

该文档为爬虫基于Redis调度方式的架构版本

若有任何疑问可以邮件我

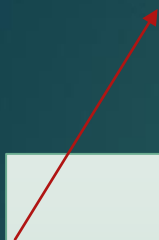
wengbenjue@163.com

- 分布式爬虫系统构架
- Zookeeper
- 爬虫模块逻辑
- 资源管理器
- 爬虫解析引擎
- 更新爬虫爬取周期
- 反被封
- 爬虫类图
- 日志记录模块
- 美团外卖爬取（二次开发方式）
- 任务层次批次
- Job管理
- Schema抽取
- 任务跟踪系统
- 监控系统

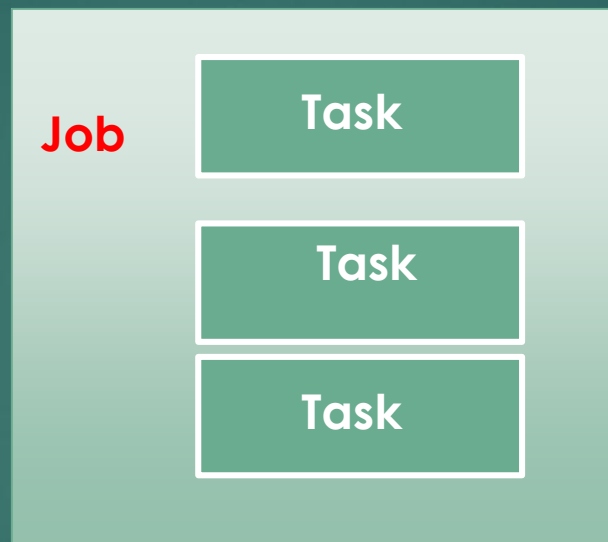
分布式爬虫系统构架

Job、Task

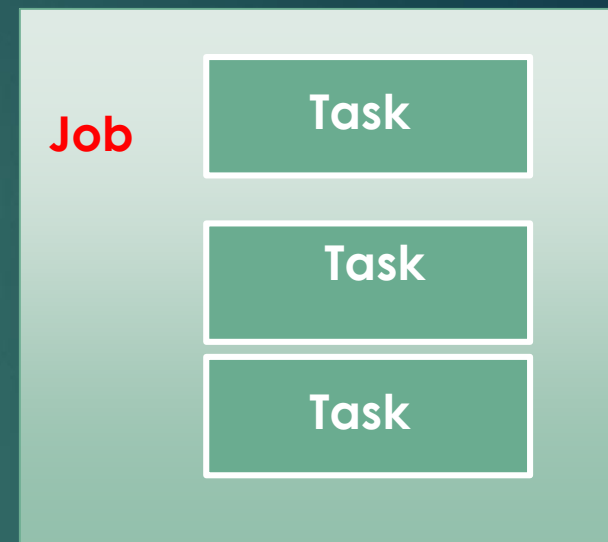
URL过滤规则接口



通用爬虫

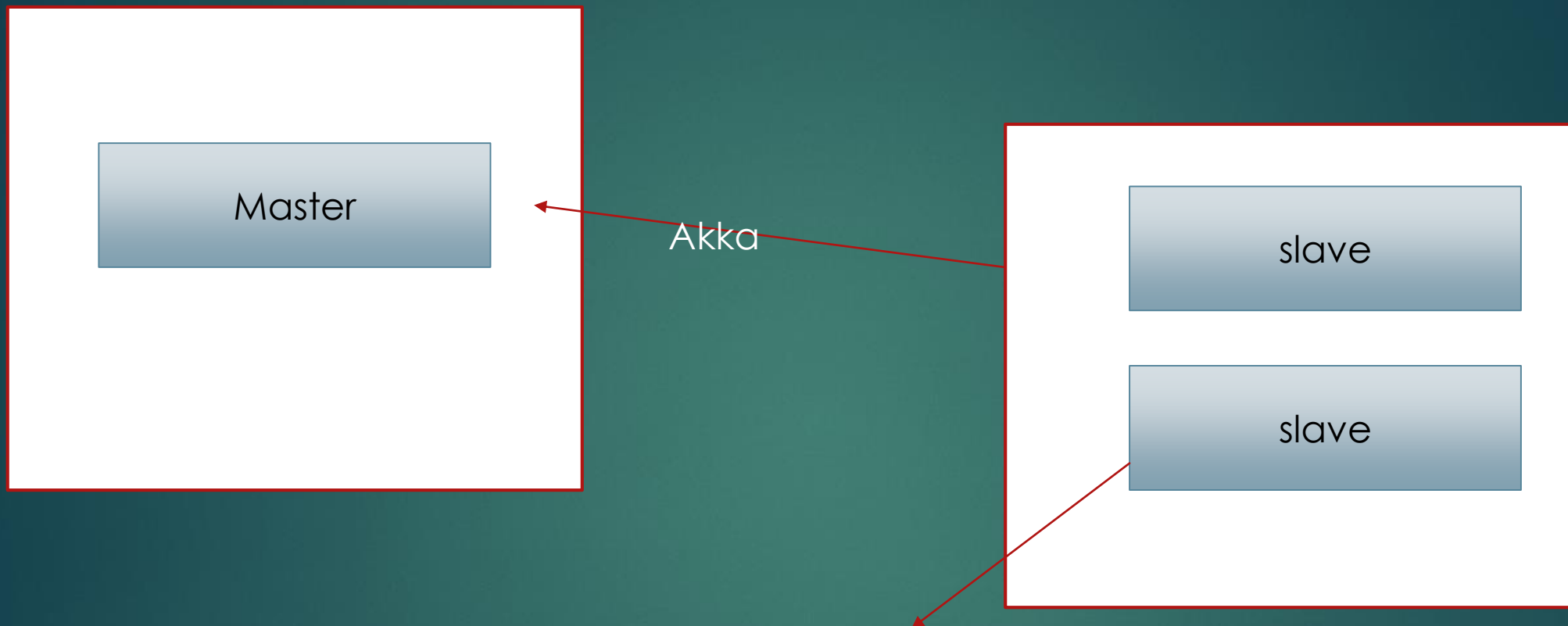


主题爬虫



指定页爬虫

规则过滤接口拓扑



对于每个job只需要第一次去Master取规则，然后自己做缓存，当Job退出的时候，首先清掉本机对应Jobid缓存，然后告诉master完成，清理掉缓存

种子入队列



需安装Redis

如果到时间点了将任务移到任务队列中，如果是主题爬虫不取出，将时间加上定期爬取时间

爬虫从队列领任务爬取

队列(Redis)



解耦

爬虫子爬虫架构

Pool

3s心跳去队列取数据并添加到事件队列，通过守护线程去触发事件，

事件队列

子爬虫
ActorRef队列

管理内存，选择子爬虫，将Url发送给子爬虫，子爬虫判断是否可以爬取，并返回状态码

事件队列

子爬虫ActorRef
队列

子爬虫

子爬虫

子爬虫

子爬虫

子爬虫

子爬虫

子爬虫日志记录及爬取成功后操作

子爬虫

子爬虫

子爬虫

子爬虫

子爬虫

子爬虫

将上次爬取时间，被封次数，爬取站域名次数，同一网页上次爬取时间等，写入文件或者Hbase等，方便通过统计分析，然后通过机器学习建模，Spark+Streaming或者Hive+Spark+Storm

将爬取后原始数据存入Hbase，因内网，提供服务接口

然后将存好后的key放入Kafka

需安装Kafka依赖于zookeeper,最好单独安装zookeeper

解耦

机器学习建模的两个应用目标

子爬虫用于判断是否能接受给予的url,进行爬取

优先级队列中，score的建模排序

解析前准备

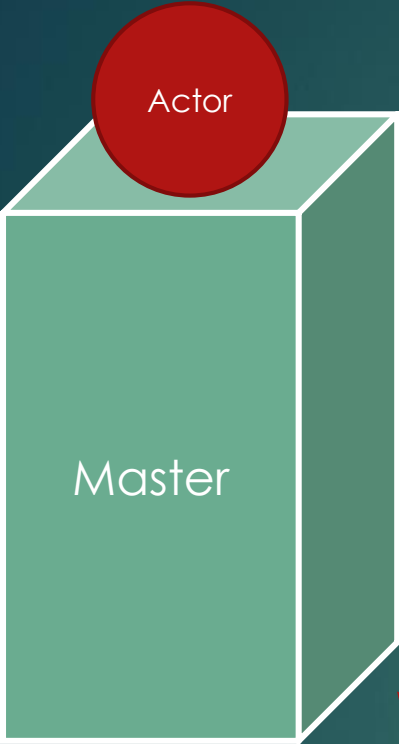
一、判断自己的资源，如内存、CPU是否够用。（通过主从分布式资源管理）

二、如果够用去Kafka拉取一条rowkey

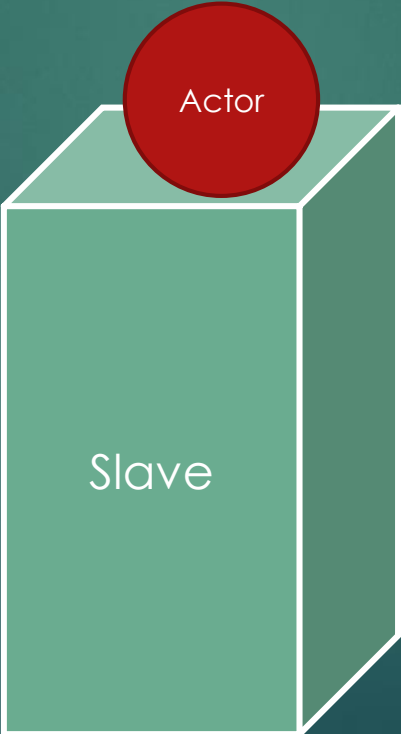
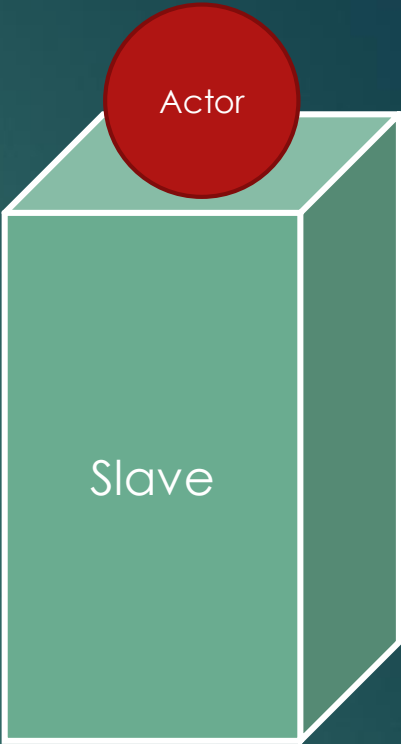
三、根据rowkey去hbase拉取该原始网页进行解析

解耦

资源管理及网页解析Topology



记录Memory和Cpu



Heartbeat

subscribe



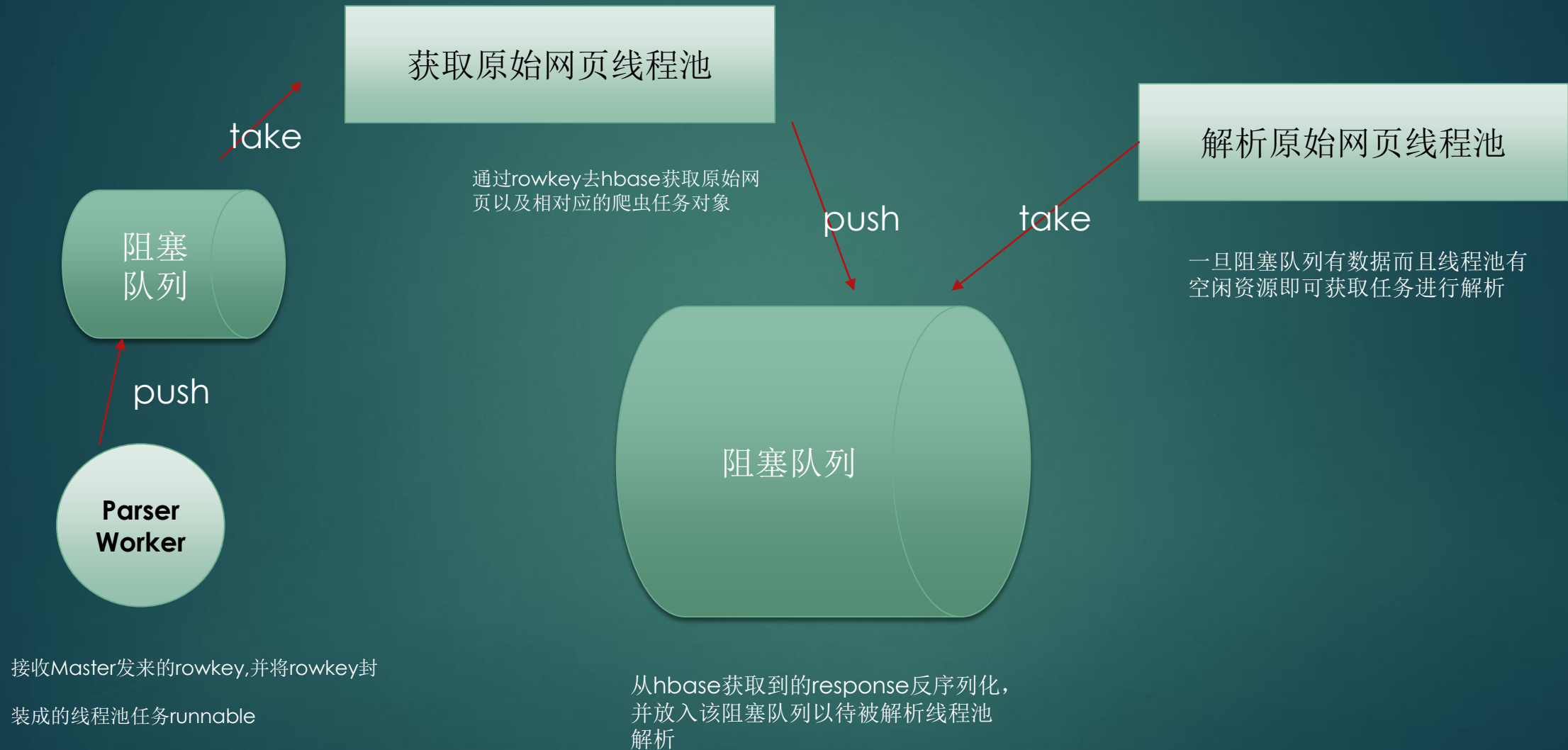
Kafka

getByRoykey



Hbase

解析线程池逻辑结构



解析网页

解析的url重新入队列

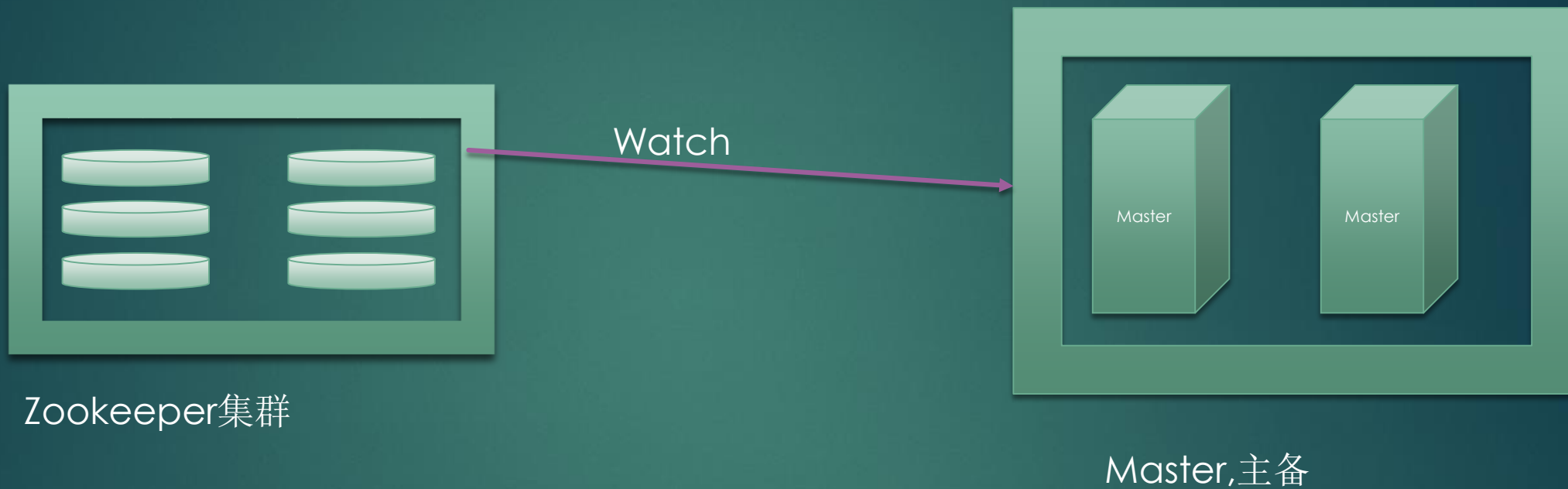
解析的网页内容进入MongoDB

解析可以用Xpath,
解析的架构和技术
需要花时间梳理下

解耦

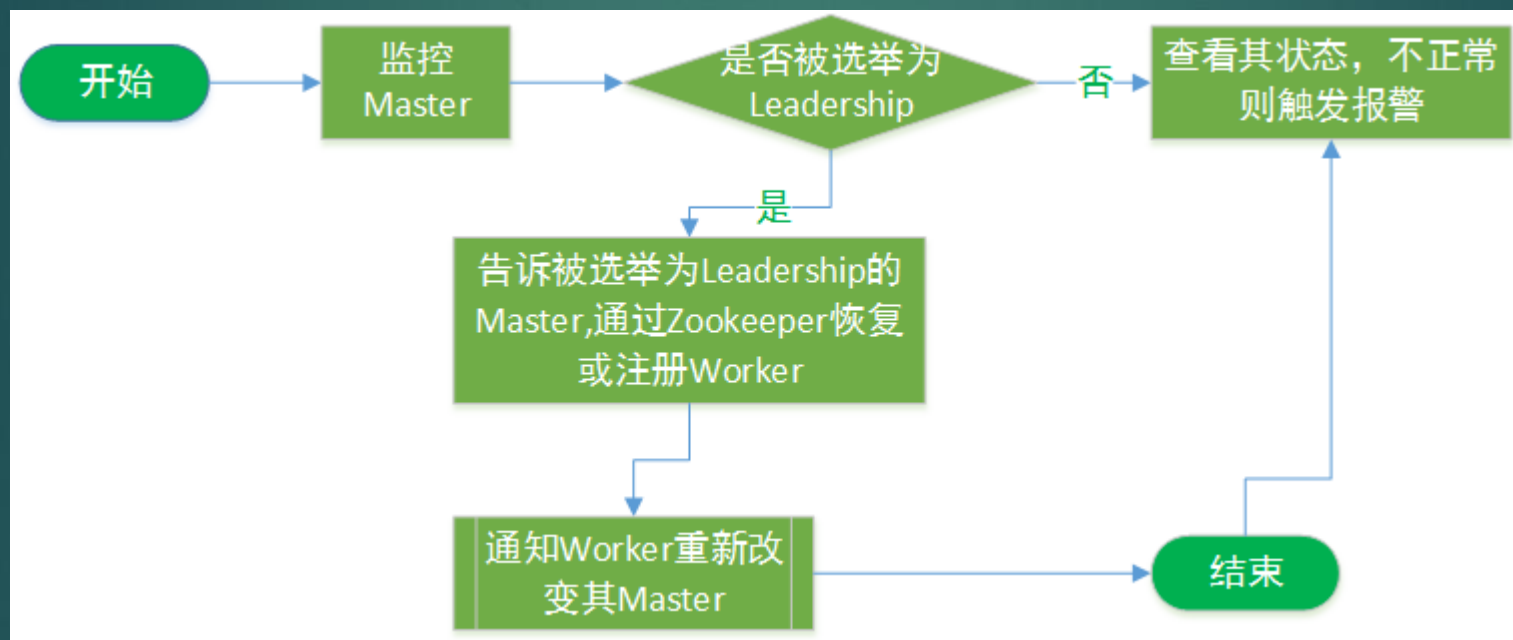
Zookeeper

Zookeeper Leadership election 图解



Zookeeper Leadership election 图解

领导选举流程图



Zookeeper

```
private[spark] class ZooKeeperLeaderElectionAgent(val masterActor: ActorRef,
    masterUrl: String, conf: SparkConf)
    extends LeaderElectionAgent with LeaderLatchListener with Logging {

    val WORKING_DIR = conf.get("spark.deploy.zookeeper.dir",
"/spark") + "/leader_election"

    private var zk: CuratorFramework = _
    private var leaderLatch: LeaderLatch = _
    private var status = LeadershipStatus.NOT_LEADER

    override def preStart() {

        logInfo("Starting ZooKeeper LeaderElection agent")
        zk = SparkCuratorUtil.newClient(conf)
        leaderLatch = new LeaderLatch(zk, WORKING_DIR)
        leaderLatch.addListener(this)

        leaderLatch.start()
    }

    override def preRestart(reason: scala.Throwable, message: scala.Option[scala.Any]) {
        logError("LeaderElectionAgent failed...", reason)
        super.preRestart(reason, message)
    }

    override def postStop() {
        leaderLatch.close()
        zk.close()
    }
}
```

