

## 31 | 规则引擎：自定义爬虫处理规则

2022-12-20 郑建勋 来自北京

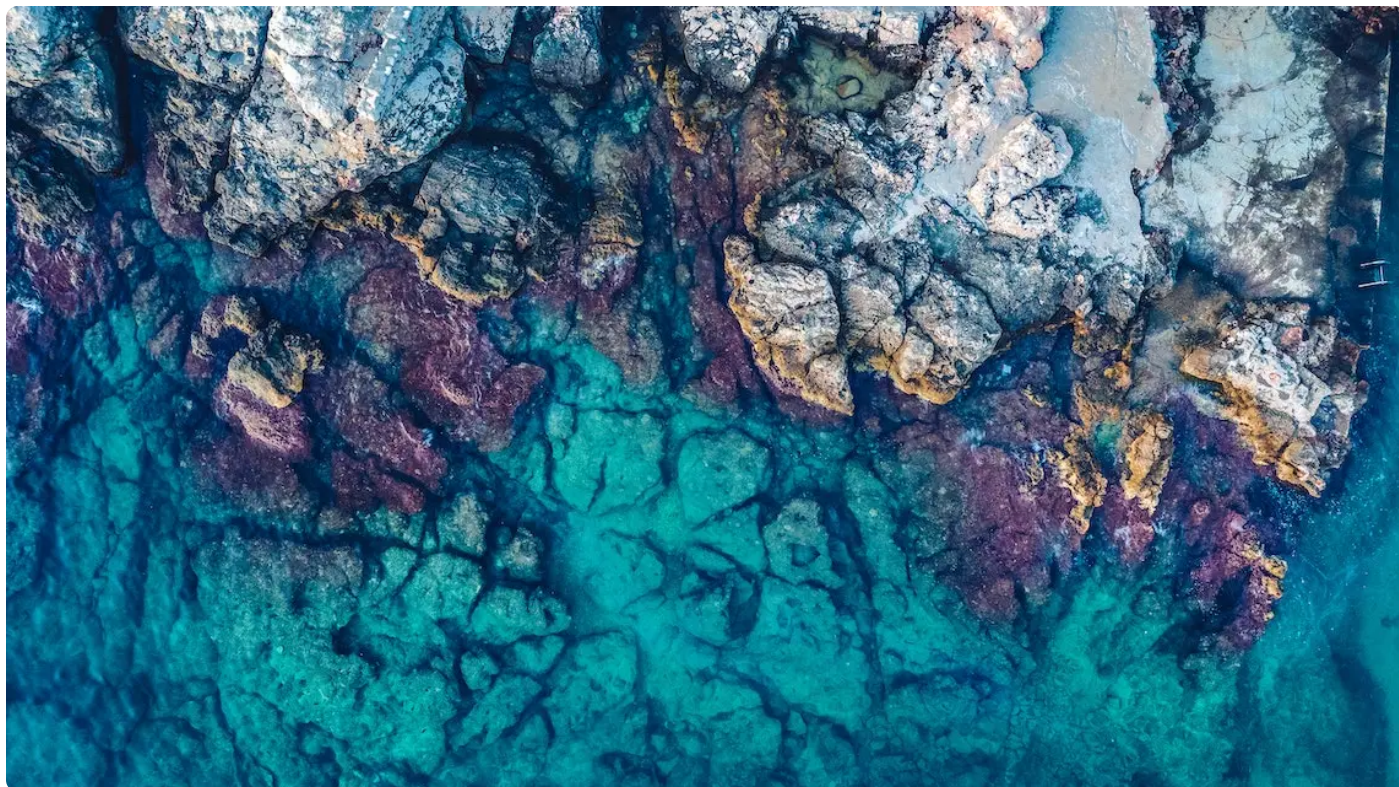


天下无鱼

<https://shikey.com/>

《Go进阶·分布式爬虫实战》

[课程介绍 >](#)



讲述：郑建勋

时长 10:28 大小 9.56M



你好，我是郑建勋。

这一节课，我们重点来看看如何更合理地设置爬虫任务规则。之前，我们在查找租房信息时，已经实现了有一定扩展性的程序。通过在每一个请求中加入 `ParseFunc` 函数，可以实现灵活的请求规则。

复制代码

```
1 // 单个请求
2 type Request struct {
3     unique    string
4     Task      *Task
5     Url       string
6     Method    string
7     Depth     int
8     Priority   int
9     ParseFunc func([]byte, *Request) ParseResult
10 }
```

