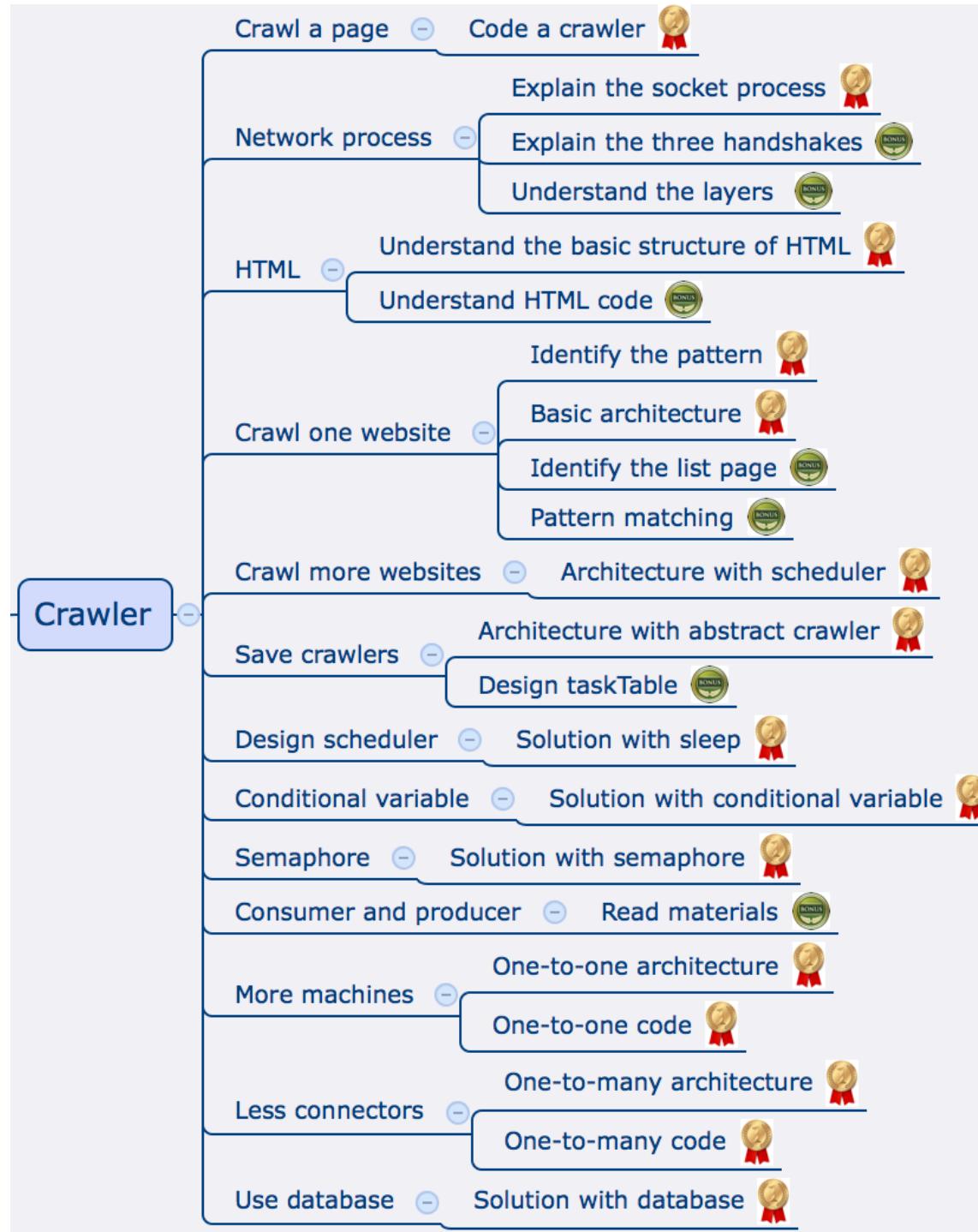


Interviewer: Can you use database?

Solution with database



Crawler Summary (15+7)



Interviewer: design Tiny URL

[Read More](#)

Novice, <https://goo.gl/>

Novice/Expert, <http://url.cn/bycYvq>

Novice/Expert/Master, <http://url.cn/aXRa8e>

Expert, <http://url.cn/euBds7>

Let's do it together ☺

Google url shortener

Paste your long URL here:

Google

Shorten URL

All goo.gl URLs and click analytics are public and can be accessed by anyone.