Zookeeper

```
/**
 * A LeaderElectionAgent keeps track of whether the current Master is the Leader, meaning it
 * is the only Master serving requests.
 * In addition to the API provided, the LeaderElectionAgent will use of the following messages
 * to inform the Master of leader changes:
 * [[org.apache.spark.deploy.master.MasterMessages.ElectedLeader ElectedLeader]]
 * [[org.apache.spark.deploy.master.MasterMessages.RevokedLeadership RevokedLeadership]]
 */
private[spark] trait LeaderElectionAgent extends Actor {
  // TODO: LeaderElectionAgent does not necessary to be an Actor anymore, need refactoring.
  val masterActor: ActorRef
}

/** Single-node implementation of LeaderElectionAgent -- we're initially and always the Leader. */
private[spark] class MonarchyLeaderAgent(val masterActor: ActorRef) extends LeaderElectionAgent {
  override def preStart() {
    masterActor ! ElectedLeader
  }

  override def receive = {
    case _ =>
  }
}
```

爬虫模块逻辑

爬虫入口逻辑

入口可以通过Web或者命令行进入

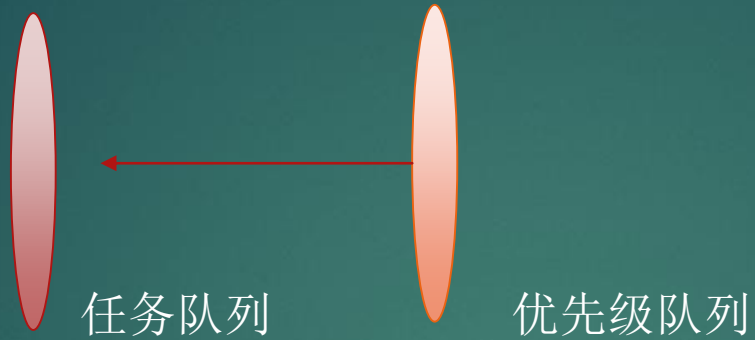
用户可以添加任务，并给予url过滤规则(比如正则过滤)，系统自动封装job，并生成JobId,关联任务及过滤规则

根据jobid将过滤对象存入job所驻留的服务器，解析的时候通过Akka取得，然后缓存到本地，当job完成或删除时候，清楚主、从所有缓存

提交Job

Job转换为Task提交task到队列，默认提供Redis队列

爬虫任务队列转移逻辑



优先级队列，按执行时间排序，时间早的先执行
判断爬虫类型，对于主题爬虫，判断爬取时间与现在时间的时差，如果过了或刚好为当前时间就取出来放入任务队列，但是不需要删除，只需要改变时间为下次执行时间，方便下次继续执行；【使用双扫线程】

爬虫到子爬虫逻辑



- 1.爬虫从redis取得任务并放入本地Blocking Queue中[使用双扫线程]
- 2.使用死循环从blocking queue中取得任务并通过子爬虫传过来的资源状况（如内存，cpu，以及等待任务数，空闲线程数等）进行判断调度
- 3.将任务进行序列化并通过Akka将任务传给对应的子爬虫

另外

- 1.子爬虫通过心跳将自身资源传给爬虫，同时爬虫通过这个来判断子爬虫是否活着
- 2.爬虫集群的单点问题通过zookeeper
- 3.资源管理抽象出来

子爬虫逻辑



子爬虫
【Akka Actor】

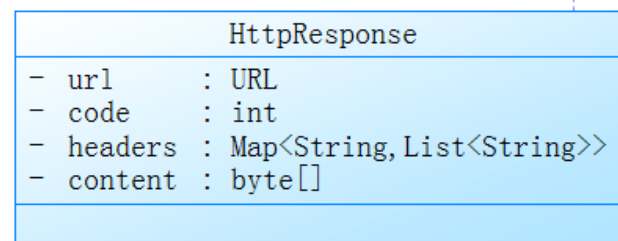
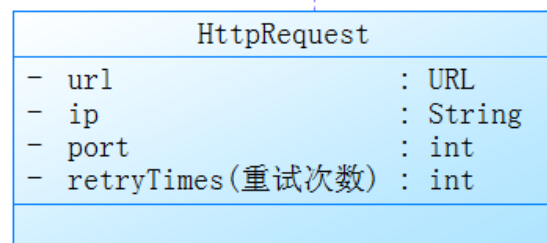
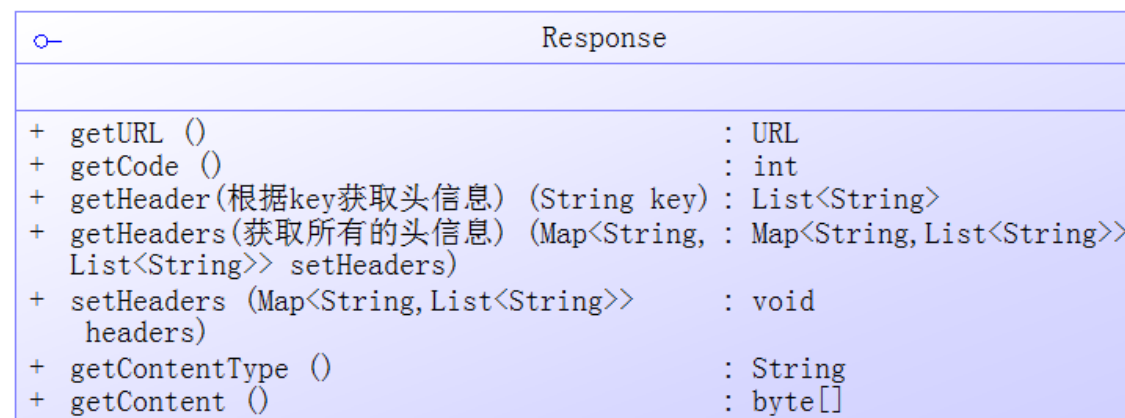
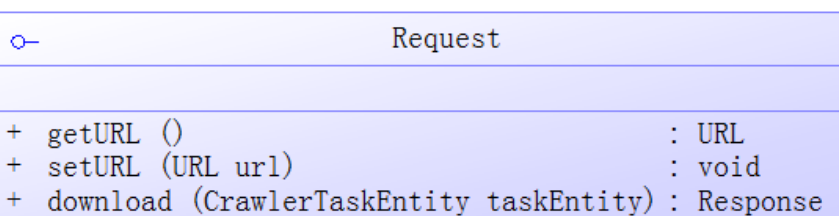
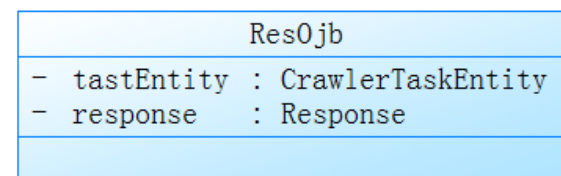
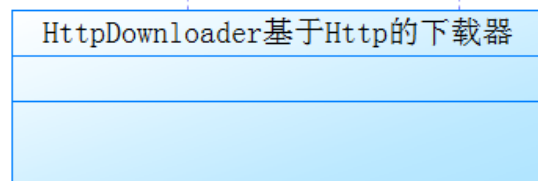
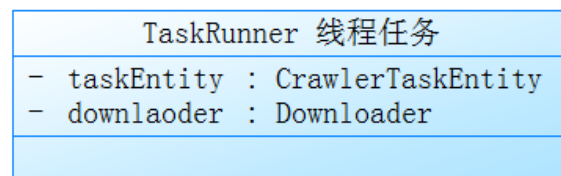
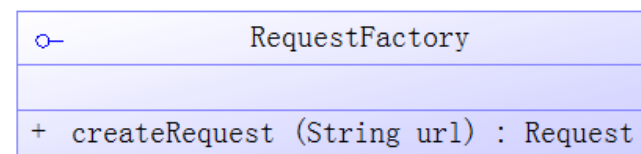
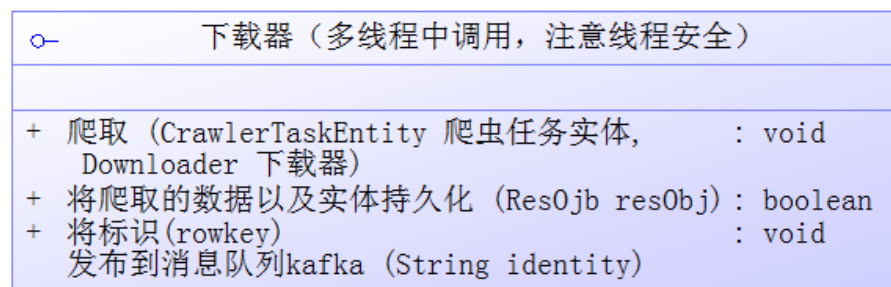


Executor

1.接收到爬虫传过来的uri后将uri进行反序列化操作，然后封装成TaskRunner传给线程池（注意线程安全问题）
线程池设置为守护线程，一直执行

2.线程池内部，即TaskRunner将uri交给下载器，下载器的初始化在TaskRunner中进行

3.下载器实现类（比如Http下载器使用HttpRequest下载网页，并将返回信息封装为HttpResponse）下载完后将实体类和Response封装为ResObj序列化并放入Hbase中，同时将rowkey放入kafka中

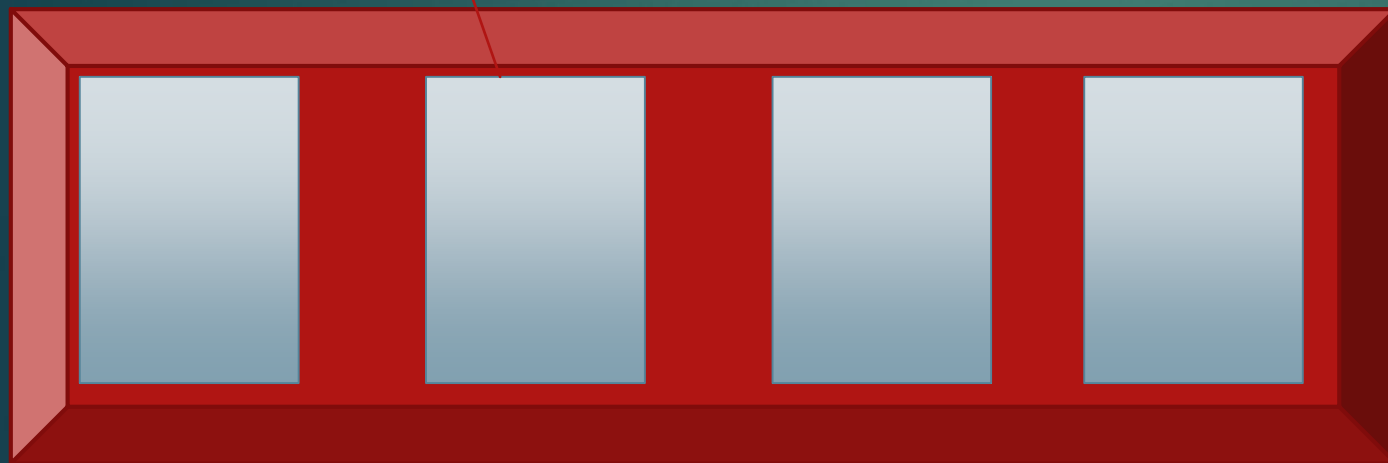


资源管理器

资源管理



优先级队列



Workers

```
case class Res(workerId: String, memUsage: Double, cpuUsage: Double, cpuCores: Int, totalMem: Double) extends TransferMsg
```

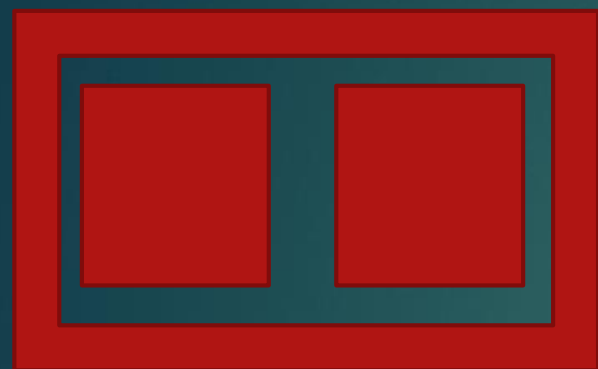
Worker定时发送本机资源信息到master，master通过优先级队列进行资源重排，若要求高性能，可采用事件机制



从服务器定时发送心跳报告资源情况
主服务器从优先级取得并删除对头元素
同时主服务器将从服务器缓存到普通队列中，普通队列远远小于优先级队列
当优先级队列为空的时候从普通队列取出

爬虫解析引擎

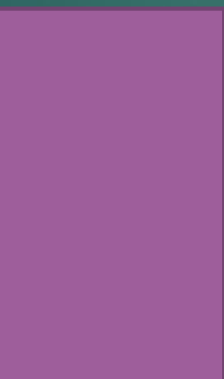
解析逻辑



解析
Master
Akka



解析
Slave
Akka



Job服务器
Pull

注意

ParserManager首先判断是否为第一层，如果是第一层则将爬虫实体更改下次爬取时间后重新放回爬虫队列

Master解析逻辑

- 1.Master通过Akka收到子爬虫发过来的rowkey
- 2.调用资源管理进行对比判断Slave的资源情况并将数据发给对应的Slave

另外

- 1.Master通过主备解决单点问题
- 2.Slave会心跳给Master（统一封装到资源管理中），更新资源情况并以此判断Slave是否挂掉

Slave解析逻辑

- 1.启动线程池
- 2.将rowkey传给ParserRunner，并进入多线程
- 3.调用抽象类ParserManager，解析对象在ParserRunner中初始化
4. ParserManager首先通过rowkey委托给**存储管理器(StorageManager)**去hbase获取数据，并进行反序列化
5. ParserManager将从存储管理器取得的数据委托给解析者Parser的工厂ParserFactory（工厂模式根据爬虫类型创建Parser），工厂创建Parser调用parse进行解析。
- 6.如果是**通用爬虫**，ParserManager解析者**Parser**调用通用爬虫引擎General Parser，如果是**API爬虫**调用APIParser，如果是**主题爬虫**TopicParser，解析方法内部根据调用主题爬虫解析者**域名**来通过反射（使用控制反转，缓存对象）调用具体的解析对象（这里使用抽象类，url抽取使用的是相同的）的解析方法，也就是解析规则可以根据域名进行二次开发，形成通用库
- 7.解析方法调用url分析器URLAnalyzer，根据url生成Jsoup Document
- 8.使用Jsoup抽取url和内容，委托给内容解析器ContentParser进行分析，提取所需的内容根据数据结构封装好返回
- 9.抽取的url委托给url过滤器ITaskFilter过滤url（url过滤规则通过JobId去Job服务器取得）
- 10.过滤后的url以及实体交给url组装器URLAssembler根据参数组装url
- 11.调用URL去重管理器URLDistinctManager（如果是重复的url不调用调度器入库Redis）
- 12.然后剩下的不重复的调用调度器，将url重新交给redis队列
13. 将网页内容委托给存储管理器持久化到MongoDB

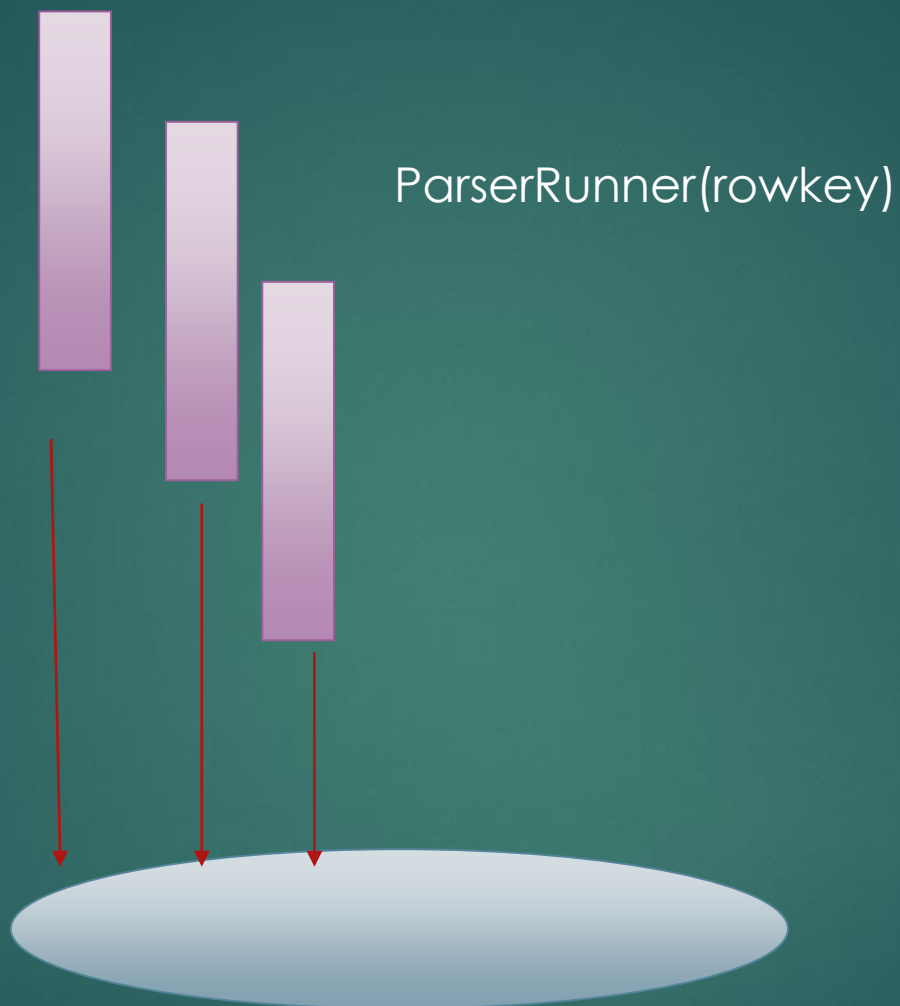
ParserMaster为任务id分配资源

- 1.通过阻塞队列去kafka获得rowkey
- 2.然后通过资源管理器判断是否有充足的资源
- 3.如果有就分配，如果都没有则休眠5s然后重新将任务放入阻塞队列

需要工作线程，一个双扫监控线程



解析
Slave
Akka



ParserRunner(rowkey)

- 1.启动线程池
- 2.将rowkey传给ParserRunner，并进入多线程

线程池数量由woker机器的CPU核数决定，也可以自行配置

守护线程池，执行
ParserRunner



- 3.调用抽象类ParserManager，解析对象在ParserRunner中初始化
4. ParseManager首先通过rowkey委托给**存储管理器(StorageManager)**去hbase获取数据，并进行反序列化