但我们现在仍然面临几个问题：

- 一个爬虫任务会针对不同的网站有不同的处理规则，但现在的处理方式导致多个规则之间是割裂的，不便于统一管理。
- 我们在添加初始爬虫网站 URL 时，这些种子任务是在 main 函数中注入进去的，与任务的规则之间是割裂的。但是我们需要将初始爬虫 URL 与处理规则进行统一的管理。
- 当前的爬虫任务还是需要手动初始化才能运行，可配置化程度比较低。我们希望这些写好的静态任务在程序初始化时能够自动加载。而通过外部接口，或者只要在配置文件中指定一个任务名就能将任务调度起来。
- 更进一步，我们当前的任务和规则都是静态的，静态指的是代码需要提前写好，重新编译运行才能够在运行中被调用。我们能否动态地增加任务和任务的规则，让程序能够动态地解析我们的规则呢？

为了解决上面这些问题，我们需要有专门负责规则管理的模块，并完成静态规则与动态规则的解析。

## 静态规则引擎

静态规则处理的确定性强，它适合对性能要求高的爬虫任务。我们可以把一个任务的规则抽象如下：

 复制代码

```
1 type RuleTree struct {
2     Root func() []*Request // 根节点(执行入口)
3     Trunk map[string]*Rule // 规则哈希表
4 }
5
6 // 采集规则节点
7 type Rule struct {
8     ParseFunc func(*Context) ParseResult // 内容解析函数
9 }
10
11 type Context struct {
12     Body []byte
13     Req *Request
14 }
```

其实规则引擎就像一棵树。**RuleTree.Root** 是一个函数，用于生成爬虫的种子网站，而 **RuleTree.Trunk** 是一个规则哈希表，用于存储当前任务所有的规则，哈希表的 **Key** 为规则名，**Value** 为具体的规则。每一个规则就是一个 **ParseFunc** 解析函数。参数 **Context** 为自定义结构体，用于传递上下文信息，也就是当前的请求参数以及要解析的内容字节数组。后续还会添加请求中的临时数据等上下文数据。

下面我们沿用之前爬取租房信息的例子，将处理规则替换为使用新的静态规则引擎。

## 第一步，定义任务与规则。

我们在 **Task** 中加入了 **Name** 字段，将其作为一个任务唯一的标识。**Task** 里除了之前具有的最大深度、等待时间等属性，我们还加入了规则条件，规则条件中 **Root** 生成了初始化的爬虫任务。**Trunk** 为爬虫任务中的所有规则。

复制代码

```
1 type Task struct {
2     ...
3     Rule      RuleTree
4 }
5
6 var DoubangroupTask = &collect.Task{
7     Name:      "find_douban_sun_room",
8     WaitTime:  1 * time.Second,
9     MaxDepth:  5,
10    Cookie:    "xxx",
11    Rule: collect.RuleTree{
12        Root: func() []*collect.Request {
13            var roots []*collect.Request
14            for i := 0; i < 25; i += 25 {
15                str := fmt.Sprintf("<https://www.douban.com/group/szsh/discussion?start
16                roots = append(roots, &collect.Request{
17                    Priority: 1,
18                    Url:      str,
19                    Method:   "GET",
20                    RuleName: "解析网站URL",
21                })
22            }
23            return roots
24        },
25        Trunk: map[string]*collect.Rule{
26            "解析网站URL": &collect.Rule{ParseURL},
27            "解析阳台房":  &collect.Rule{GetSunRoom},
28        },
29    },
30 }
```

当前任务规则包括了“解析网站 URL”和“解析阳台房”这两个规则，分别对应了处理函数 ParseURL 和 GetSunRoom，如下所示。

 复制代码

```
1 const urlListRe = `(<https://www.douban.com/group/topic/[0-9a-z]+/>)"[^>]*>([^\<
2 const ContentRe = `<div class="topic-content">[\s\S]*?阳台[\s\S]*?<div class="as
3
4 func ParseURL(ctx *collect.Context) collect.ParseResult {
5     re := regexp.MustCompile(urlListRe)
6
7     matches := re.FindAllSubmatch(ctx.Body, -1)
8     result := collect.ParseResult{}
9
10    for _, m := range matches {
11        u := string(m[1])
12        result.Requesrts = append(
13            result.Requesrts, &collect.Request{
14                Method:    "GET",
15                Task:      ctx.Req.Task,
16                Url:       u,
17                Depth:    ctx.Req.Depth + 1,
18                RuleName: "解析阳台房",
19            })
20    }
21    return result
22 }
23
24 func GetSunRoom(ctx *collect.Context) collect.ParseResult {
25     re := regexp.MustCompile(ContentRe)
26
27     ok := re.Match(ctx.Body)
28     if !ok {
29         return collect.ParseResult{
30             Items: []interface{}{},
31         }
32     }
33     result := collect.ParseResult{
34         Items: []interface{}{ctx.Req.Url},
35     }
36     return result
37 }
```