Clicks for the past: [two hours](#) | [day](#) | [week](#) | [month](#) | [all time](#)

<input type="checkbox"/>	LONG URL	CREATED	SHORT URL	CLICKS
<input type="checkbox"/>	www.jiuzhang.com/course/2	3 days ago	goo.gl/9OTcuV	Details 563

Crack a design in 5 steps

- Scenario: case/interface
- Necessary: constrain/hypothesis
- Application: service/algorith
- Kilobit: data
- Evolve



Interviewer: What is the Scenario?

Scenario



short_url Insert(long_url)
long_url Lookup(short_url)

Interviewer: What is the **Necessary**?

Necessary for 1million daily users



- Daily active users
 - 1,000,000
- Insert
 - Per day: $1,000,000 * 1\%(\text{function usage}) * 10(\text{function frequency}) = 100,000$
 - Per year: $100,000 * 365 = 36,500,000$
 - Per second: $100,000 / 86400 = 1.2$
- Lookup
 - Per day: $1,000,000 * 100\%(\text{function usage}) * 3(\text{function frequency}) = 3,000,000$
 - Per second: $3,000,000 / 86400 = 35$

Interviewer: What is the Algorithm?

Insert



LongToShort

longURL	shortURL
www.jiuzhang.com/course/2/	...
www.inoreader.com/all_articles	...
blog.csdn.net/longyulu/article/details/9159589	...
...	...

ShortToLong

shortURL	shortURL
...	www.jiuzhang.com/course/2/
...	www.inoreader.com/all_articles
...	blog.csdn.net/longyulu/article/details/9159589
...	...

```
class Shortener{
```

```
    map<string,string> mLongToShort
```

```
    map<string,string> mShortToLong
```

```
    string Insert( string longURL )
```

```
    If( mLongToShort.Find( longURL ) == NULL )
```

```
        string shortURL = GenerateShortURL();
```

```
        mLongToShort[ longURL ] = shortURL;
```

```
        mShortToLong[ shortURL ] = longURL;
```

```
    Return shortURL;
```

```
};
```

GenerateShortURL



LongToShort

longURL	shortURL
www.jiuzhang.com/course/2/	0
www.inoreader.com/all_articles	1
http://blog.csdn.net/longyulu/article/details/9159589	2
...	...

```
string GenerateShortURL()  
Return string( mLongToShort.size() );
```

ShortToLong

shortURL	shortURL
0	www.jiuzhang.com/course/2/
1	www.inoreader.com/all_articles
2	http://blog.csdn.net/longyulu/article/details/9159589
...	...

Interviewer: How to reduce
the size of shortURL?

Selection of characters



	Before	After
Yearly URL	36,500,000	36,500,000
Usable characters	$[0-9] = 10$	$[0-9a-zA-Z] = 62$
Encoding length	$\log_{10}(36,500,000) = 7.6 = 8$	$\log_{62}(36,500,000) = 4.2 = 5$
Example	goo.gl/36500000	goo.gl/2t9jG

GenerateShorterURL



```
string GenerateShortURL( string longURL )  
    Return ConvertTo62( LongToShort.size() );
```

```
string ConvertTo62(int number)  
    char Encode[62] = {'0',...,'9', 'a', ...,'z', ..., 'A',...,'Z'};  
    string ret = "";  
    While(number)  
        ret = Encode[ number%62 ] + ret;  
        number/=62;  
    Return ret;
```

Interviewer: What is the Kilobit?



Average size of longURL = 100 bytes

Average size of shortURL = 4 bytes (int)

State = 4 byte

Daily new URL = $100,000 * 108 = 10.8\text{MB}$

Yearly new URL = $10.8 * 365 = 4\text{GB}$

URL table (MySQL or NoSQL)

longURL (varchar)	shortURL (int)	state (int)
www.jiuzhang.com/course/2/	0	0
www.inoreader.com/all_articles	1	0
http://blog.csdn.net/longyulu/article/details/9159589	2	1
...	...	0