存储管理器

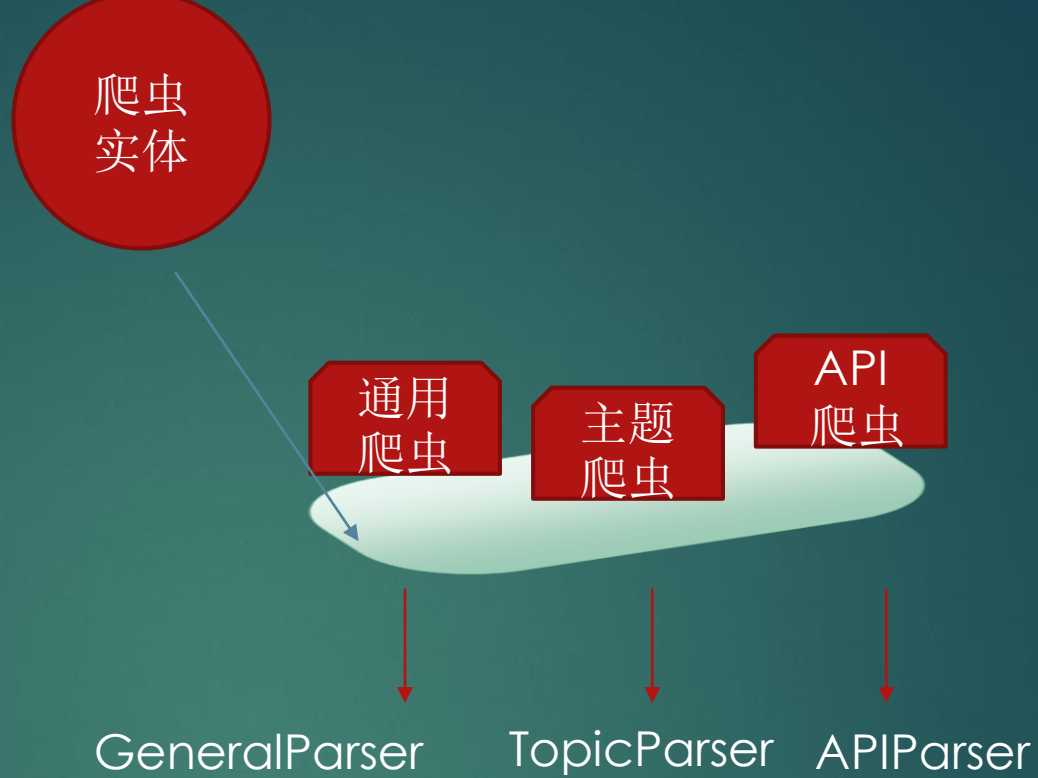


getByKey(key:String): T

getByKey(key:String): Seq[T]

putByKey(key:String,value:T): Boolean

put (entity:T): Boolean

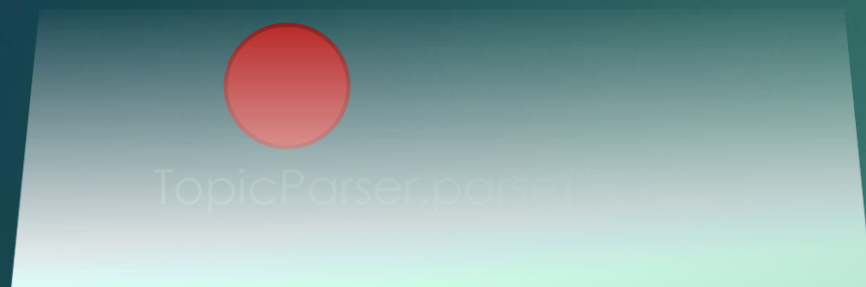


5. ParserManager将从存储管理器取得的数据委托给解析者Parser的工厂ParserFactory (crawlerTaskEntity) (工厂模式根据爬虫类型创建Parser,实际系统中根据域名进行创建 [注意从缓存中取, 保证有且仅有一个对象, 因为在多线程环境下, 需要加锁]),工厂创建Parser调用parse进行解析.

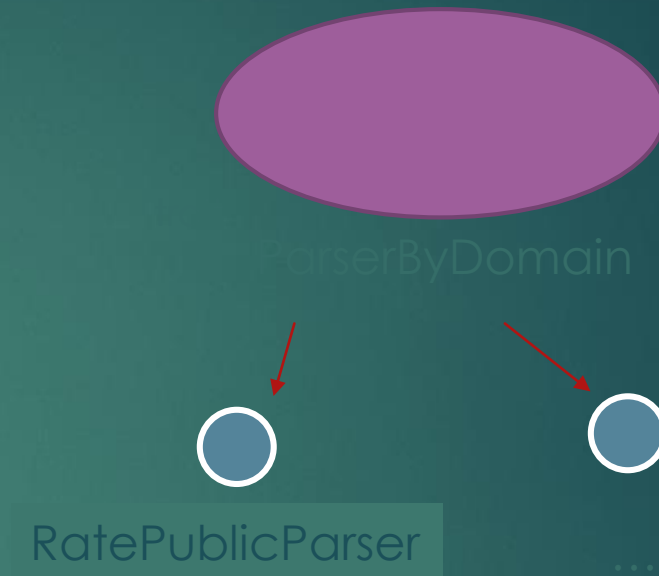
根据域名创建对象通过从映射内存表中获取类全名, 如果没有, 爬虫类型为主题爬虫时候, 抛出异常, 反之使用通用解析引擎

主题爬虫根据域名选择解析实体

- 1.通过解析类工厂创建传入参数域名
- 2.通过域名从本地缓存取对应的实体HashMap（domain->主题爬虫实体）
- 3.如果没有取到实体则根据爬虫实体中的属性，主题爬虫类路径反射创建实体并放入 本地缓存（注意捕获异常）

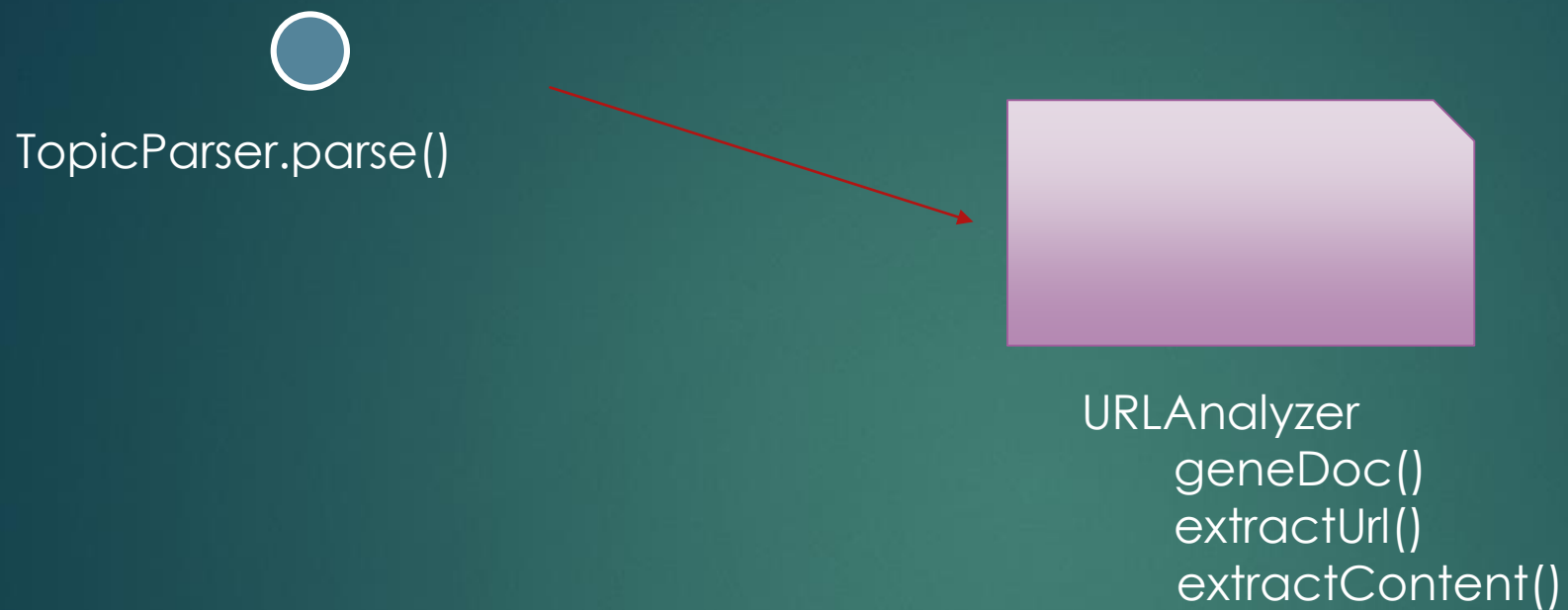


工厂模式
根据域名产生具体解析对象



~~Version0.1~~

6.如果是通用爬虫，ParserManager调用通用解析器General Parser，如果是API爬虫调用APIParser,如果是主题爬虫TopicParser,解析方法内部根据调用主题爬虫解析者域名来通过反射（可以使用解析对象池，使用队列）调用具体的解析对象ParserByDomain.parse（这里使用抽象类，url抽取使用的是相同的）的解析方法，也就是解析规则可以根据域名进行二次开发，形成通用库



- 7.解析方法调用url分析器URLAnalyzer，根据url生成Jsoup Document
- 8.使用Jsoup抽取ur和内容，委托给内容解析器ContentParser进行分析，提取所需的内容根据数据结构封装好返回



TopicParser.parse()

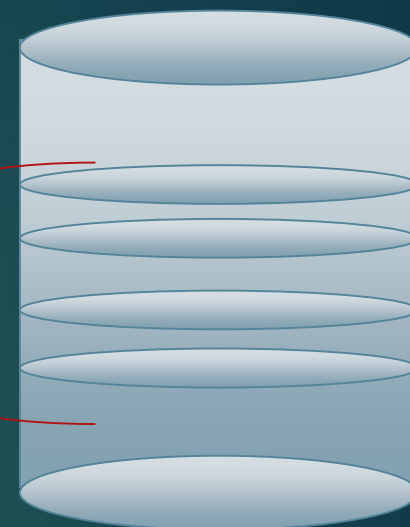
TopicParser.geneNewTaskEntitySet()

比如说分页，没有返回null,产生新任务，但是批次不变
还有，如果深度是第一级别的需要重新放回队列



ITaskFilter.matchRule

rules



Job服务器

parse方法内部逻辑

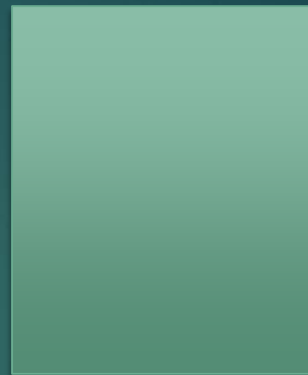
taskFilter通过工厂创建，默认从本地缓存取
Map(jobid->TaskFilter) TaskFilter是pipeleline形式，
使用装饰模式，一个job对应一个嵌套TaskFilter，只
需要一个触发条件

9.抽取的url委托给url过滤器ITaskFilter过滤url（url过滤规则通过Jobld去Job服务器取得）



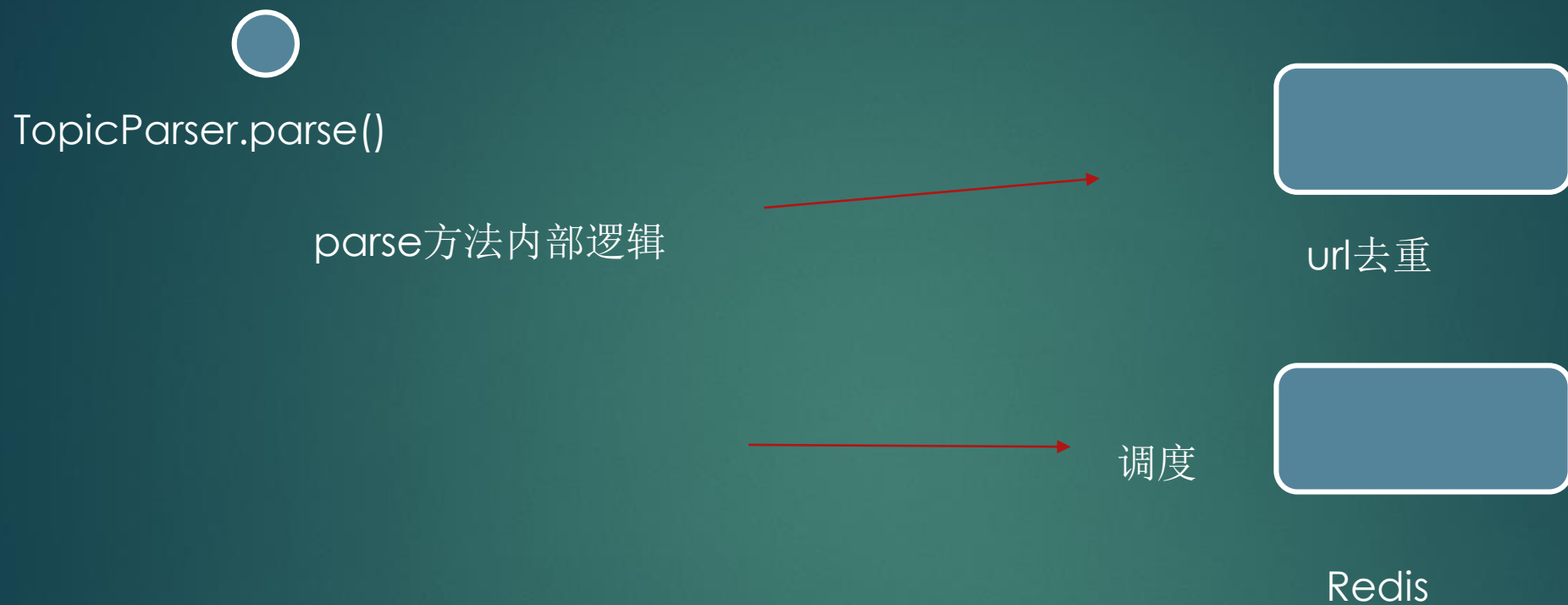
TopicParser.parse()

parse方法内部逻辑



参数组
装器

10.过滤后的url以及实体交给url组装器URLAssembler器根据参数组装url



- 11.调用URL去重管理器URLDistinctManager（如果是重复的url不调用调度器入库Redis）
- 12.然后剩下的不重复的调用调度器，将url重新交给redis队列

url解析与内容解析（入库）解耦

Parse(resobj){

1.生成parser对象

2.自动判断任务终止条件

3.抽取url (CralwerTaskEntity)

4.异步抽取内容并入库

5.对抽取的urls(CralwerTaskEntity)进行过滤，去重...

6.将新的url(CralwerTaskEntity)重新放入队列，等待调度

}

使用阻塞队列BlockingQueue,插入不阻塞，读取进行阻塞，读取的时候使用双扫线程，将读入的数据放入另外一个解析线程池进行解析

解耦使用引用传递，总共只有一个解析对象以及一个document

解析url并重新调度线程池

解析内容并封装成对象，
交给存储管理器（默认使用mongo）

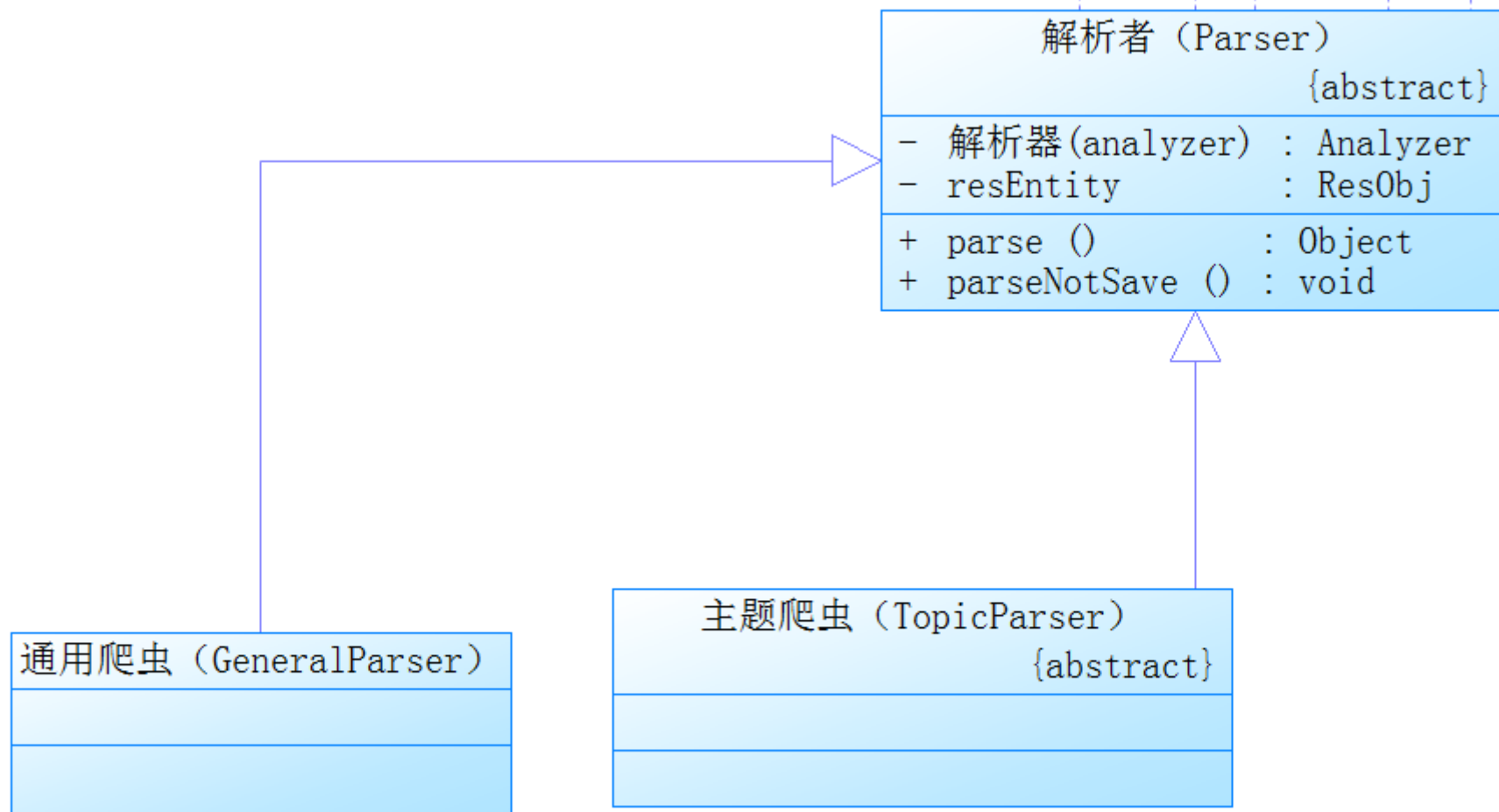
解析管理者 (ParserManager)

{abstract}

- storageManager : StorageManager
 - parser : Parser
 - 解析器(analyzer) : Analyzer
 - Url过滤器(taskFilter) : ITaskFilter
 - url组装器(urlAssembler) : URLAssembler
 - url去重(urlDistinctManager) : URLDistinctManager
 - 调度器(ITaskQueue) : ITaskQueue
-
- + 启动(String taskKey) : void
 - + 获取爬虫待分析 () : ResObj
 - + 抽取url连接(ResObj 待分析实体) : List<Link>
 - + 抽取href Urls(ResObj 待分析实体) : List<String>
 - + url过滤(ResObj 待分析实体, List<String> 抽取的url链接) : List<String>
 - + url去重(根据类型, 通用爬虫需要去重) (ResObj 待分析实体, List<String> 抽取的url链接) : List<String>
 - + 组装url(ResObj 待分析实体, List<String> 抽取的url链接) : List<CrawlerTaskEnti
 - + 调度, 重新放入队列(CrawlerTaskEntity 爬虫实体对象) : boolean
 - + 保存实体 () : boolean

爬虫解析线程 (ParserRunner)