第二步，初始化任务与规则。

在 `engine/schedule.go` 文件中，`init` 函数中的 `Store.Add` 函数将任务加载到全局的任务队列中。在 `Go` 中，`init` 函数是一个特殊的函数，它会在 `main` 函数之前自动执行。注意，当添加的任务越来越多之后，代码会变得臃肿，这不是一种优雅的写法。后面我们还会优化它。

天下无鱼  
<https://shikey.com/>

复制代码

```
1 // engine/schedule.go
2 func init() {
3     Store.Add(doubangroup.DoubangroupTask)
4 }
5
6 func (c *CrawlerStore) Add(task *collect.Task) {
7     c.hash[task.Name] = task
8     c.list = append(c.list, task)
9 }
10
11 // 全局爬虫任务实例
12 var Store = &CrawlerStore{
13     list: []*collect.Task{},
14     hash: map[string]*collect.Task{},
15 }
16
17 type CrawlerStore struct {
18     list []*collect.Task
19     hash map[string]*collect.Task
20 }
```

### 第三步，启动任务。

启动爬虫任务的方式可以分为两种，一种是加载配置文件，另一种是在调用用户接口时，传递任务名称和参数。不过在这里我们先用硬编码的形式来实现。而通过配置文件和用户接口来操作任务的方式我们会有专门的课程来实现。

复制代码

```
1 func main(){
2     ...
3     seeds = append(seeds, &collect.Task{
4         Name:      "find_douban_sun_room",
5         Fetcher: f,
6     })
7     s := engine.NewEngine(
8         engine.WithFetcher(f),
9         engine.WithLogger(logger),
10        engine.WithWorkCount(5),
11        engine.WithSeeds(seeds),
```



```
12     engine.WithScheduler(engine.NewScheduler()),
13     )
14     s.Run()
15 }
```



天下无鱼

<https://shikey.com/>

## 第四步，加载任务。

在调度器启动时，通过 `task.Rule.Root()` 获取初始化任务，并加入到任务队列中。

 复制代码

```
1 func (e *Crawler) Schedule() {
2     var reqs []*collect.Request
3     for _, seed := range e.Seeds {
4         task := Store.hash[seed.Name]
5         // 获取初始化任务
6         rootreqs := task.Rule.Root()
7         reqs = append(reqs, rootreqs...)
8     }
9     go e.scheduler.Schedule()
10    go e.scheduler.Push(reqs...)
11 }
```

在 Worker 处理请求时，需要从 `Rule.Trunk` 中获取当前请求的解析规则，并将内容和请求包装到 `Context` 中，调用 `ParseFunc` 对内容进行解析。

 复制代码

```
1 func (s *Crawler) CreateWork() {
2     for {
3         ...
4         //获取当前任务对应的规则
5         rule := req.Task.Rule.Trunk[req.RuleName]
6         // 内容解析
7         result := rule.ParseFunc(&collect.Context{
8             body,
9             req,
10        })
11        // 新的任务加入队列中
12        if len(result.Requesrts) > 0 {
13            go s.scheduler.Push(result.Requesrts...)
14        }
15    }
16 }
```

## 动态规则引擎

对 Go 语言这样的静态语言来说，我们需要在编译前就明确规则。这保证了程序的安全性与高性能，但是也失去了一些灵活性。像 Javascript、Python、Lua 这样的动态语言与 Go 有显著不同，它们不需要提前进行编译，能够比较灵活地书写并执行动态的规则。这一功能依赖于一种语言解释器，这种解释器一般是静态的语言编写的，例如 C/C++，解释器会解析这些动态语言的语法，然后执行相应规则。

和你分享一个实际企业中的例子。一家人工智能企业的核心产品之一是对视频流进行人脸识别。完成视频中人脸的解析涉及到视频的解码、人脸的识别、人脸的矫正、特征的提取等多个阶段。这个过程被称为一个 pipeline。pipeline 中的每一个阶段可能是串行的也可能是并行的。在过去，人脸、人群、物体的识别都需要单独来开发，这样开发的周期比较长，也缺乏灵活性。