Follow up

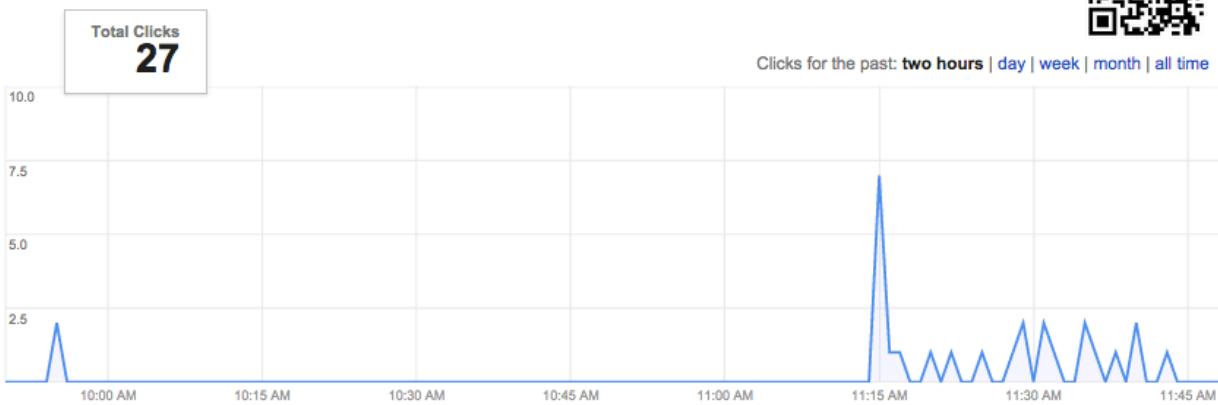


- How to support random?
 - Random(0,range)
- How to avoid conflicting?
 - Try again
- How to implement time-limited service?
 - Expire/state
- How to cache?
 - Pre-load
 - Replacement

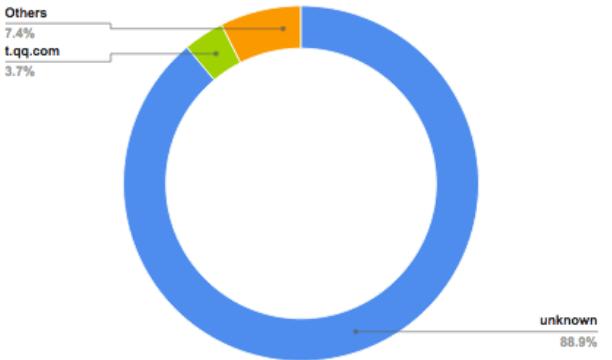
[Read More](#)

Master, <http://url.cn/dr4uux>

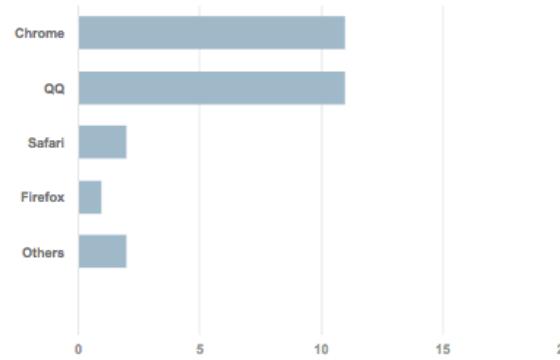
Interviewer: How to support analysis?