- rowkey : String
- parserManager : ParserManager



更新爬虫爬取周期

反被封

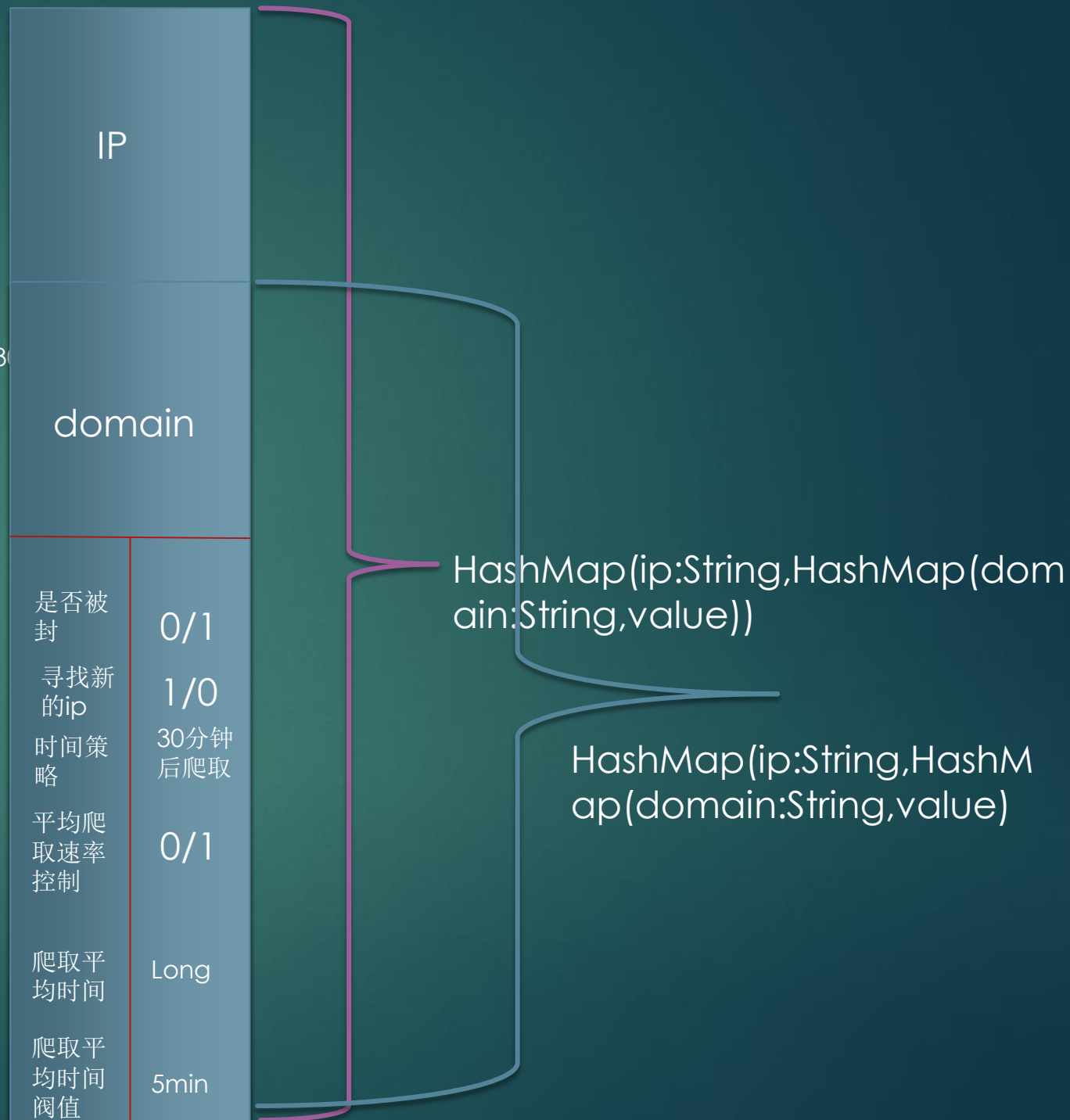
爬取什么

什么时候爬取

能不能爬取，什么时候能爬取

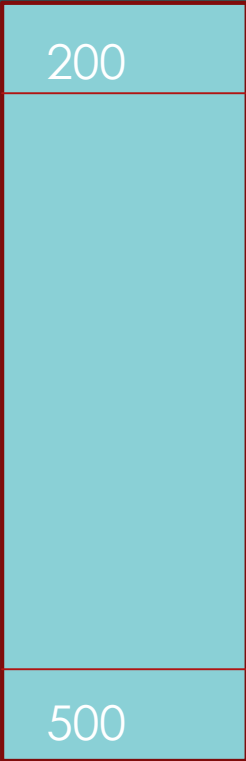
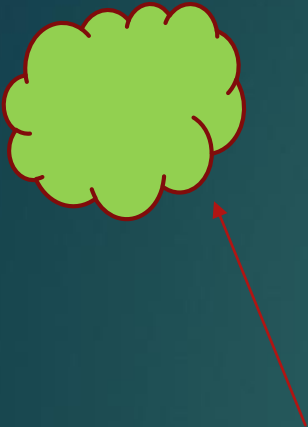
@1.爬取IP被远程域给封了 -》 1) 重新选取IP爬取 2) 3
@2.爬取IP爬取远程域很慢（由于地域分布问题）

怎么爬取



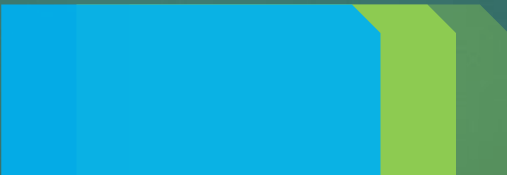
机器学习训练网页是否被封

(类似垃圾邮件过滤)



2.离线抽取网页源码，人工进行训练集标注，然后通过机器学习（可以使用native bayes、SVM或者logistic regression）

子爬虫ip	请求域名	Response 状态码	被封状态	网页源码内容
-------	------	-----------------	------	--------



日志记录器

1.爬虫请求页面，如果是404则直接丢掉，如果是在300到500中间，则将网页内容记录下来，如果是302（开启默认不让浏览器跳转），则重新手动请求Location并记录下跳转后的页面源码

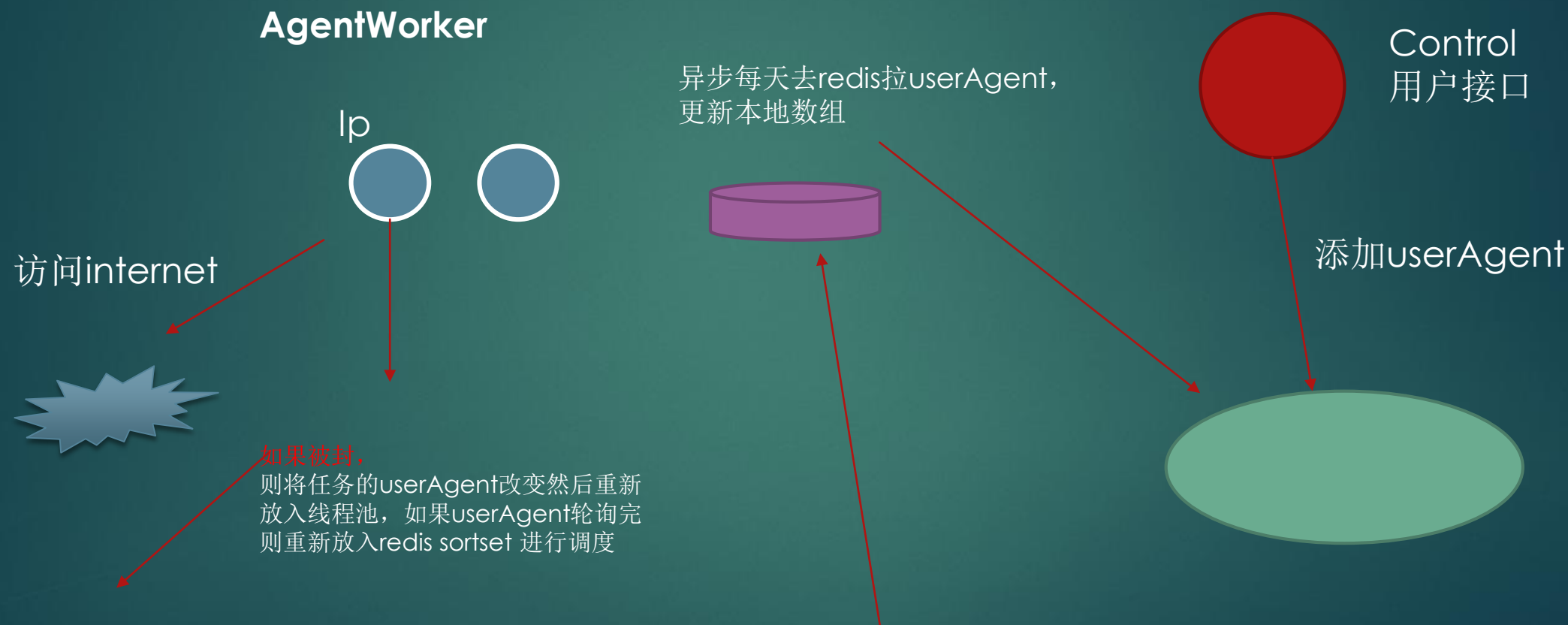
3.每次请求的时候如果是200到500之间，则得到源码并分词抽取向量输入模型进行判断

4.如果是被封状态则将redis的状态改为被封，方便分配资源的时候进行ip轮训

一般网站封ip的维度为

domain	IP	userAgent
--------	----	-----------

关于浏览器的动态分配



怎么判断是否被封, 通过机器学习

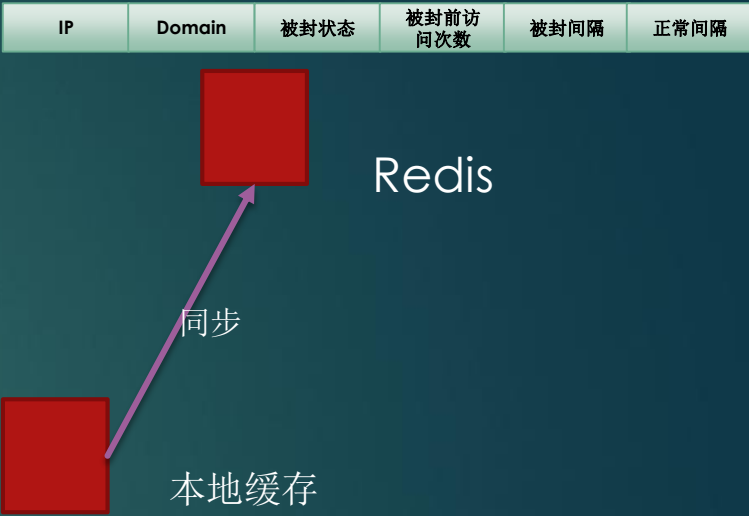
domain	当前浏览器	当前使用的是第几个浏览器	指针
--------	-------	--------------	----

被封指标----单位时间内访问域的次数 次数/单位时间

目标，通过一台机器上因为被封学习到的经验分享到其他机器，以免其他机器被封掉

- 1.怎么才能算被封
- 2.被封后该怎么处理
- 3.怎么预测被封（一台机器学习到的规则，资源共享）

一台机器学习到的经验，全局共享



AgentWorker

domain	(被封前请求开始时间，被封前请求结束时间，访问次数，开始被封时间，结束被封时间，状态（被封或者没有被封）)
--------	---

可以简化为下面，被封前请求开始时间就等于被封结束时间，同理，被封前请求结束时间就等于开始被封时间

domain	(访问次数，开始被封时间，结束被封时间，状态（被封或者没有被封）)
--------	-----------------------------------

在callback里面更新，当ip被封的时候，将本地缓存放入redis,放之前判断（请求结束时间-请求开始时间）是不是小于redis中已经存在的时间，如果是则更新，反之不更新；然后将本地缓存的请求开始时间和请求结束时间归零，记录开始被封时间，同时将被封状态改为被封；每次请求200的时候判断是否被封状态为被封如果是则更新结束被封时间，同时更新被封状态为解封

Callback中，每次请求完进行判断，同一域名，时间是不是接近于redis中的时间，如果是则取出请求次数，判断本地请求次数是否接近于redis中记录的，基于性能考虑可以每天去redis中获取该记录同步到本地缓存，如果请求次数接近Redis则将任务重新放入redis sortset，并设置好相应的Interval阈值等待调度，可以设置为3分钟

在资源调度的基础上，基于反封策略进行调度

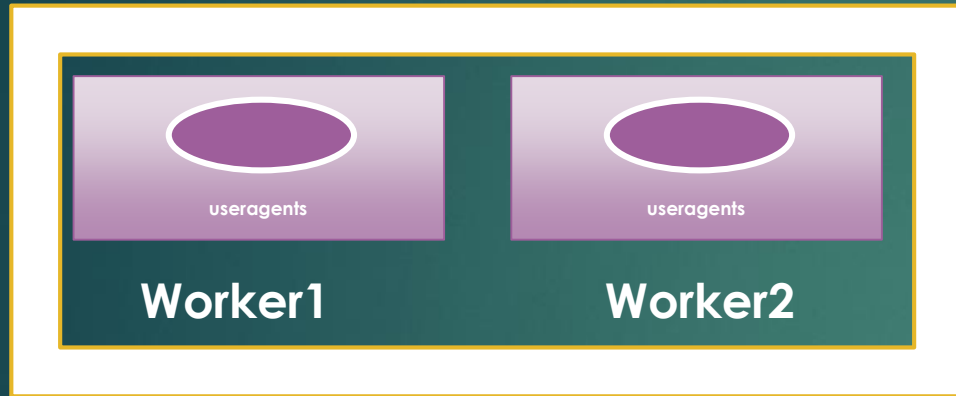
1.ip轮询

2.时间间隔策略(如果所有ip轮询完后都处于被封状态则启用该策略)

domain	(请求开始时间，请求结束时间，访问次数，开始被封时间，结束被封时间，状态（被封或者没有被封）)
--------	---

Redis里面取得该对象，得到(被封结束时间-被封开始时间)，为了性能考虑，可以定时异步从**redis**获取与本地缓存同步

代理Proxy模块



被封

如果所有Worker上的物理IP都被封了，没加入代理Proxy之前是将任务放进sortset等待调度，等待时间是被封的间隔时间；

加入代理Proxy之后则不是直接等待调度，而是首先启用代理，代理的有效期只在物理ip被封的时间段内，如果在该时间段，所有的Proxy也被封了，则重新放入sortset等待调度

AgentMaster

proxyIp+domain

(proxy启用时间，当前proxy索引,已经执行的proxy个数)

(proxyHost,Port)

Proxy数组，定期与
redis set合并

Logic

- 1.如果所有worker ip被封，启用proxy,并选择第一个proxy,更新缓存
- 2.如果被封时间间隔(当前时间-proxy启用时间 与 实体里面被封时间字段对比)到了，关掉proxy开关，启用worker 物理ip
- 3.如果所有的proxy都被封了，而还在被封状态，重新放入sortset score为(实体里面被封时间字段-当前时间-proxy启用时间)

爬虫类图

模块结构



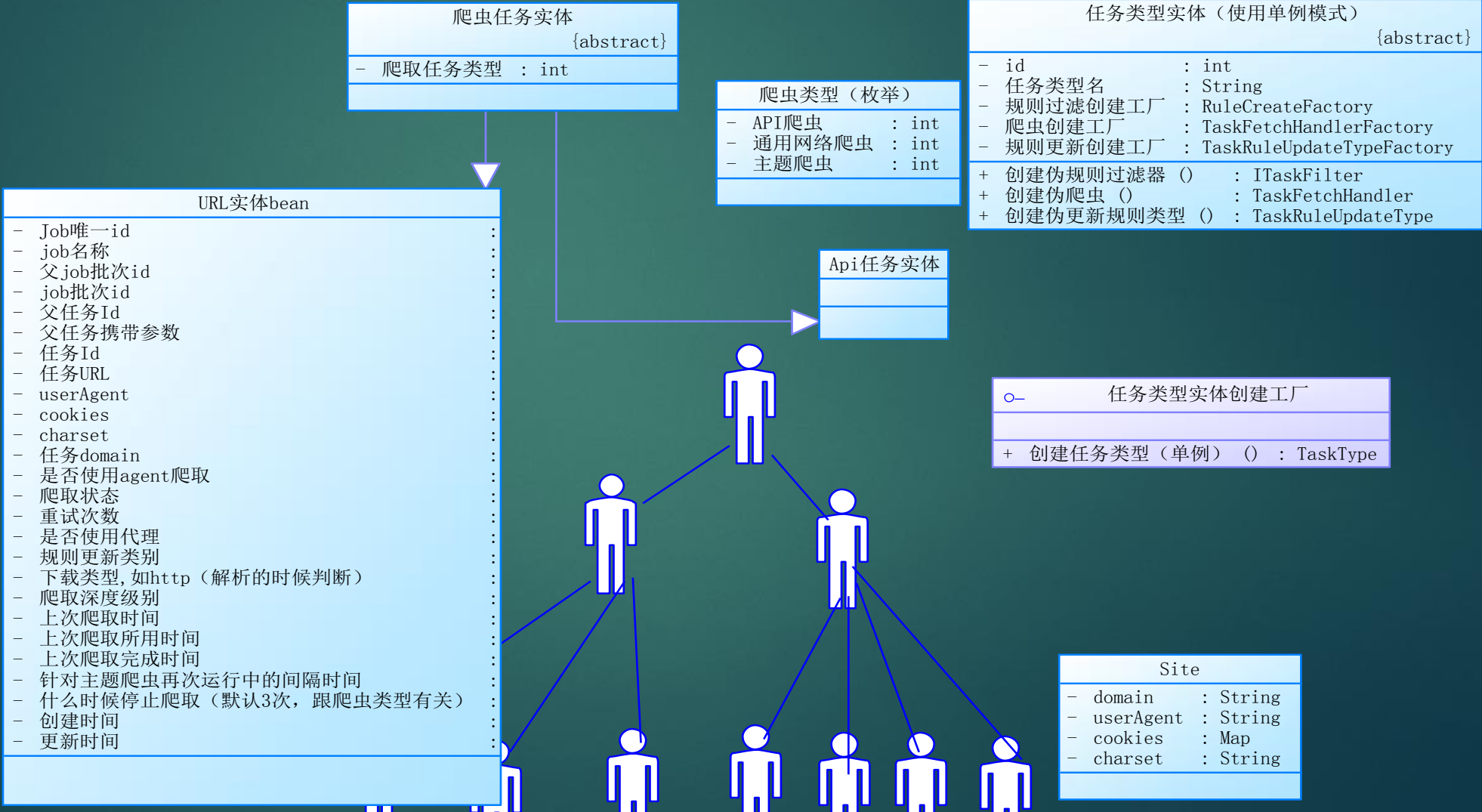
总控

去工作() {
 分析url
 匹配规则
 一致性hash调度
}

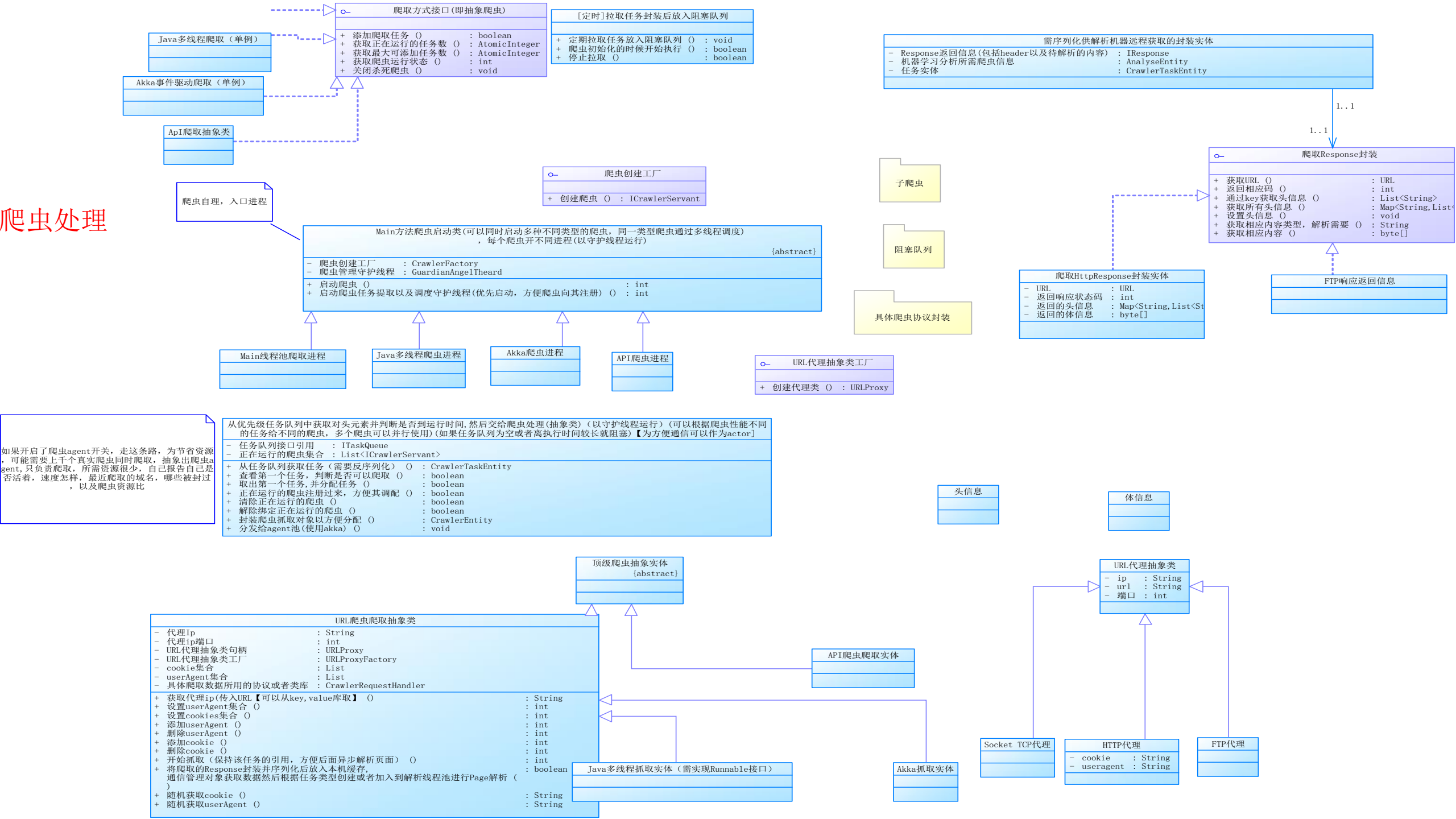
爬虫总控抽象类	
{abstract}	
- 爬虫任务实体	: CrawlerTaskEntity
- 调度接口	: 调度抽象类
+ 初始化 ()	: void
+ 封装种子任务并分配jobid ()	: CrawlerTaskEntity
+ 提交Job, 使用调度器将种子添加到对应的队列（一致性hash）[将过滤规则类及参数通过akka广播到各个worker维护，一个jobid对应多个规则过滤器] ()	: void
+ job转taskSet ()	: TaskSet
+ 提交Task ()	: void

