

# CSE305 Project: Parallel crawling

Gleb Pogudin, [gleb.pogudin@polytechnique.edu](mailto:gleb.pogudin@polytechnique.edu)

Crawling refers to the process of going to a webpage, following all the links from it, following all the links in the found pages, etc. This process is used, for example, by search engines to index the content of the Internet. Or it can be applied to a particular multi-page website to get its content. The goal of this project will be to crawl Wikipedia (or some other website of interest) to create its “index”, a list of pages.

The resulting code will consist of two major parts:

1. Multithreaded crawler, where each thread takes the next webpage to visit, downloads it, and extracts all the links from it. You may want to use `LIBCURL`<sup>1</sup> for this.
2. A data structure storing already found webpages and allowing the multiple crawling threads to add new elements. For a team of one, the `SETLIST` would be sufficient. For a team of two or more, a concurrent hash-table should be implemented (see, for example, Chapter 13 in the Herlity book<sup>2</sup>; `StripedHashSet` or something more involved it required).

Furthermore, it is expected that you:

1. Analyze the behavior of your implementation depending on the number of threads involved and propose a heuristic to choose such a number.
2. Run the code on Wikipedia or comparable website to generate its index, and analyze which parts of the code takes the largest portions of time.

---

<sup>1</sup><https://curl.se/libcurl/>

<sup>2</sup><https://dl.acm.org/doi/book/10.5555/2385452>