为了应对未来灵活多变的检测需求，例如监测人是否摔倒，工人是否佩戴安全帽等，我们需要更短的开发周期，需要用更灵活的方式把这些阶段串联起来。这时这家公司就在 Go 中使用了 Lua 虚拟机。开发者遇到一个新的长尾需求时，通过书写 Lua 脚本来定义新的规则。在 Go 程序中通过动态加载 Lua 脚本来实现灵活性。

我们现在的爬虫项目其实也面邻着一样的问题。网站和规则是多种多样的，我们无法穷尽所有的规则，如果每次遇到新网站都要重新写代码，写爬取规则然后重启程序，这会比较繁琐，所以我们希望能够动态地在程序运行过程中加载规则。

动态规则带来的另一个好处是，降低了书写代码规则的门槛，它甚至可以让业务人员也能书写简单的规则。说到在爬虫项目中实现动态规则的引擎，我们首先想到的就是使用 **Javascript 虚拟机**了。因为使用 **JS 操作网页有天然的优势**。

自己要在短时间内实现一个工业级的虚拟机可能比较困难，我们可以使用一些开源的项目，例如 [🔗otto](#)。otto 是用 Go 编写的 JavaScript 虚拟机，用于在 Go 中执行 Javascript 语法。

下面是用 otto 编写的一个简单的例子。在这个例子中，script 字符串即为要执行的 Javascript 语法，console.log 是 JS 中的函数，用于打印变量。

```

1 package main
2
3 import (
4     "fmt"
5     "github.com/robertkrimen/otto"
6 )
7
8 func main() {
9     vm := otto.New()
10    script := `
11        var n = 100;
12        console.log("hello-" + n);
13        n;
14    `
15    value, _ := vm.Run(script)
16    fmt.Println("value:", value.String())
17 }

```



最终结果为：

```

1 hello-100
2 value: 100

```

这样，我们就实现了在 Go 语言中执行 JS 脚本的目的。实际上，otto 内部解析了这一串字符串，并按照 JS 语法中对应的规则进行了相应的处理，例如脚本中的 console.log 函数最终其实也调用了 Go 中的 fmt 函数，实现将文本打印到控制台。不过我们要注意的是，不一定 JS 的所有语法 JS 虚拟机都是支持的，是否支持取决于当前虚拟机的实现。例如当前 otto 支持 JS5 语法，但是不支持 JS6 语法。

下面我们仍然以爬取租房信息为例，借助 otto 库实现的 JS 虚拟机来实现动态的规则引擎。

## 第一步，构建动态规则模型 TaskModle。

```

1 type (
2     TaskModle struct {
3         Name      string    `json:"name"` // 任务名称，应保证唯一性
4         Url       string    `json:"url"`

```



```

5     Cookie    string    `json:"cookie"`
6     WaitTime  time.Duration `json:"wait_time"`
7     Reload    bool      `json:"reload"` // 网站是否可以重复爬取
8     MaxDepth  int64      `json:"max_depth"`
9     Root      string    `json:"root_script"`
10    Rules     []RuleModle `json:"rule"`
11 }
12 RuleModle struct {
13     Name        string `json:"name"`
14     ParseFunc    string `json:"parse_script"`
15 }
16 )

```



为什么这里要单独构建一个任务的结构体呢？主要原因在于，现在我们的规则都是字符串了，这和之前的静态规则引擎有本质的不同。其中，`TaskModle.Root` 为初始化种子节点的 JS 脚本，`TaskModle.Rules` 为具体爬虫任务的规则树。

## 第二步，书写动态爬虫规则。

示例代码如下。其中，`Root` 脚本就是我们要生成的种子网站 URL。在这里我们构建了一个 JS 数组 `arr`，将生成的请求数组添加到 `arr` 之后，又调用了 `AddJsReq` 函数。`AddJsReq` 函数其实是一个 Go 函数，用于最终生成 Go 中的请求数组。在这里我们可以看到，在 Go 的 JS 虚拟机中，还可以灵活地调用原生的 Go 函数。

复制代码

```

1 var DoubanggroupJSTask = &collect.TaskModle{
2     Property: collect.Property{
3         Name:        "js_find_douban_sun_room",
4         WaitTime: 1 * time.Second,
5         MaxDepth: 5,
6         Cookie:     "xxx",
7     },
8     Root: `
9         var arr = new Array();
10        for (var i = 25; i <= 25; i+=25) {
11            var obj = {
12                Url: "<https://www.douban.com/group/szsh/discussion?start=>" + i,
13                Priority: 1,
14                RuleName: "解析网站URL",
15                Method: "GET",
16            };
17            arr.push(obj);
18        };
19        console.log(arr[0].Url);
20        AddJsReq(arr);