Job类	
- jobId	: String
- job名称	: String
- Job类型(主题爬虫、通用爬虫、指定ur爬虫)	: int
- 过滤规则列表(提交种子任务的时候使用Akka广播规则，可以使用反射，先实例化，然后传入规则)	: List<ITaskFilter>
- 爬虫任务实体列表	: List<CrawlerTaskEntity>
- 规则更新	: RuleUpdateType
- URL过滤规则（定义接口，可以使用主从管理方式，通过Akka通信）	: ITaskFilter
- 父job批次id	: int
- job批次id	: int

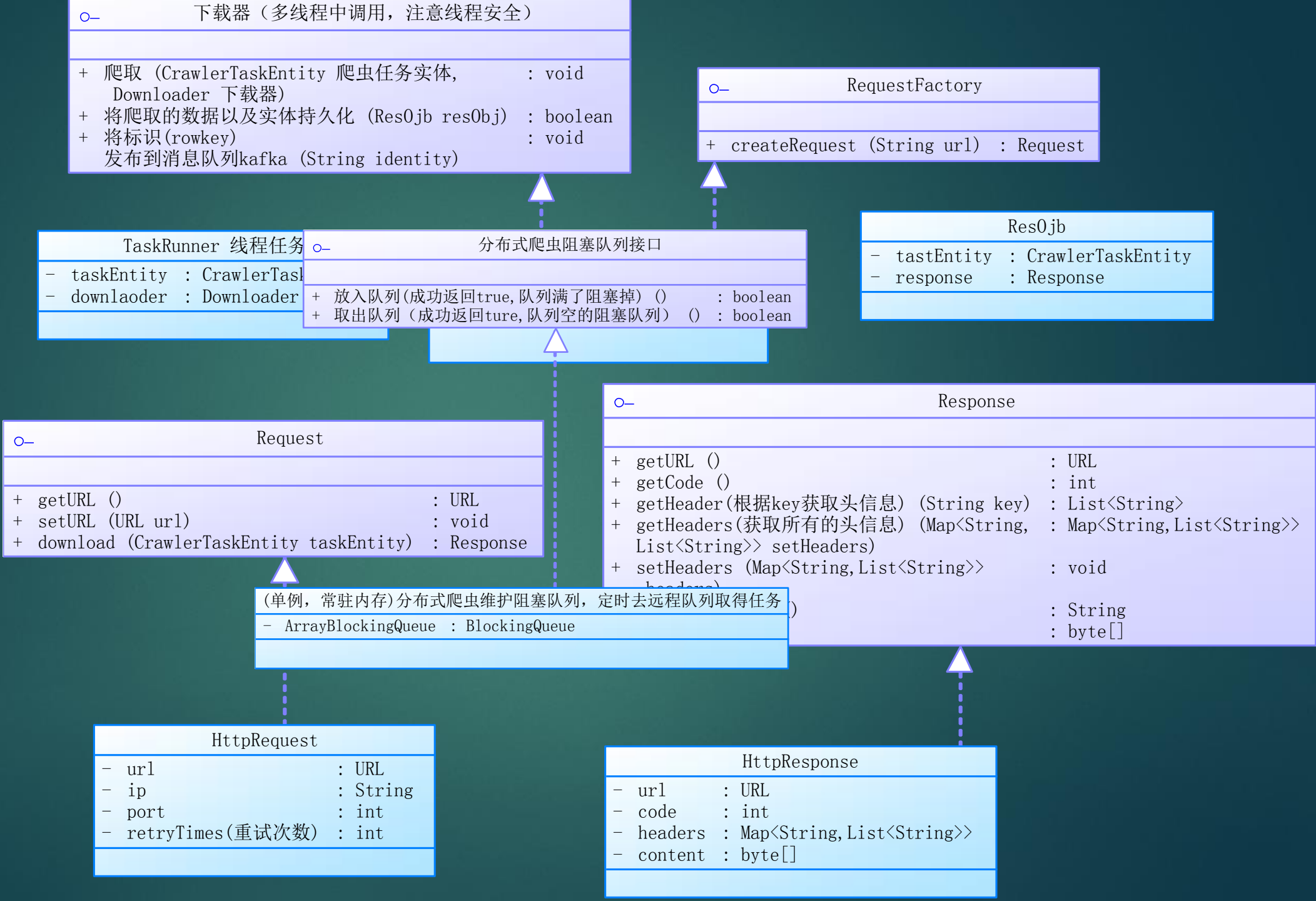
任务实体



爬虫处理



子爬虫



阻塞队列

○—

分布式爬虫阻塞队列接口

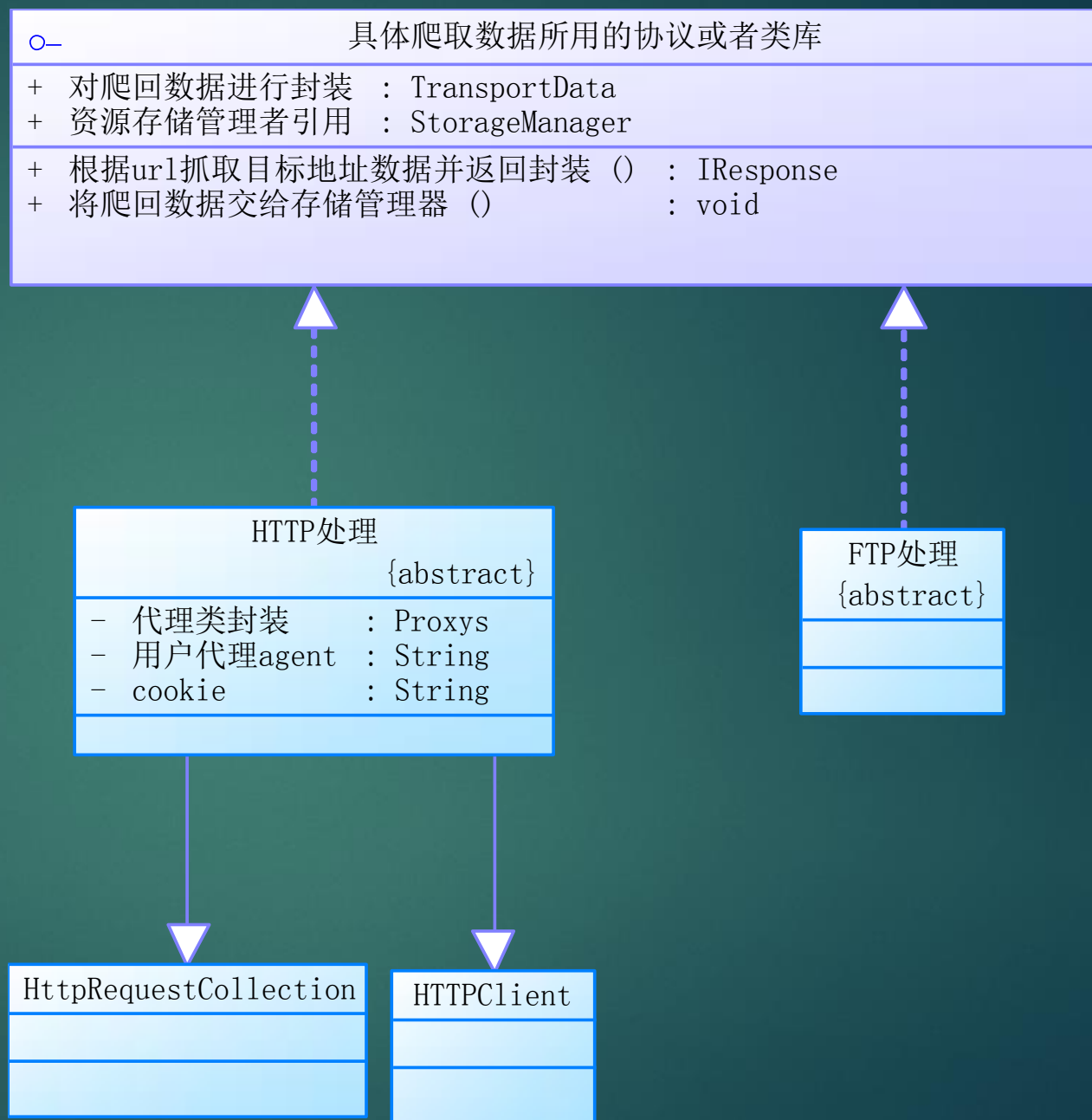
- + 放入队列 (成功返回true, 队列满了阻塞掉) () : boolean
- + 取出队列 (成功返回ture, 队列空的阻塞队列) () : boolean



