

Architecting a Scalable Web Scraping Solution



Neha Setia Nagpal Developer Advocate at Zyte

" My role is to serve the web scraping developer's community."







8 steps to architecting a scalable web scraping solution.



Step 1: Clarify The Goal - Why?

What is the business problem you want to solve with the data?

How will you extract that data?

What kind of data do you need?

Where would you find that data?

zyte | 5

Step 1: Clarify The Goal - Why?

The Purpose is to build

that will help you to solve

problem

use case

for which i need to extract

define the data schema to solve the problem

of approx size - define how much

fromdefine Target Websites

In time and

refresh in days.

Step 2: Analyze The Website



Are there any APIs available?

2 Does the crawl involves dealing with dynamic elements on the web pages?

3 +

How powerful the anti-bot mechanisms are on the website?

4

Is it a crawler for text-only web pages only? Or should we fetch and store other types of media type as well?



What is the expected number of pages we will crawl?



What is the average size of the web page to be downloaded?

Step 3: Prioritize The Project Attributes





Scalability

Distributed Crawling.

Extensibility

Single Responsibility Principle.



Availability

Low Latency and resilient monitoring systems.



Velocity

Strategize on the basis of the project needs.

Step 4: Highlight The Constraints



Step 5: Design The Crawl



How to crawl? Which algorithm to use for scanning the pages? BFS/DFS or Path ascending crawls



How do we prevent requesting the same URL?

- Same Crawling Session- Scrapy handles it by default.
- Different Crawling Session- Frontier Service



How to make the crawler polite? Fewer Requests



Validate the design for Single responsibility principle- look for "and"



Step 6: Plan The Quality Assurance Process

We have determined 6 data quality dimensions



Step 7: Choose The Tech Stack



Base Technology/Framework- Scrapy



Deploy and Maintain Spiders- Scrapy Cloud



Rotating Proxy and Antiban Solution-Zyte API - including Smart Proxy Manager/Smart Browser.



Browser Automation Tools- Headless Browser libraries



Maintenance and Monitoring- Spidermon and Matt



Step 8: Brace For Impact

Brace the uncertainty by preparing for the challenges in maintaining a large scale web scraping project by

Low latency Monitoring Tools like spidermon.

To Summarize



Thank you