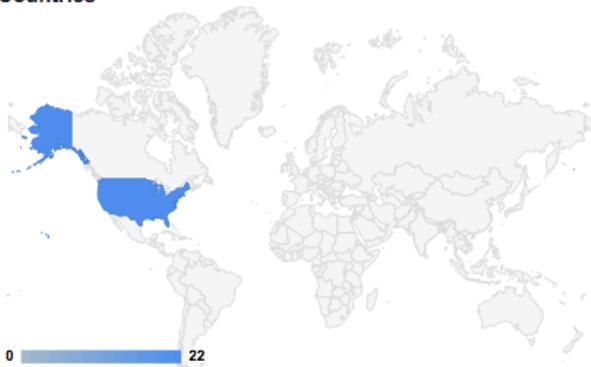
Referrers



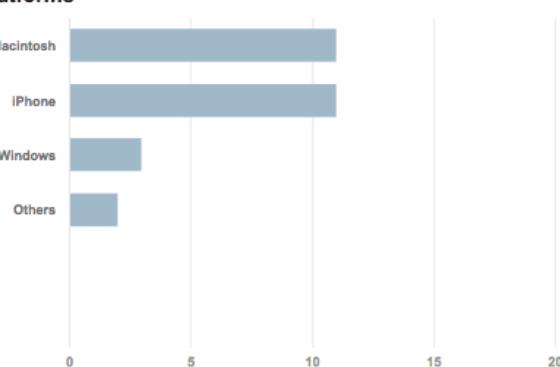
Browsers

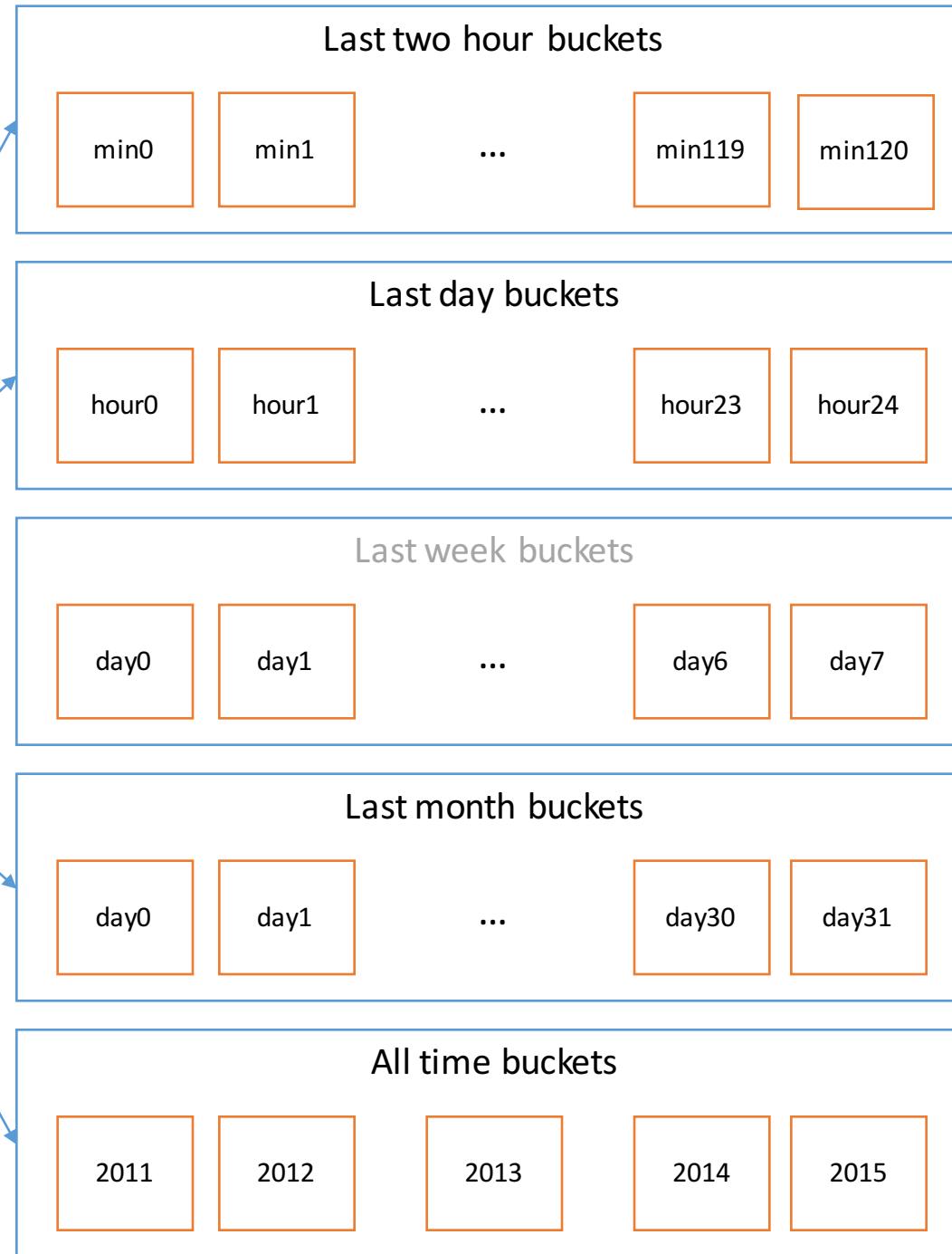


Countries

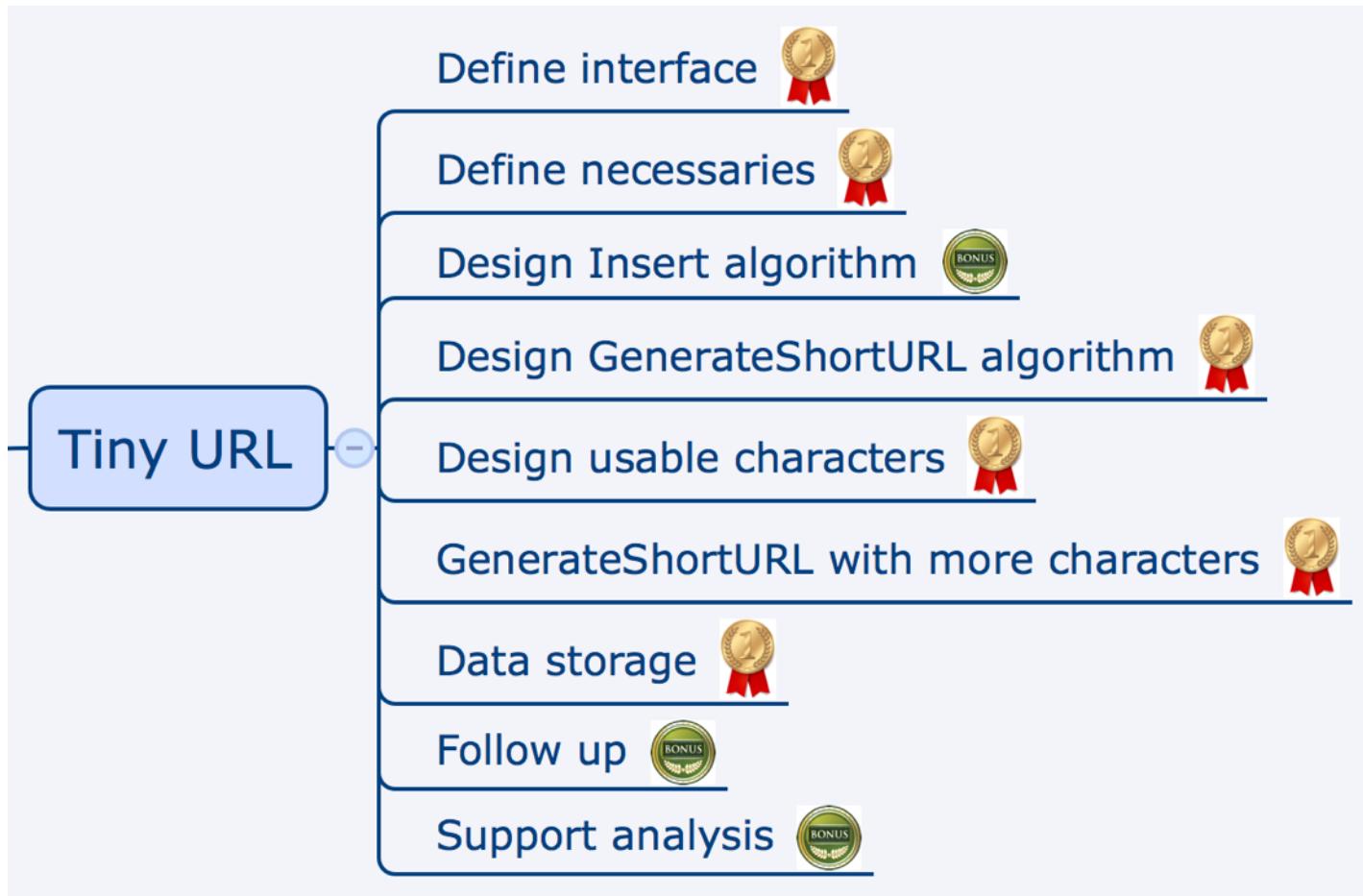


Platforms





Summary of Tiny URL





c o n g r a t s , h e r e ' s y o u r l i n k

<http://虎牙表情包.网址>

MOJI