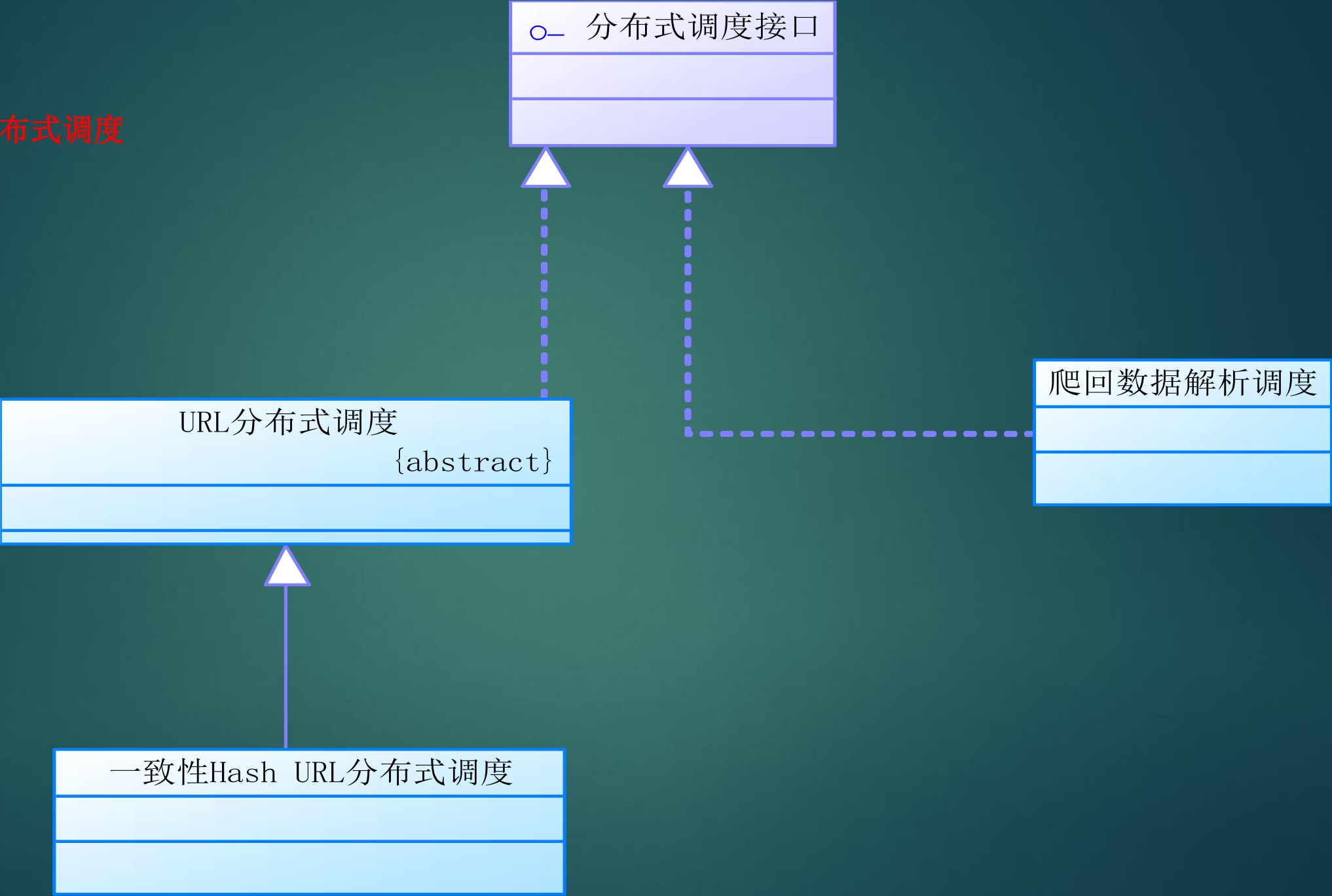
(单例，常驻内存) 分布式爬虫维护阻塞队列，定时去远程队列取得任务

- ArrayBlockingQueue : BlockingQueue

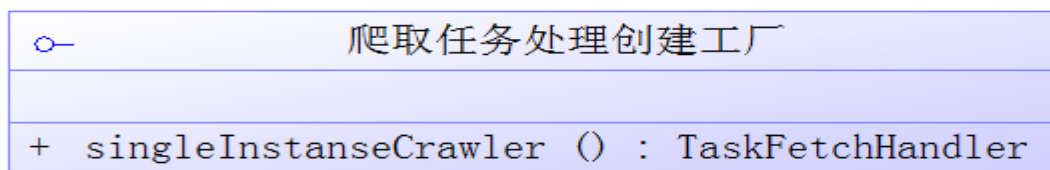
具体爬虫协议



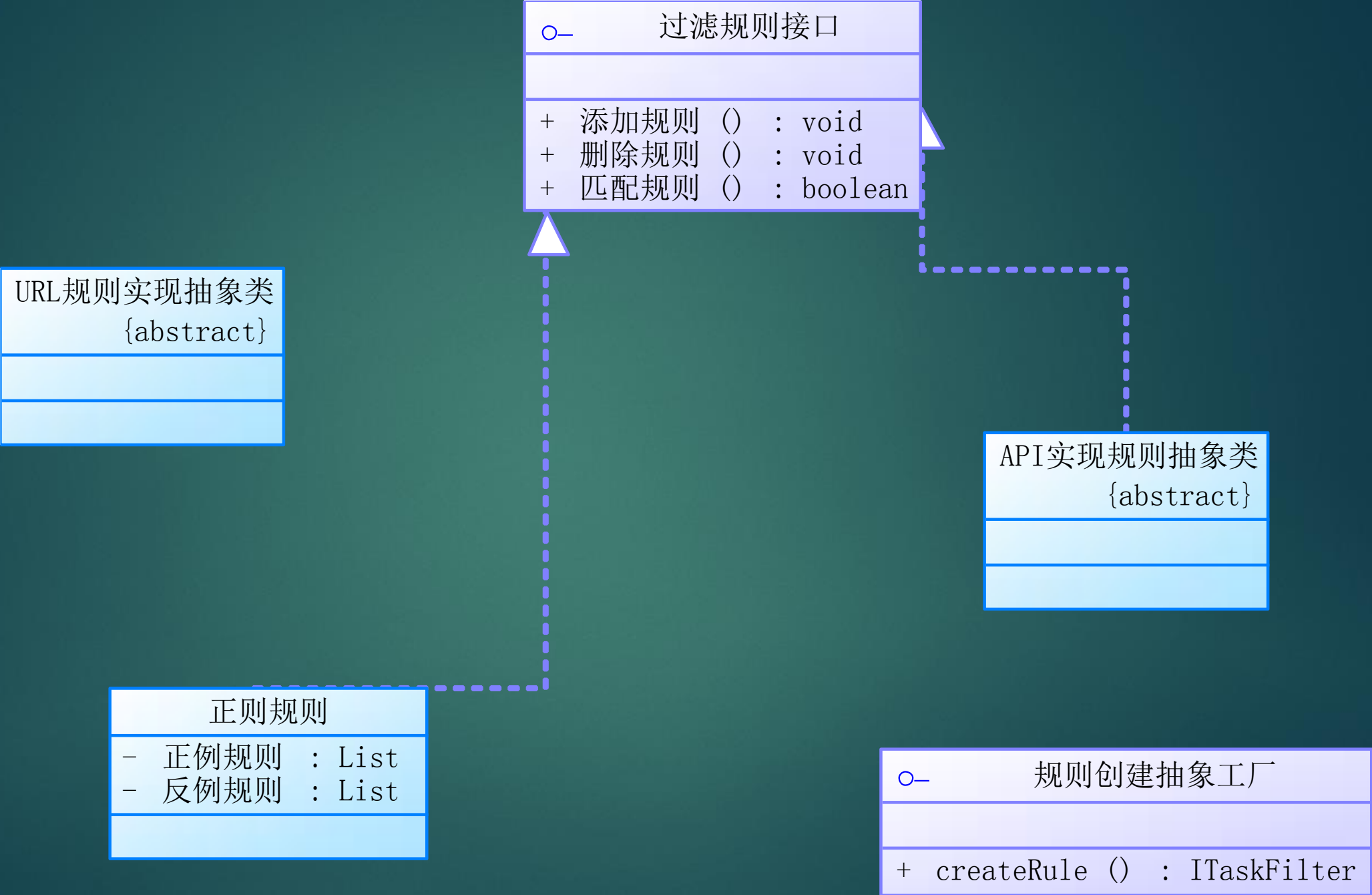
分布式调度



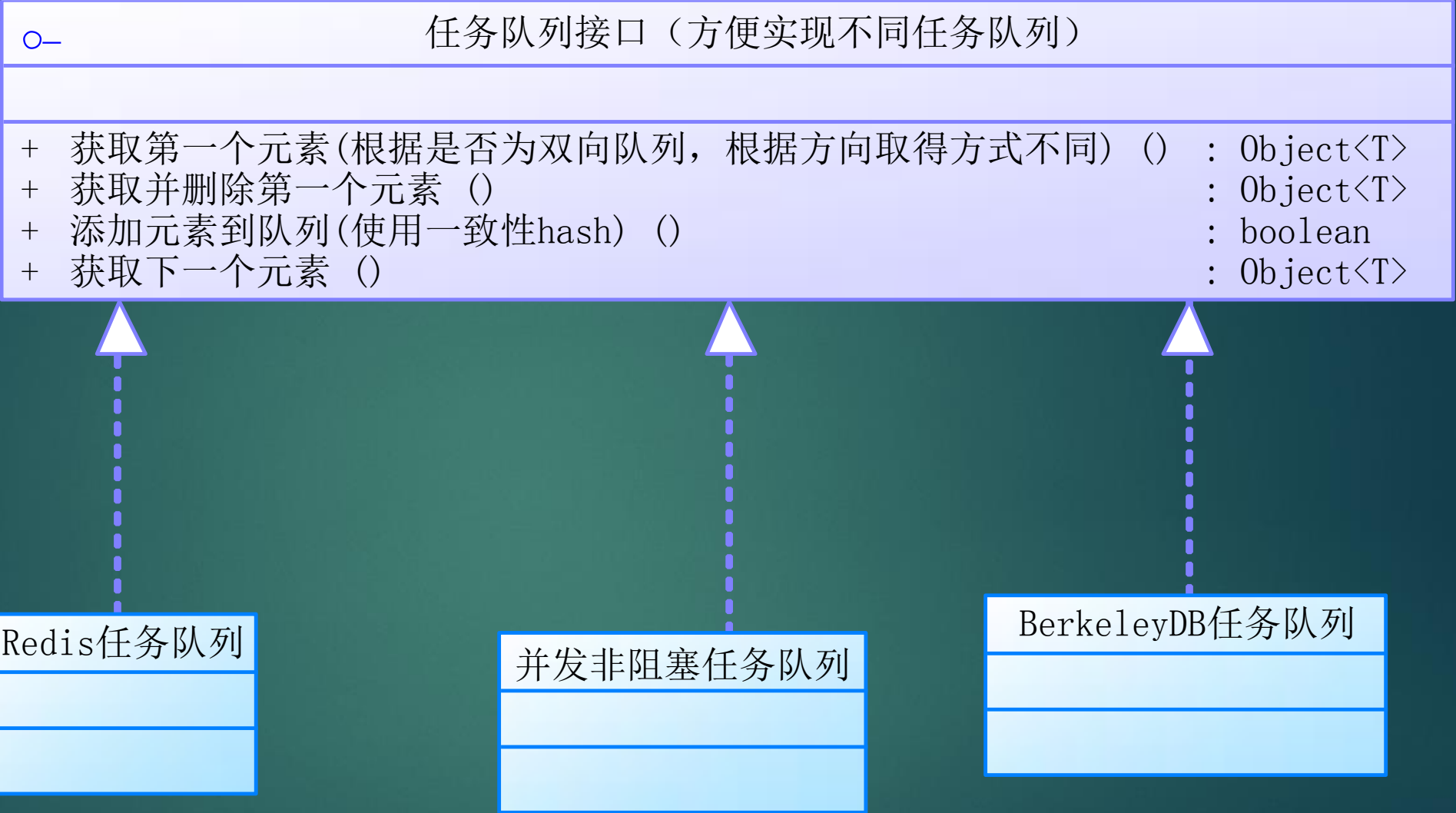
分布式爬取处理



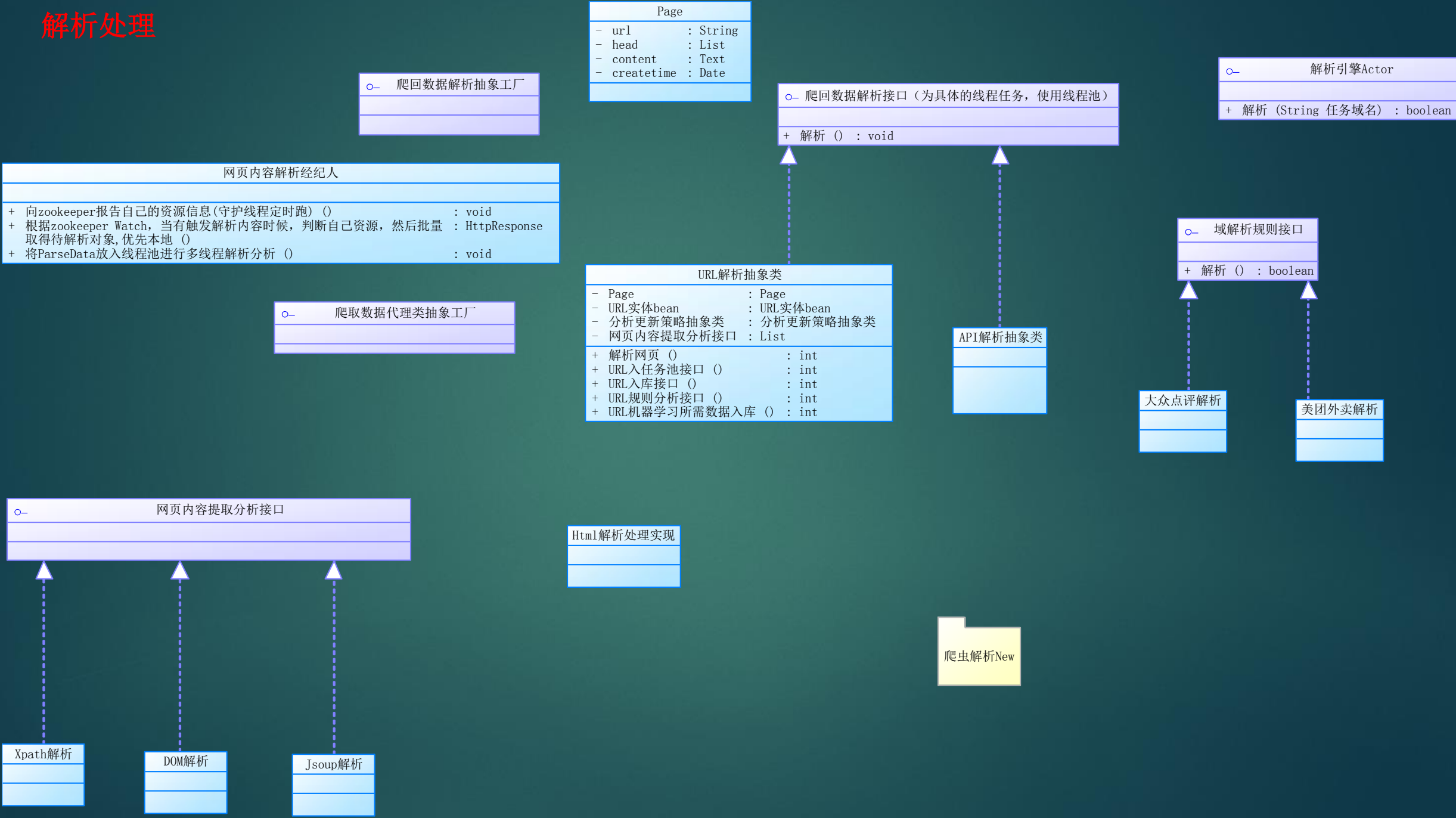
任务过滤规则

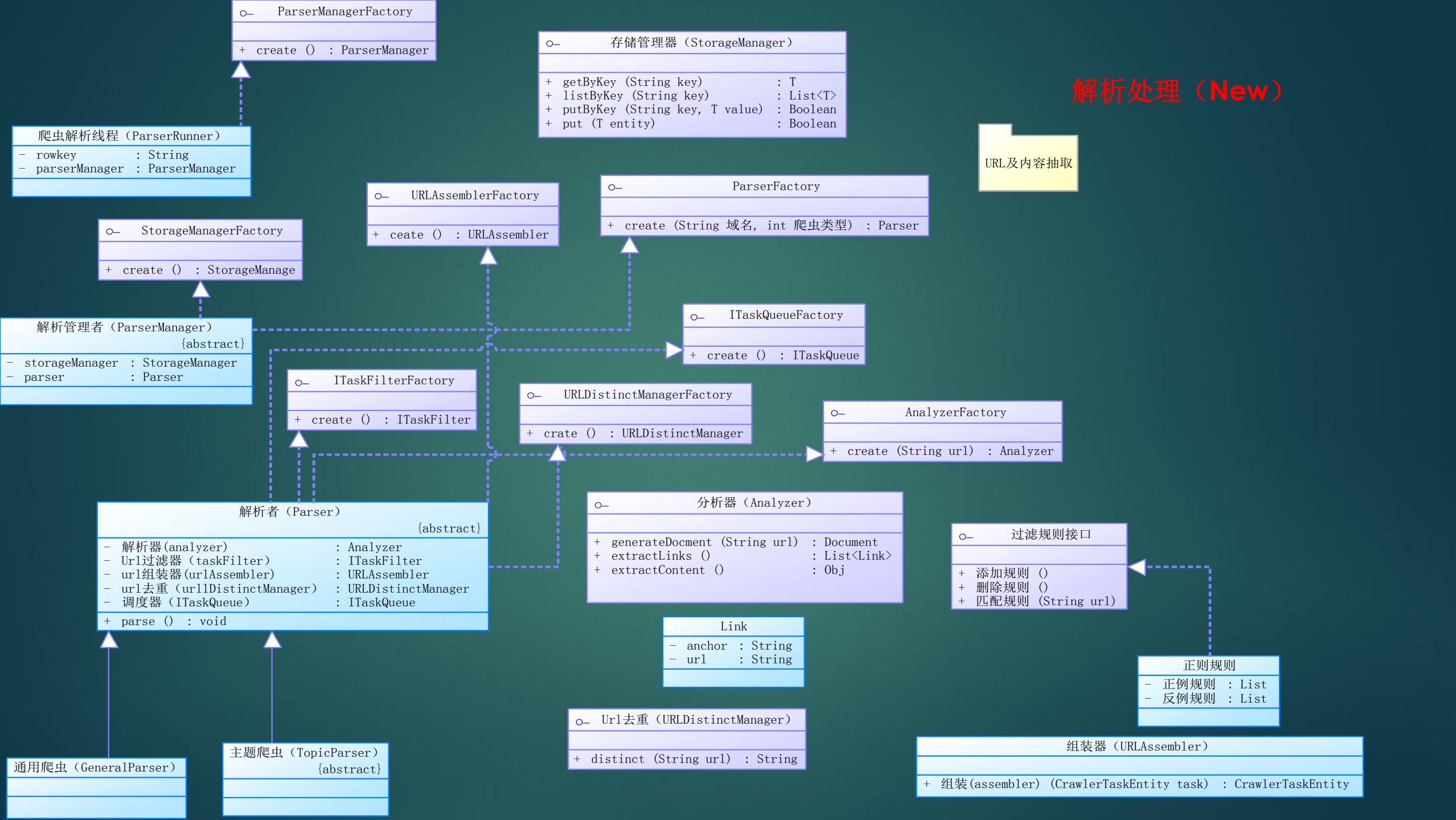


任务队列



解析处理

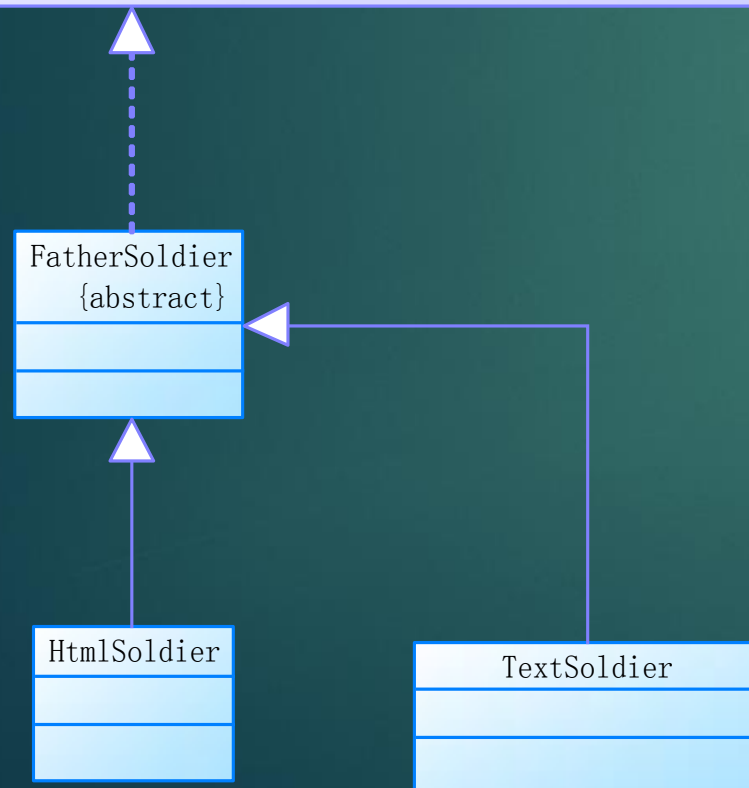




Soldier		
+ xpath (String xpath)	:	Soldier
+ \$ (String soldier)	:	Soldier
+ \$ (String soldier, String attName)	:	Soldier
+ css (String soldier)	:	Soldier
+ css (String soldier, String attrName)	:	Soldier
+ links ()	:	Soldier
+ regex (String regex)	:	Soldier
+ regex (String regex, int group)	:	Soldier
+ value ()	:	String
+ match ()	:	Boolean
+ list ()	:	List<String>
+ nodes ()	:	List<Soldier>
+ all ()	:	List<String>
+ search (TextMonster monster)	:	Soldier
+ searchList (TextMonster monster)	:	Soldier

TextMonster		
+ search (String text)	:	String
+ searchList (String text)	:	List<String>

ElementMonster		
+ search (Element element)	:	String
+ searchList (Element element)	:	List<String>



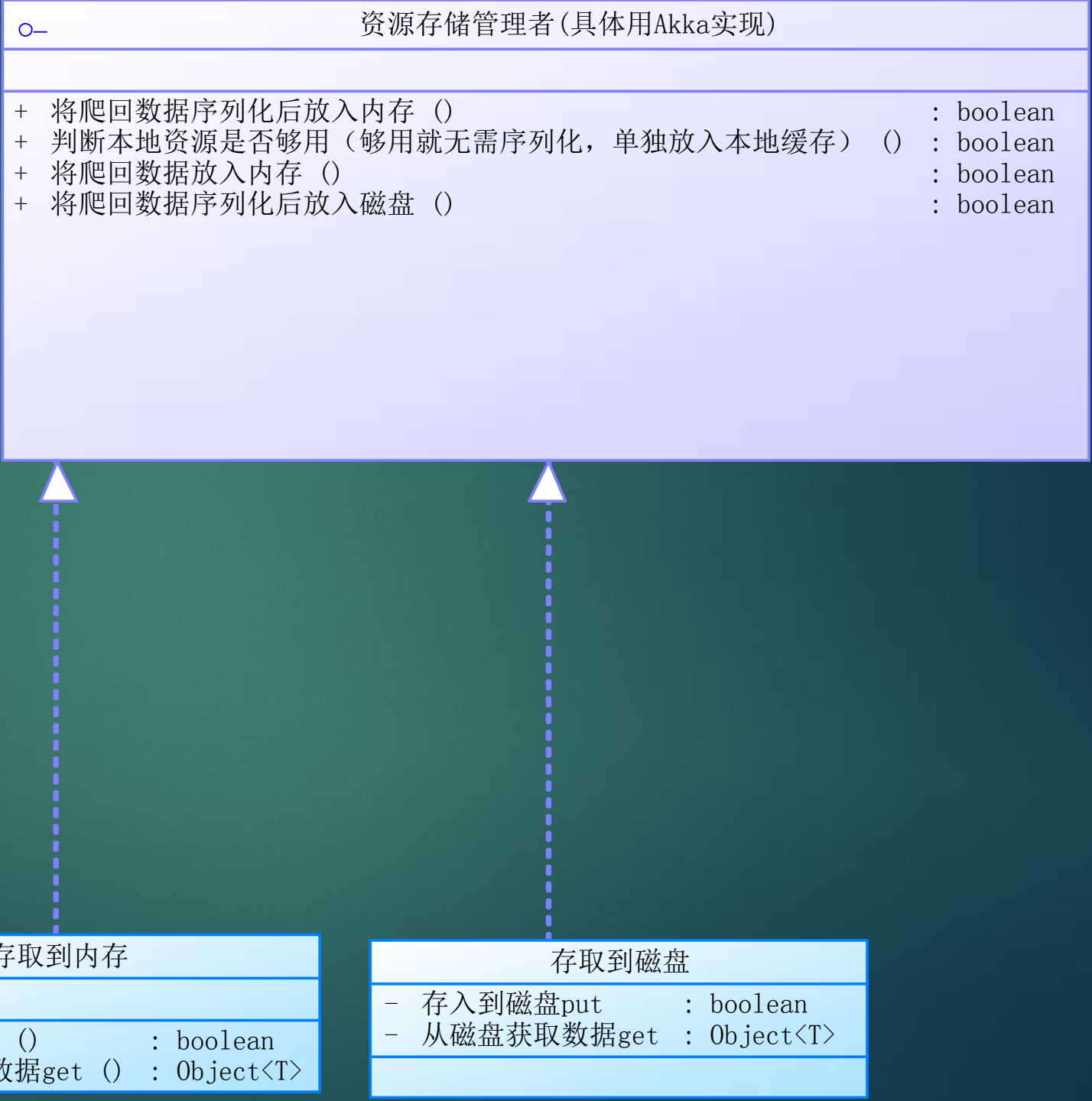
RegexMonster

CssMonster

XpathMonster

爬虫网页抽取

存储管理



爬取对象生成器

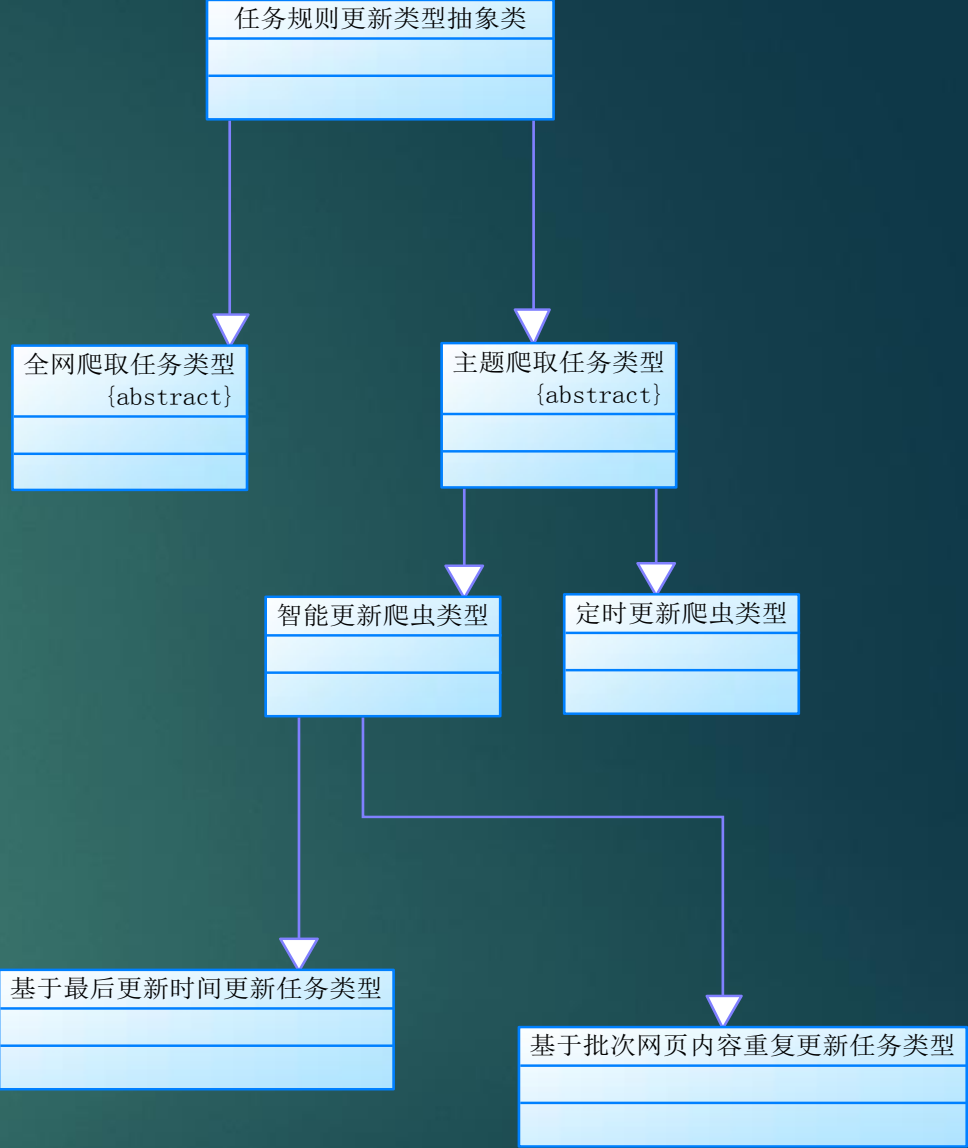
○ 爬虫url生成器 (是对爬取所需资源的一种算子抽象)

