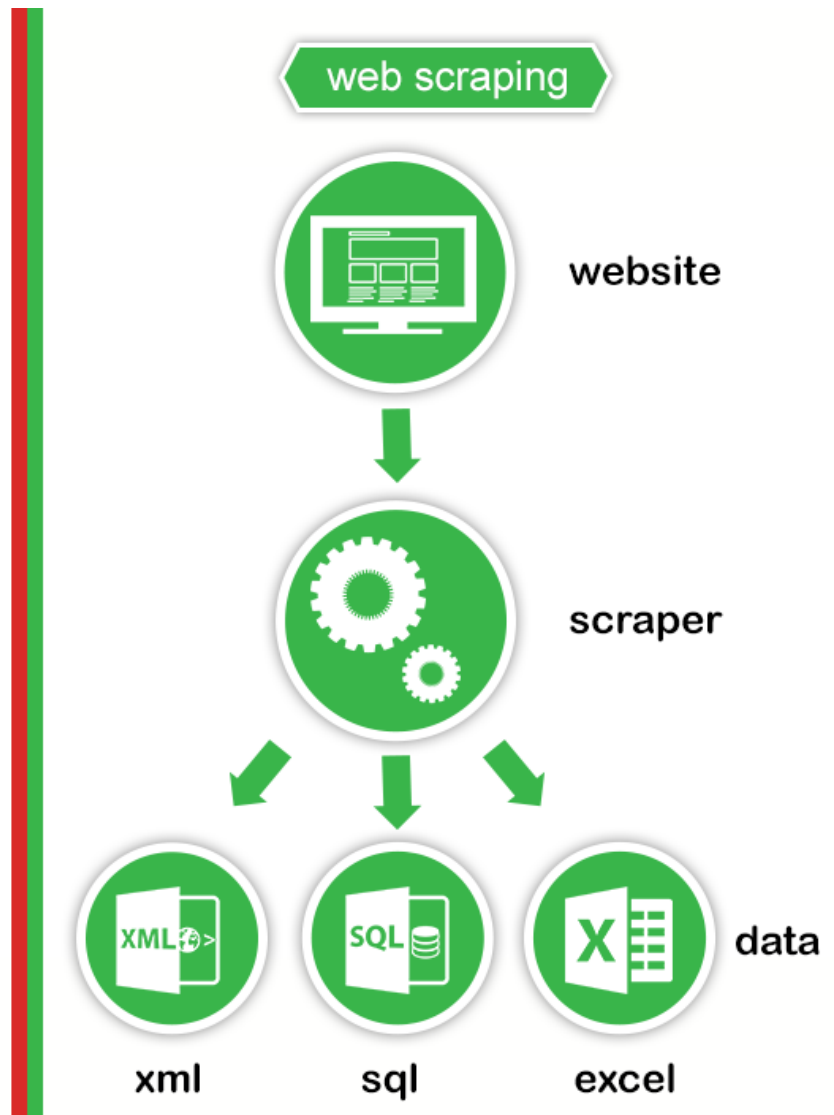


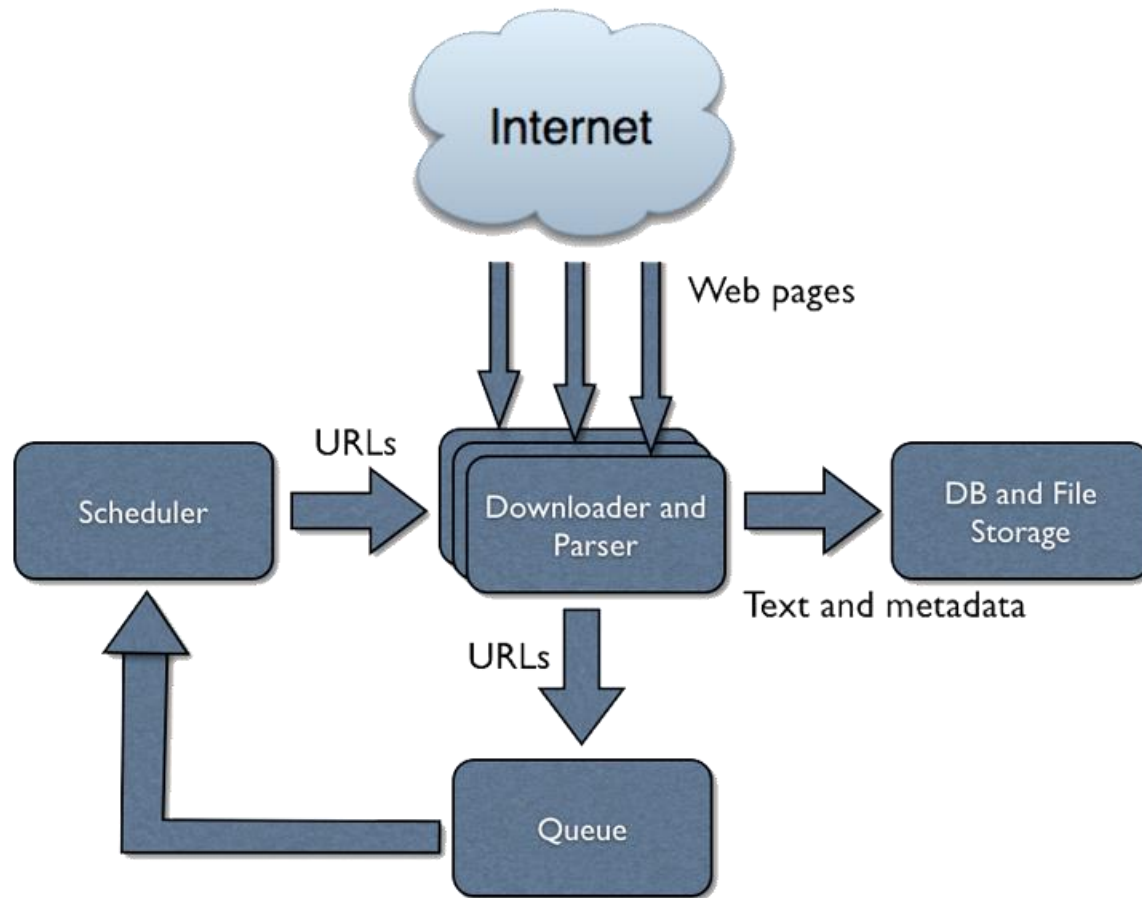
# Scraping

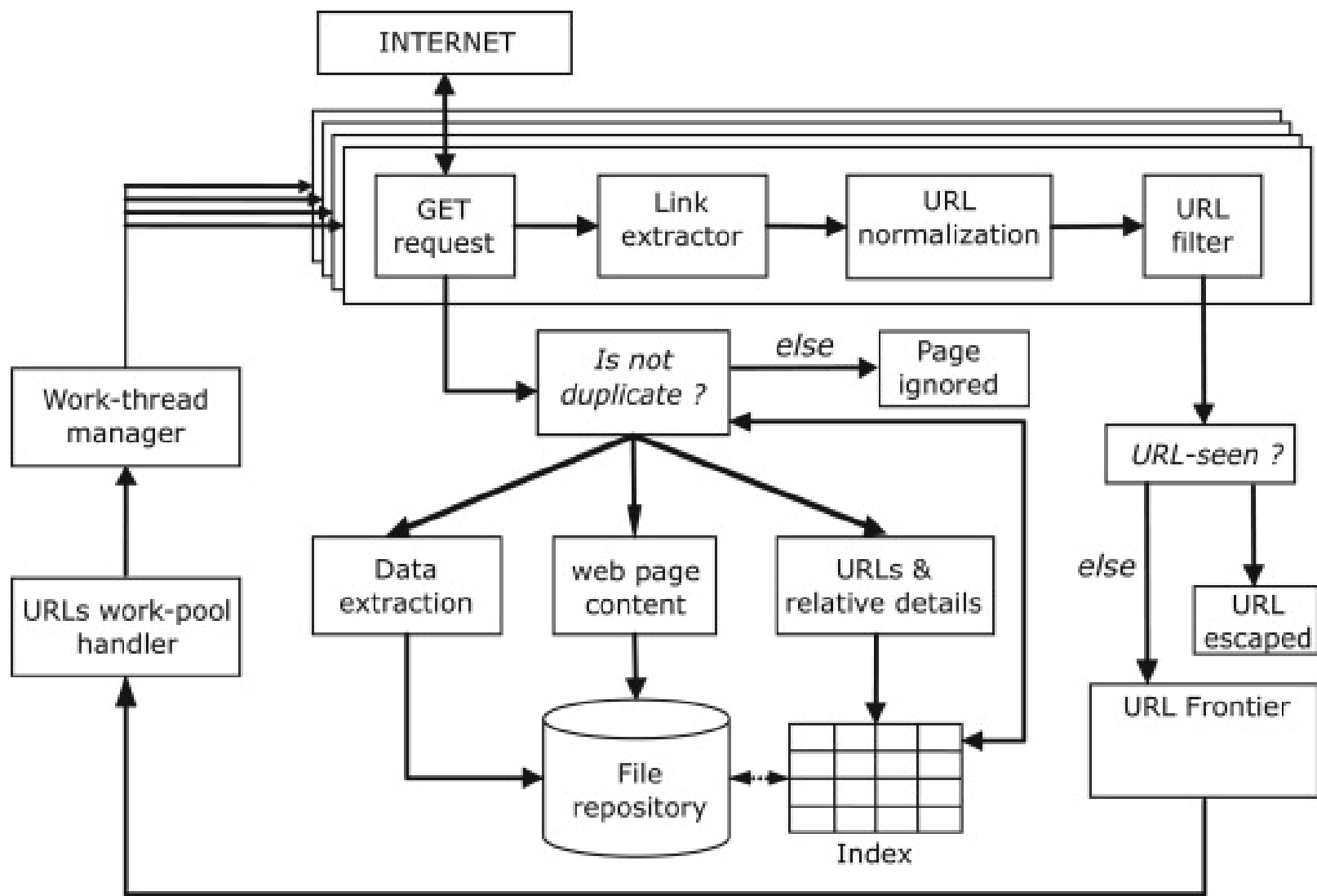
<b>Data scraping</b>	<b>Data Crawling</b>
<b>Involves extracting data from various sources including web</b>	<b>Refers to downloading pages from the web</b>
<b>Can be done at any scale</b>	<b>Mostly done at a large scale</b>
<b>Deduplication is not necessarily a part</b>	<b>Deduplication is an essential part</b>
<b>Needs crawl agent and parser</b>	<b>Needs only crawl agent</b>



# Architecture

Process of **automatically** requesting a web document and **collecting** information from it





# exercises

**[https://validator.w3.org/unicorn/?ucn\\_lang=ko](https://validator.w3.org/unicorn/?ucn_lang=ko)**