

Crawling Framework

Open Source

Apache Nutch

[Apache Nutch](#)

Programming Language - Java

Pros

Highly extensible and Flexible system for web crawling Implements search when combined with open source search platforms like Apache Lucene or Apache Solr Dynamically scalable with Hadoop

Cons

Difficult to setup Poor documentation Some operations take longer, as the size of crawler grows

Heritrix

- [Heritrix](#)

Programming Language - Java

Pros

Excellent user documentation and easy setup Extensible, good performance and decent support for distributed crawls Respects robot.txt

Cons

Not dynamically scalable

- <https://www.scrapehero.com/best-web-crawling-tools-and-frameworks/>

Plugin Backlinks:

From:

<http://moro.kr/> - **Various Ways**

Permanent link:

<http://moro.kr/open/crawling-framework>

Last update: **2021/01/27 01:52**