+ 获取下一个URI () : CrawlerTaskEntity

规则更新

(守护线程Main) 获得历史数据，通过机器学习建模更新任务队列score	
- 数据池	: AnalyseEntity
+ 建模到score () : Double	

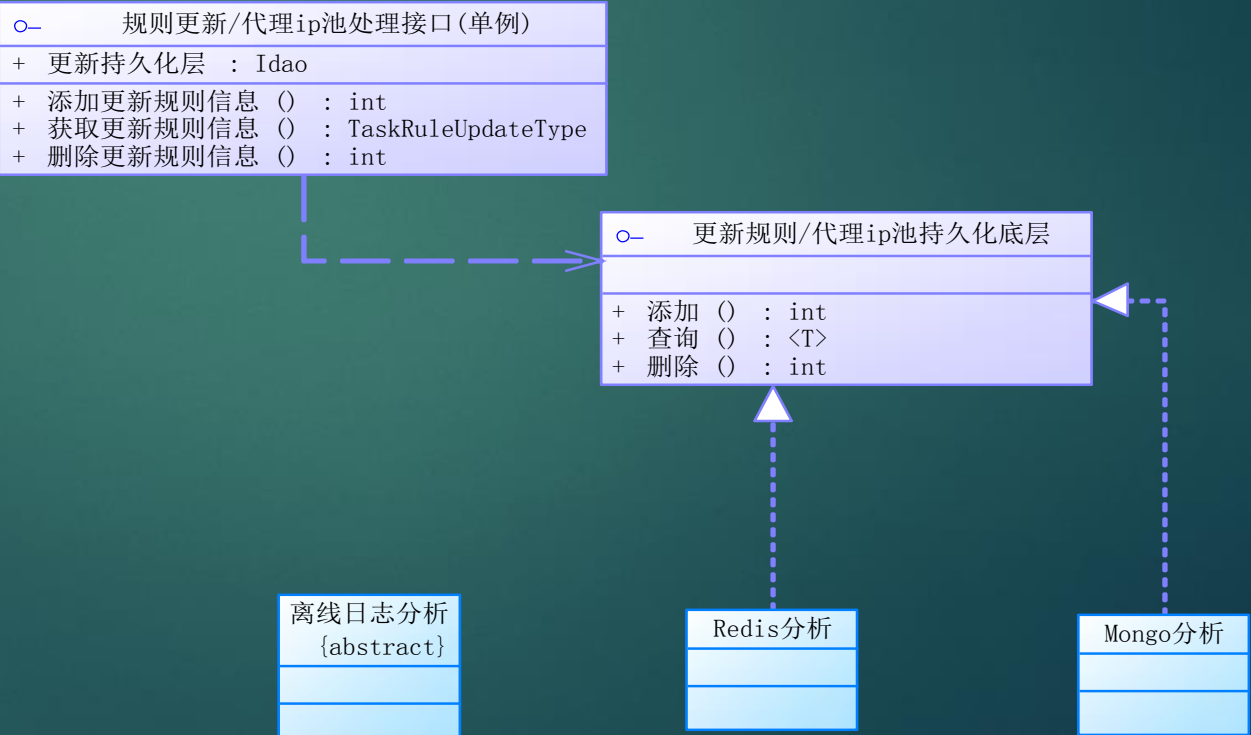
规则更新类型（枚举）	
- 无需类型	: int
- 基于最后更新时间更新任务类型（机器学习）	: int
- 基于批次网页内容重复更新任务类型（机器学习）	: int
- 定时更新爬虫类型	: int



○-	规则更新类型创建工厂
+ 创建伪任务规则更新类型 () : TaskRuleUpdateType	

机器学习数据积累

分析更新策略/代理ip池抽象类(使用守护线程分析, 最后通过线性加权转换为优先级队列的score)	
{abstract}	
- Id	: int
- 任务批次ID	: int
- 任务需要的更新类型	: String
- 任务网站上次更新时间	: date
- 抓取所使用的代理Ip	: String
- 抓取所使用的代理Ip端口	: int
- 代理ip是否可用	: int
- 上次爬取时间	: date
- 爬取所用时间	: int
- 更新频度	: int
- 爬虫类型	: String

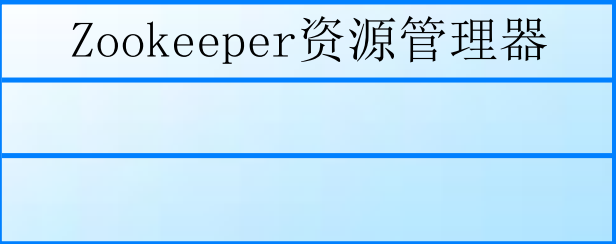
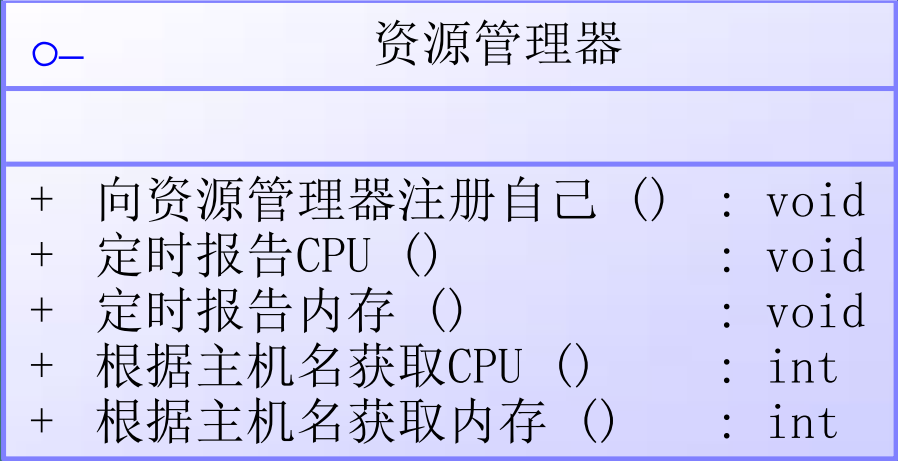


爬虫代理资源池

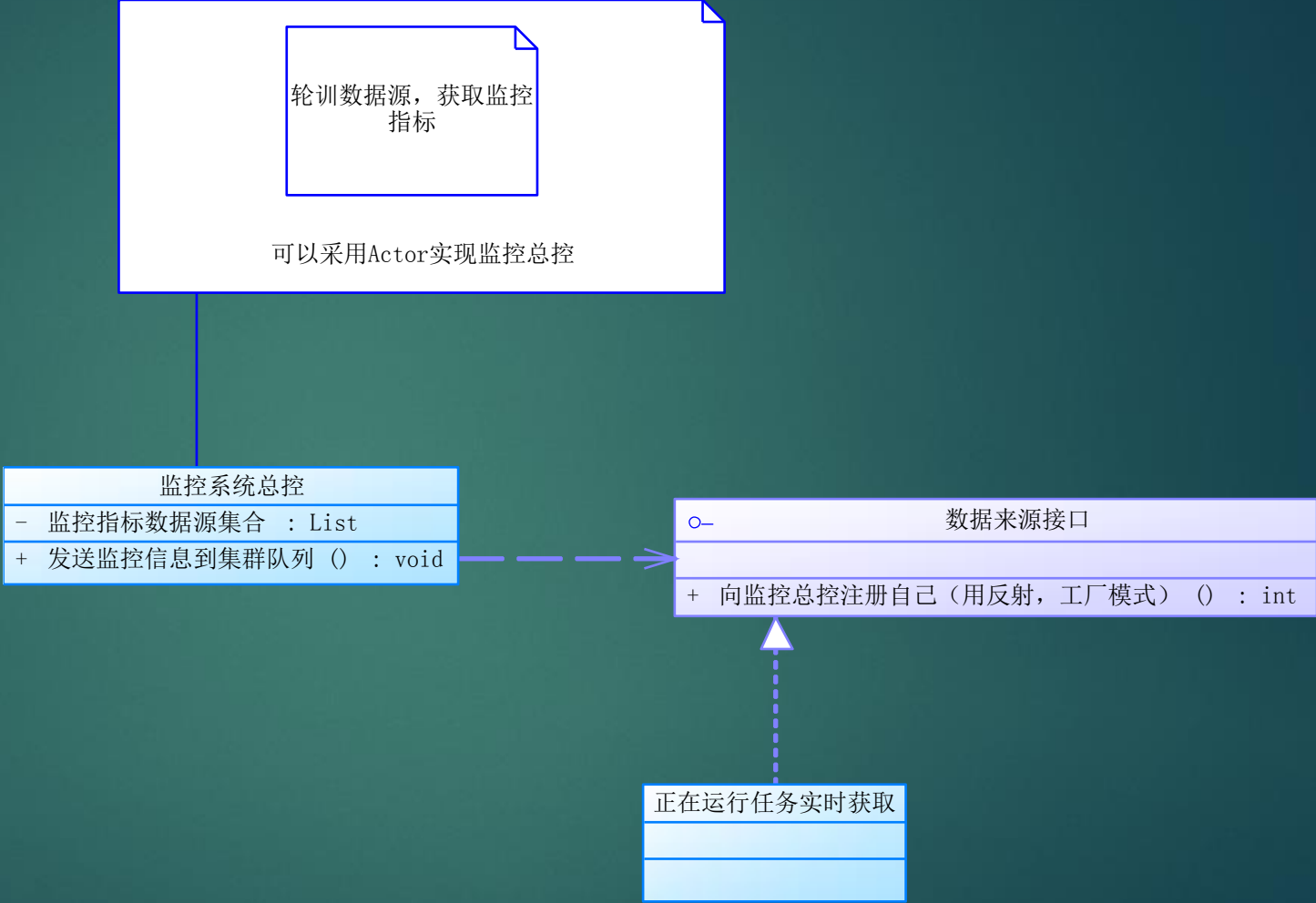
代理Agent管理者(使用Akka的actor)
+ 接收任务 () : void
+ hash函数将UrlId变为字节发送到各个小agent爬虫 () : byte[]
+ 广播发送baye[]id(即规则id) () : void
+ actor接收小爬虫过来请求 () : void

agent爬取Actor
+ 爬虫抓取管理 : GuardianAngelTheard
+ 接收规则id () : void
+ 根据规则id, 检查能否爬取(自己管理规则) () : boolean
+ 调用具体爬取方式爬取 () : void

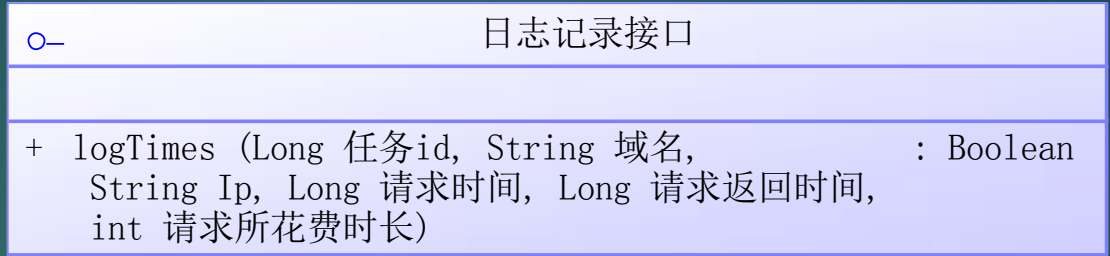
资源管理器



监控



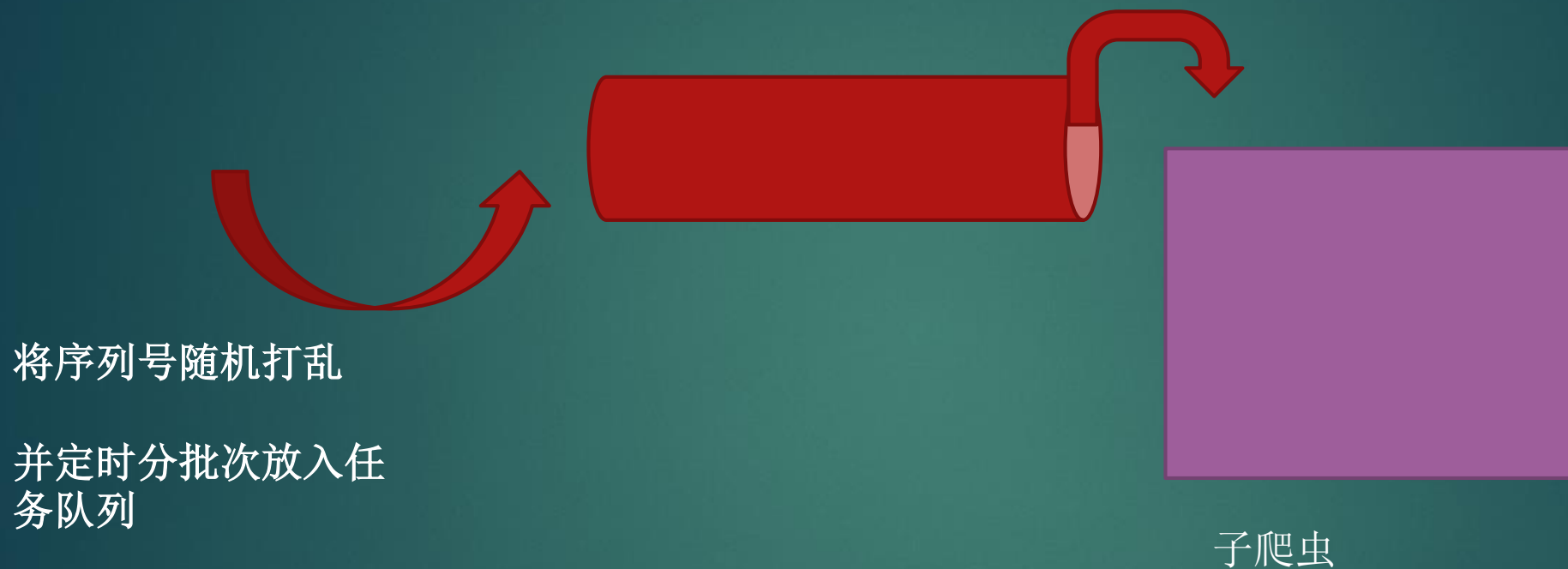
日志记录模块



美团外卖爬取（二次开发方式）

因为美团外卖url是序列递增的,可以借助这一规则穷举爬取得到餐厅信息，但是得不到电话号码
所以得到餐厅信息和地址信息后可以再去百度外卖手机和大众点评网站去爬

当输入的序列号不存在的时候，会被美团外卖重定向，而如果ip被封了也会被重定向，怎么解决



子爬虫在每次请求获得302重定向code后，定义一个全局线程安全的连续302记录器，如果超过某一阈值则认为为被封，否则认为为美团的正常餐厅下线

同时，记录下所有302所对应的url和子爬虫ip地址以及我们认为的被封的状态，为什么这么做？

因为我们虽然随机为正态分布，但是也有可能误判，我们重新拿被认为被封的url去为爬取过该url或者过了某一定时间后再去爬

美团外卖获取城市和经纬度策略

首先爬取一遍美团网，得到地址和餐厅名字

然后通过地址名字去搜狗爬取该地址所对应的经纬度

或者直接请求百度api

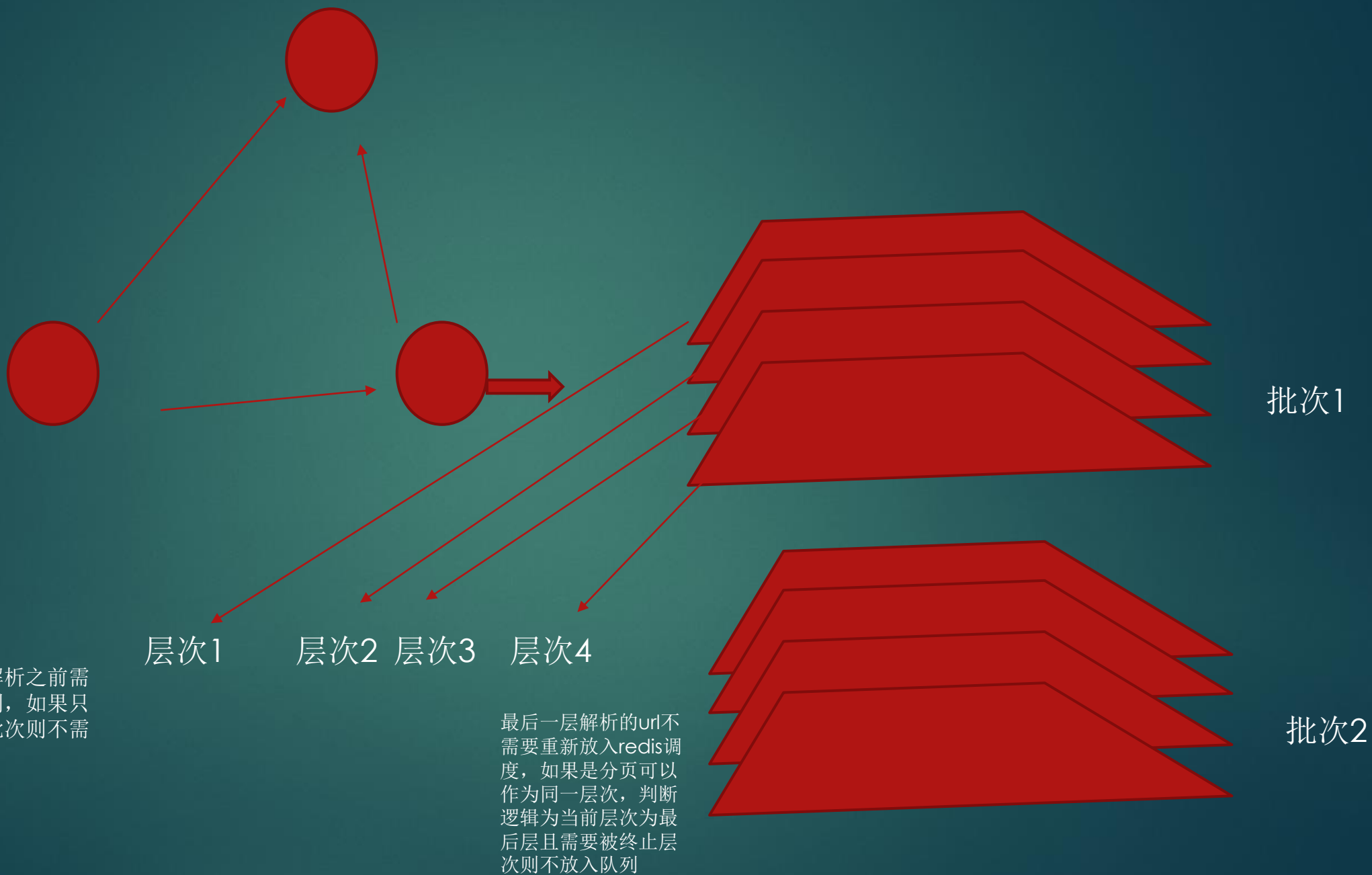
[http://api.map.baidu.com/geocoder/v2/?address="浦东大道](http://api.map.baidu.com/geocoder/v2/?address=%E6%B8%A5%E6%97%B9%E5%A4%A7%E9%81%93)