```

```

21     `,
22     Rules: []collect.RuleModle{
23     {
24         Name: "解析网站URL",
25         ParseFunc: `
26         ctx.ParseJSReg("解析阳台房", "(<https://www.douban.com/group/topic/[0-9a-z]+
27         `,
28     },
29     {
30         Name: "解析阳台房",
31         ParseFunc: `
32         //console.log("parse output");
33         ctx.OutputJS("<div class=\"topic-content\">[\\\\s\\\\S]*?阳台[\\\\s\\\\S]
34         `,
35     },
36     },
37 }

```



而在 Rules 脚本中，我们加入了两个爬虫规则，分别是“解析网站 URL”和“解析阳台房”，他们都可以使用非常简单的规则来实现。在这里我们调用了 ctx.ParseJSReg 来解析请求，调用了 ctx.OutputJS 来解析并输出找到的内容，注意这里的 ctx.ParseJSReg 与 ctx.OutputJS 也是 Go 原生的函数，下面我们会看到他们的实现。

### 第三步，书写动态规则中的 Go 函数。

AddJsReqs 函数将在 JS 脚本中的请求数据变为 Go 结构中的数组 []\*collect.Request，而 ctx.ParseJSReg 方法则会动态解析 JS 中传递的正则表达式并生成新的请求，ctx.OutputJS 负责解析传递过来的正则表达式并完成结果的输出。注意 JS 虚拟机会自动将 JS 脚本中的参数转换为函数参数中对应的结构。

复制代码

```

1 // 用于动态规则添加请求。
2 func AddJsReqs(jreqs []map[string]interface{}) []*collect.Request {
3     reqs := make([]*collect.Request, 0)
4
5     for _, jreq := range jreqs {
6         req := &collect.Request{}
7         u, ok := jreq["Url"].(string)
8         if !ok {
9             return nil
10        }
11        req.Url = u
12        req.RuleName, _ = jreq["RuleName"].(string)
13        req.Method, _ = jreq["Method"].(string)
14        req.Priority, _ = jreq["Priority"].(int64)

```

```

15     reqs = append(reqs, req)
16 }
17 return reqs
18 }
19
20 // 动态解析JS中的正则表达式
21 func (c *Context) ParseJSReg(name string, reg string) ParseResult {
22     re := regexp.MustCompile(reg)
23
24     matches := re.FindAllSubmatch(c.Body, -1)
25     result := ParseResult{}
26
27     for _, m := range matches {
28         u := string(m[1])
29         result.Requesrts = append(
30             result.Requesrts, &Request{
31                 Method:    "GET",
32                 Task:        c.Req.Task,
33                 Url:         u,
34                 Depth:       c.Req.Depth + 1,
35                 RuleName: name,
36             })
37     }
38     return result
39 }
40
41 // 解析内容并输出结果
42 func (c *Context) OutputJS(reg string) ParseResult {
43     re := regexp.MustCompile(reg)
44     ok := re.Match(c.Body)
45     if !ok {
46         return ParseResult{
47             Items: []interface{}{},
48         }
49     }
50     result := ParseResult{
51         Items: []interface{}{c.Req.Url},
52     }
53     return result
54 }

```



天下无鱼

<https://shikey.com/>

## 第四步，初始化任务与规则。

初始化动态规则这一步更复杂一些，因为我们需要将 JS 脚本放入 `paesrFunc` 函数中，供 `otto` 库解析，代码如下所示。

 复制代码

```
1 func init() {
```

```

2     ...
3     Store.AddJSTask(doubangroup.DoubangroupJSTask)
4 }
5
6 func (c *CrawlerStore) AddJSTask(m *collect.TaskModle) {
7     task := &collect.Task{
8         Property: m.Property,
9     }
10
11     task.Rule.Root = func() ([]*collect.Request, error) {
12         vm := otto.New()
13         vm.Set("AddJsReq", AddJsReqs)
14         v, err := vm.Eval(m.Root)
15         e, err := v.Export()
16         return e.([]*collect.Request), nil
17     }
18
19     for _, r := range m.Rules {
20         paesrFunc := func(parse string) func(ctx *collect.Context) (collect.ParseRe
21             return func(ctx *collect.Context) (collect.ParseResult, error) {
22                 vm := otto.New()
23                 vm.Set("ctx", ctx)
24                 v, err := vm.Eval(parse)
25                 e, err := v.Export()
26                 return e.(collect.ParseResult), err
27             }
28         }(r.ParseFunc)
29         if task.Rule.Trunk == nil {
30             task.Rule.Trunk = make(map[string]*collect.Rule, 0)
31         }
32         task.Rule.Trunk[r.Name] = &collect.Rule{
33             paesrFunc,
34         }
35     }
36
37     c.hash[task.Name] = task
38     c.list = append(c.list, task)
39 }

