

Functional Requirements

- Should be able to parse html content and extract urls from a web page.
- Should be able to do malware detection.
- Should be able to detect duplicate urls.
- repeat the whole process infinitely.

Non Functional Requirements

- Web crawler should be scalable because there are billions of pages to be scaled. Web crawler should efficiently use parallelization.
- Should not overwhelm the website servers with too much requests otherwise it would lead to DDOS attacks.
- Robust enough to handle poorly formatted HTMLs, unresponsive servers and web crashes.
- Ability to include future changes, flexible and extensible enough.
- Should have a robust logging and monitoring feature to keep a track of metrics and overall health of web crawler.

QPS Estimation

- Assuming 1 billion pages to be crawled every month.
- $1000000000 / 30 \text{ days} / 24 \text{ hours} / 60 \text{ minutes} / 60 \text{ seconds} = 385 \text{ QPS}$
- Peak value assumed to be a little more than twice that is 800 QPS.

Storage Estimation

- Assuming average page size is 500 KB,
- 1 billion pages require 500 TB storage a month
- If the pages need to be stored for a year , then we require $500 \text{ TB} * 12 = 6 \text{ PB}$ of storage a year.