"&output=json&ak=UW8gadYahcUD3etpegMrsVhW

得到具体城市

然后再通过餐厅名在百度手机外卖和大众点评反搜得到餐厅的其他信息，比如电话号码

任务层次批次

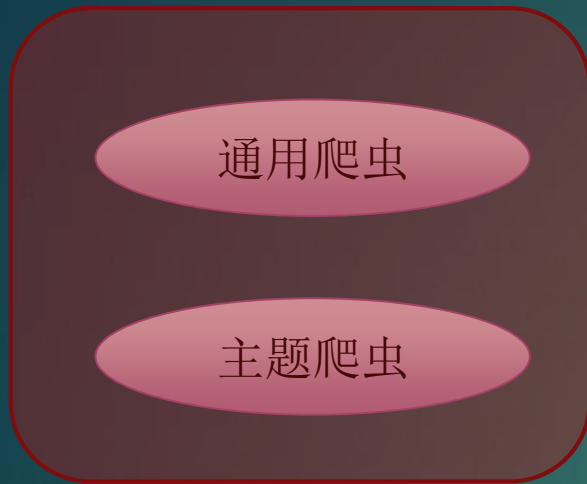


第一层在进行解析之前需要重新放入队列，如果只需要爬取一个批次则不需要

最后一层解析的url不需要重新放入redis调度，如果是分页可以作为同一层次，判断逻辑为当前层次为最后层且需要被终止层次则不放入队列

Job管理

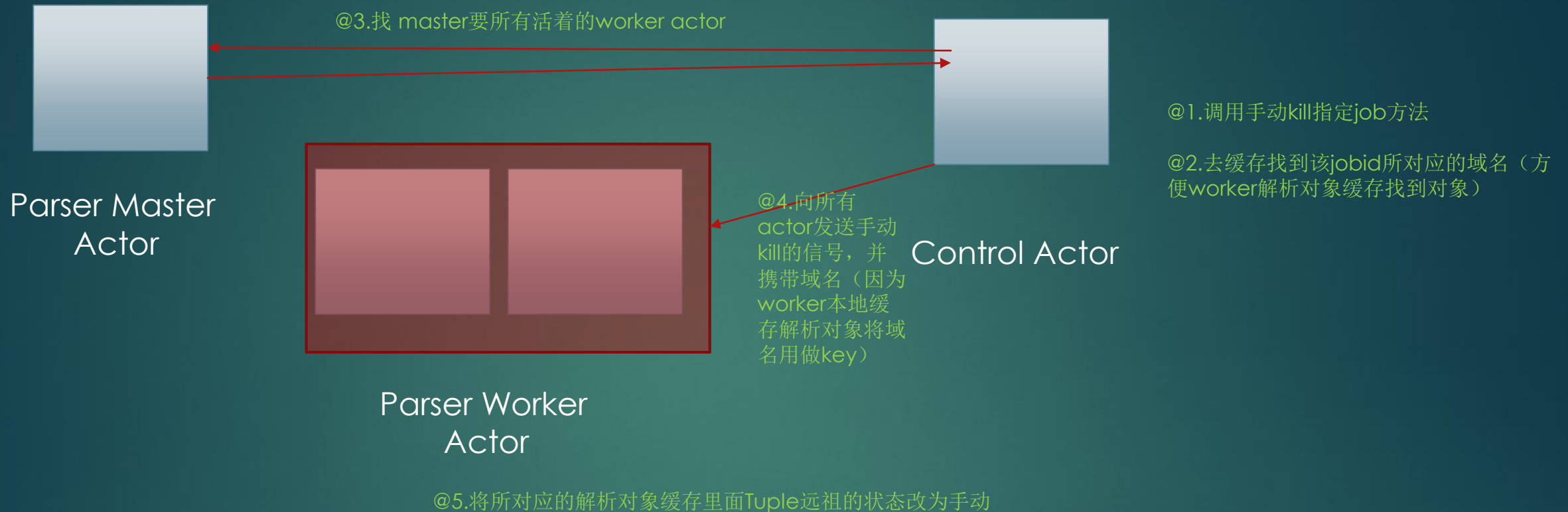
自动kill job



在ParseManager里面的start方法里面（后面）调用autoKillJob(ResObj),根据CrawlerTaskEntity的爬虫类型判断、

- 1.如果是通用爬虫，则根据爬虫深度判断，为-1则不限制，否则，任务到此为止，不需要重新放入队列
- 2.如果是主题爬虫，则根据爬虫深度+批次进行判断，都为-1则不限制，否则任务到此为止，不需要重新放入队列

手动kill job



在ParseManager里面的start方法里面（后面）调用autoKillJob(ResObj),根据CrawlerTaskEntity的爬虫类型判断、

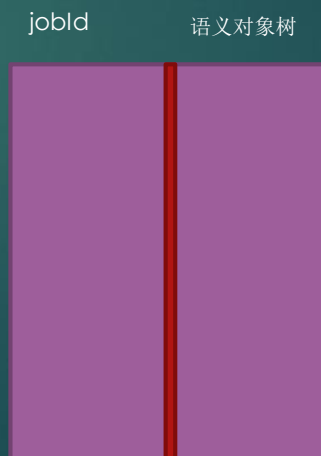
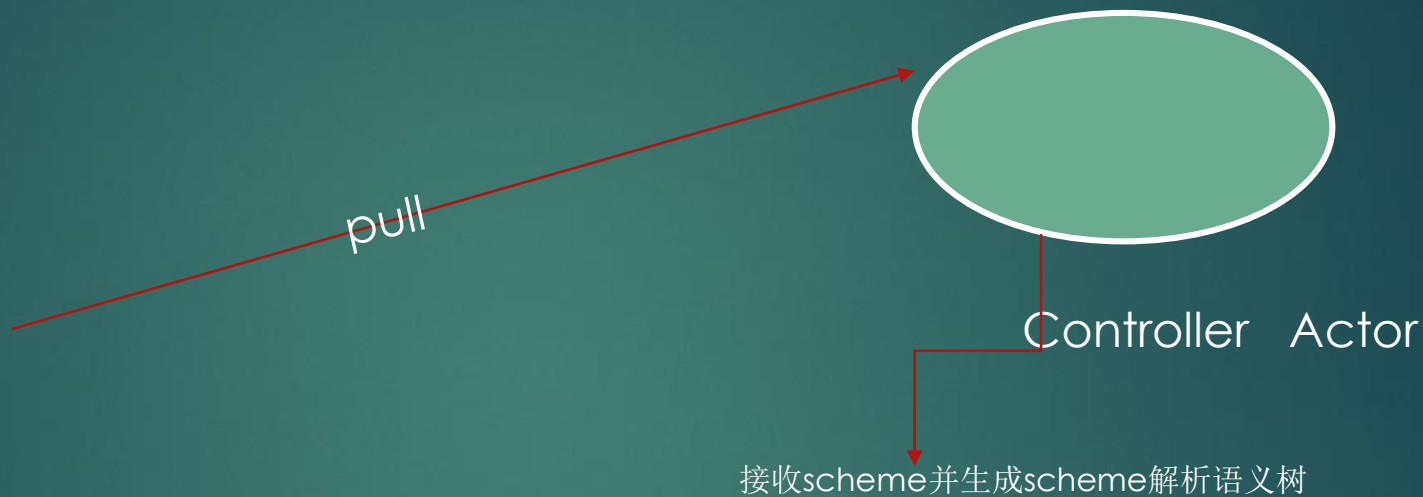
- 1.如果是通用爬虫，则根据爬虫深度判断，为-1则不限制，否则，任务到此为止，不需要重新放入队列
- 2.如果是主题爬虫，则根据爬虫深度+批次进行判断，都为-1则不限制，否则任务到此为止，不需要重新放入队列

Schema抽取



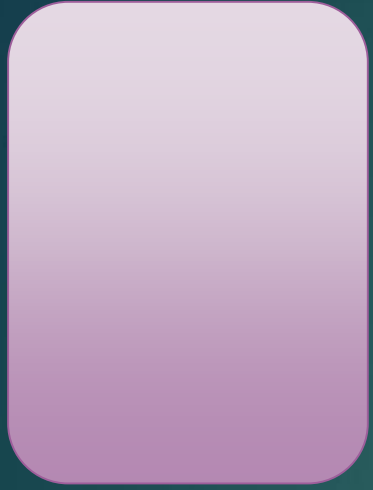
ParserWorker Actor

解析的时候先从本地缓存取，取不到到
Controller取，本地缓存一天取一次



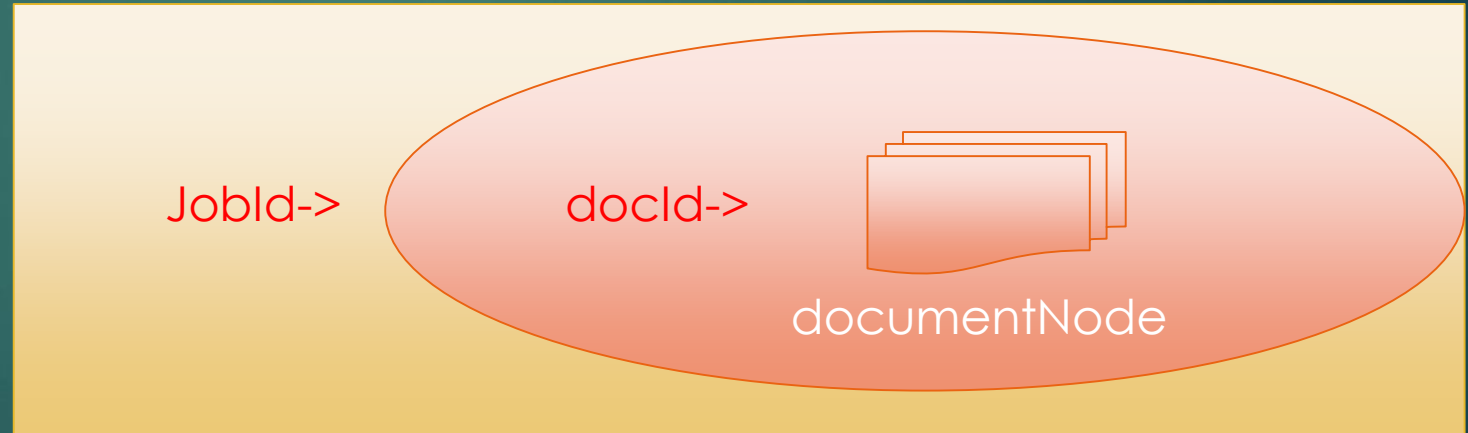
Map

Parser Worker缓存数据结构



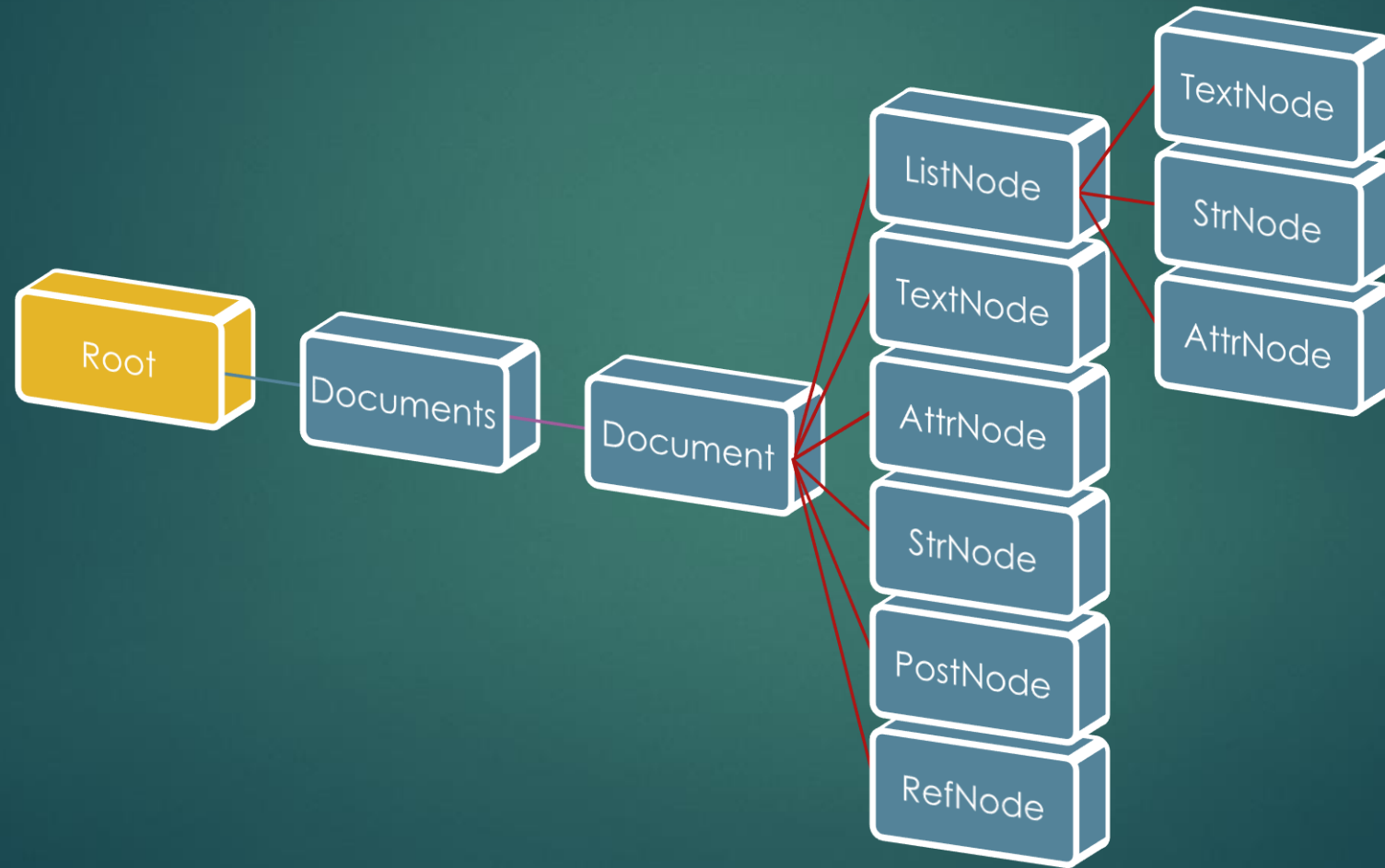
ParserWorker

$\text{Map}(\text{jobid} \rightarrow \text{Map}(\text{docId} \rightarrow \text{DocumentNode}))$



SemanticTree

pipeline



docTree逻辑

通过语义树进行内容解析后(exec())

1. 如果语义树解析完内容为空，则说明该docNode没内容解析，则直接使用引用(递归)该docNode的docNode的docId重新取得被应用的语义树
2. 根据树节点逻辑进行处理，并根据page或者pageSize判断是否需要分页，请求支持post,get，解析支持Json,html
3. Schema支持正则，xpath表达式混合编写

Sheme样例

```
<!-- 一个doc 对应一个网页, url 表示网页 -->
<doc docid="dipingcity"
  url="http://www.dianping.com/ajax/json/index/citylist/getCitylist?do=allCitylist&nr_force=1439875530364"
  method="get" firstdoc="true">

  <list jsonselector="msg->html" name="result" notsave="true">
    <text name="cityid" parameter="true" separator="|" index="3" notsave="true"/>
    <!-- jsonselector="cityid" you can omit it ,default we choose name as jsonselector, also you can use regex el-->
  </list>

  <ref docid="dipingnavcity"/>
</doc>

<!--
http://www.dianping.com/search/category/${cityid}/10#nav-tab|0|1
http://www.dianping.com/search/category/${cityid}/10#nav-tab%7c0%7c1
-->
<doc docid="dipingnavcity" url="http://www.dianping.com/search/category/${cityid}/10" method="get">
  <list selector="div.nav-category.nav-tabs div.nc-contain div#J_nt_items div#region-nav a" name="result"
    isfiled="false">
    <attr attr="href" name="href" parameter="true" notsave="true" filterregex="([\\w/]+)#" filterGroup="1"/>
  </list>

  <text selector="body#top div.section.Fix div.bread.J_bread span.num" name="total" filterRegex="(\\d+)"
    filterGroup="1"
    notsave="true"/>

  <ref docid="dipingsubnavcity"/>
</doc>

<doc docid="dipingsubnavcity" url="http://www.dianping.com${href}" method="get">
  <list selector="div.nav-category.nav-tabs div.nc-contain div#J_nt_items div#region-nav-sub a" name="result">
    <attr attr="href" name="href" parameter="true" notsave="true"/>
  </list>
  <text selector="body#top div.section.Fix div.bread.J_bread span.num" name="total" filterRegex="(\\d+)"
    filterGroup="1"
    notsave="true"/>

  <!-- if there is no total item number, you also can use totalpagesize -->

  <ref docid="dipinglist"/>
</doc>

<doc docid="dipinglist" firsturl="http://www.dianping.com${href}"
  url="http://www.dianping.com${href}/p${page}?aid=bf3133f6cf5079c59b8c8b315bccac026dc74ae45e529cc10368c3562c4b17f2715152b6bdfb29834f8e761cbec7f47c8cc713db7df194db5557cff3892124763190b89756f9e4156a7f2016980b11045&tc=2"
  method="get" page="true" maxPage="50" pagesize="15" firstpagesize="17">
  <list selector="div.content-wrap div.shop-wrap div.content div.shop-list.J_shop-list.shop-all-list ul li"
    name="result">
    <attr selector="div.pic a" name="image" attr="href" index="0"/>
    <!-- index default 0 -->
    <text selector="div.txt div.tit a h4" name="title"/>
    <text selector="div.tag-addr a span" name="tag" index="0"/>
    <text selector="div.tag-addr a span" name="address" index="1"/>
    <!-- default index=0 -->
    <attr attr="href" selector="div.txt div.tit a" name="href" parameter="true" notsave="true"/>
  </list>

  <text selector="body#top div.section.Fix div.bread.J_bread span.num" name="total" filterRegex="(\\d+)"
    filterGroup="1"
    notsave="true"/>
  <ref docid="dipingdetail"/>
</doc>

<doc docid="dipingdetail" url="http://www.dianping.com${href}">
  <text selector="div.body-content.clearfix div.main div#basic-info p.expand-info.tel span.item" name="tel"
    type="multiple"
    separator="|"/>
</doc>

<!-- post test -->
<!--<doc docid="dipingdetail" url="http://www.dianping.com/${href}">
  <post>
    <str name="a" value="haha" notsave="true"/> [default notsave=false ]
    <str name="title" notsave="true">${page}</str>
  </post>
-->
```

任务跟踪系统


```
/**
 * Created by soledede on 2015/9/17.
 */
trait TraceListener {
  def onJobStart(jobstart: JobStarted) = {}
  def onJobTaskFailed(jobTaskFailed: JobTaskFailed) = {}
  def onJobTaskCompleted(jobTaskCompleted: JobTaskCompleted) = {}
  def onJobTaskAdded(jobTaskAdded: JobTaskAdded) = {}
}
```

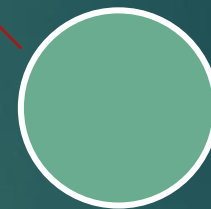


Listeners

Add listener

异步线程，实时执行，阻塞通过信号量通信

Add event



ListenerEvents

同步锁

REDIS

to

```
/**
 * Created by soledede on 2015/9/17.
 */
sealed trait TraceListenerEvent
case class JobStarted(jobId: String, seedNum: Int) extends TraceListenerEvent
case class JobTaskFailed(jobId: String, num: Int) extends TraceListenerEvent
case class JobTaskCompleted(jobId: String, num: Int) extends TraceListenerEvent
case class JobTaskAdded(jobId: String, num: Int) extends TraceListenerEvent
```

该job当前总任务数

该job当前已完成任务数

该job当前失败任务数

监控系统