```



天下无鱼

<https://shikey.com/>

在这里，用于生成种子网站的 Root 函数中的 vm.Eval(m.Root) 执行了配置中的 root 脚本，然后返回了生成的请求数组。vm.Set("AddJsReq", AddJsReqs) 可以将 Go 原生函数注册到 JS 虚拟机中，这样我们才能在 JS 脚本中调用 Go 函数。paesrFunc 函数也是一样，在这里我们使用了闭包，方便后续执行 parse 脚本并最后返回解析后的 collect.ParseResult 结果。

**第五步，启动并加载任务。**

启动任务和加载任务与静态规则引擎代码完全一致，完全复用就可以了。到这里，我们只需要指定一个任务名 `js_find_douban_sun_room` 就可以利用动态规则引擎将爬虫任务跑起来了。



复制代码

```
1 func main(){
2     ...
3     seeds = append(seeds, &collect.Task{
4         Property: collect.Property{
5             Name: "js_find_douban_sun_room",
6         },
7         Fetcher: f,
8     })
9
10    s := engine.NewEngine(
11        engine.WithFetcher(f),
12        engine.WithLogger(logger),
13        engine.WithWorkCount(5),
14        engine.WithSeeds(seeds),
15        engine.WithScheduler(engine.NewSchedule()),
16    )
17
18    s.Run()
19 }
```

完整的代码位于 [🔗v0.2.5](#)。

## 总结

在本节课程中，为了更好地对爬虫任务进行管理，我们构建了规则引擎模块来管理静态和动态的规则。

静态的规则指的是需要在 Go 代码中提前写好的解析规则，这些规则固定、性能更高。但是由于我们无法穷尽所有的网站和规则，所以每加入一个网站都需要修改代码，还要重新启动程序，过程比较繁琐。动态的规则构建了 Javascript 的虚拟机，它会动态解析 JS 的语法规则。

其实从我们现在的实现来看，用这些简单的规则就可以支持许多简单的爬取规则。在 JS 语法中，我们实现了调用 Go 函数的能力，这让书写 JS 脚本变得更加简单。我们甚至可以想象这样一个场景，用户只要在页面中选择自己希望爬取的内容，就足够我们生成动态的规则满足客户的爬取需求了。这种商业模式是成立的，用我们现在的技术完全可以实现。

不过，动态的规则当前的性能确实比不上静态的规则，因为它内部有大量运算、内存分配还有反射的使用。但是在爬虫项目中，很多时候真正的瓶颈是来自于网络 I/O。因此通常动态规则带来的灵活性收益要大于其性能的损耗。



## 课后题

你认为下面这个最基本的 JS 脚本，在 JS 虚拟机内部是如何实现的？

复制代码

```
1 script := `  
2   var n = 1+2;  
3   console.log("hello-" + n);  
4   n;
```

欢迎你在留言区与我交流讨论，我们下节课再见！

分享给需要的人，Ta购买本课程，你将得 20 元

生成海报并分享

赞 1 提建议

© 版权归极客邦科技所有，未经许可不得传播售卖。页面已增加防盗追踪，如有侵权极客邦将依法追究其法律责任。

上一篇 30 | 辅助任务管理：任务优先级、去重与失败处理

下一篇 32 | 存储引擎：数据清洗与存储

## 精选留言 (2)

写留言



Realm

2022-12-20 来自浙江

课程渐入佳境了，知识越来越有趣味了。老师，有时间可以把每节课的思考题答疑下，那就更好了。



思考题：

猜测是用AST抽象语法树+反射，解析js，转成Golang语法。



疑问：

...

```
type CrawlerStore struct {  
    list []*collect.Task  
    hash map[string]*collect.Task  
}  
...
```

这个list的设计是用于做什么？好像程序中没有提到。



翡翠虎

2022-12-20 来自广西

设计一个这样的引擎有什么好处？如果把xpath规则或者正则放到数据库，运行的时候随着任务传递给程序，程序按几个预设配置处理，会不会更好？