click it and ctrl + c to copy

MAKE ANOTHER

00 views

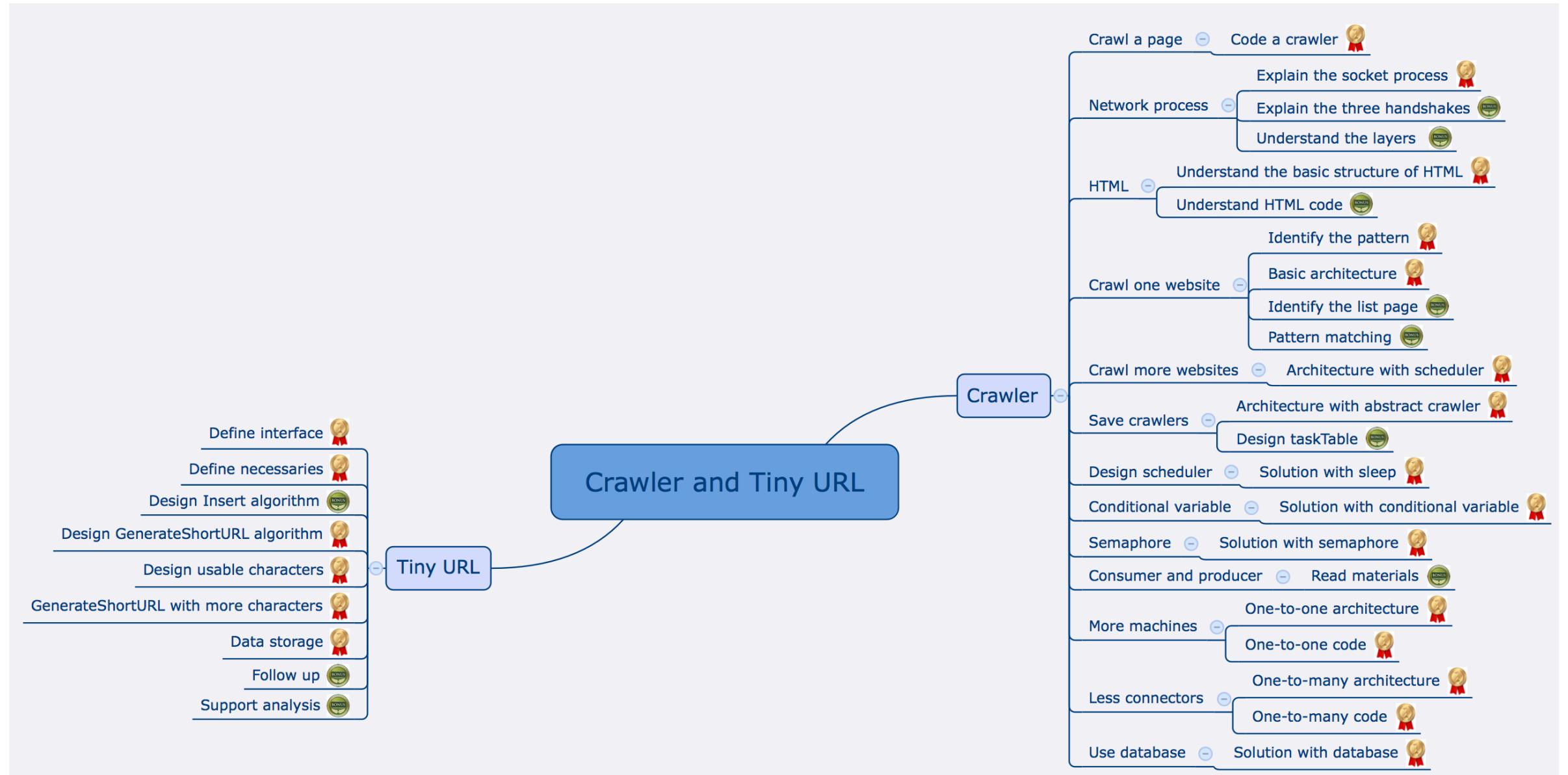
1

<http://www.jiuzhang.com/course/2/>

Read More

Novice, <http://www.xn--vi8hiv.ws/>

Class3 summary



Keyword

- Crawler
- HTTP
- Socket
- HTML
- TCP
- Thread
- Lock
- Conditional variable
- Semaphore
- URL
- Insert
- Lookup
- Encoding

QA



关注微信/微博，获取最新面试题及权威解答

微信: [ninechapter](#)

微博: <http://www.weibo.com/ninechapter>

官网: www.jiuzhang.com

One more thing

- How to ask?
 1. Answer the question by yourself
 2. Google solutions/unknown concepts
 3. Summarize your solution in details in QA
 4. Ask your question directly and clearly



能工摹其形，巧匠摄其魂