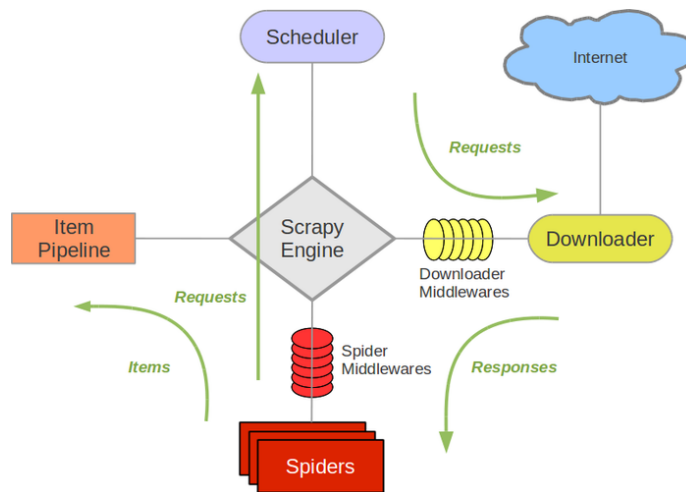


Scrapy 爬虫框架 (基于 Twisted 编写-异步)

架构图



组件

- Scrapy Engine** — Scrapy 引擎负责控制数据流在系统中所有组件中流动，并在相应动作发生时触发事件
- 调度器 (Scheduler)** — 调度器从引擎接受 request 并将他们入队，以便之后引擎请求他们时提供给引擎。
- 下载器 (Downloader)** — 下载器负责获取页面数据并提供给引擎，而后提供给 spider。
- Spiders** — Spider 是 Scrapy 用户编写用于分析 response 并提取 item (即获取到的 item) 或额外跟进的 URL 的类。每个 spider 负责处理一个特定 (或一些) 网站。
- Item Pipeline** — Item Pipeline 负责处理被 spider 提取出来的 item。典型的处理有清理、验证及持久化 (例如存取到数据库中)。
- 下载器中间件 (Downloader middlewares)** — 下载器中间件是在引擎及下载器之间的特定钩子 (specific hook)，处理 Downloader 传递给引擎的 response。其提供了一个简便的机制，通过插入自定义代码来扩展 Scrapy 功能。
- Spider 中间件 (Spider middlewares)** — Spider 中间件是在引擎及 Spider 之间的特定钩子 (specific hook)，处理 spider 的输入 (response) 和输出 (items 及 requests)。其提供了一个简便的机制，通过插入自定义代码来扩展 Scrapy 功能。

流程

- Scrapy Engine 打开一个网站，找到处理该网站的 Spider 并向该 spider 请求第一个要爬取的 URL
- Scrapy Engine 从 Spider 中获取到第一个要爬取的 URL 并在调度器 (Scheduler) 以 Request 调度。
- Scrapy Engine 向调度器请求下一个要爬取的 URL
- 调度器返回下一个要爬取的 URL 给引擎，引擎将 URL 通过下载中间件 (请求 (request) 方向) 转发给下载器 (Downloader)。
- 一旦页面下载完毕，下载器生成一个该页面的 Response，并将其通过下载中间件 (返回 (response) 方向) 发送给 Scrapy Engine
- Scrapy Engine 从下载器中接收到 Response 并通过 Spider 中间件 (输入方向) 发送给 Spider 处理。
- Spider 处理 Response 并返回爬取到的 Item 及 (跟进的) 新的 Request 给 Scrapy Engine。
- Scrapy Engine 将 (Spider 返回的) 爬取到的 Item 给 Item Pipeline，将 (Spider 返回的) Request 给调度器
- (从第二步) 重复直到调度器中没有更多地 request，Scrapy Engine 关闭该网站。