

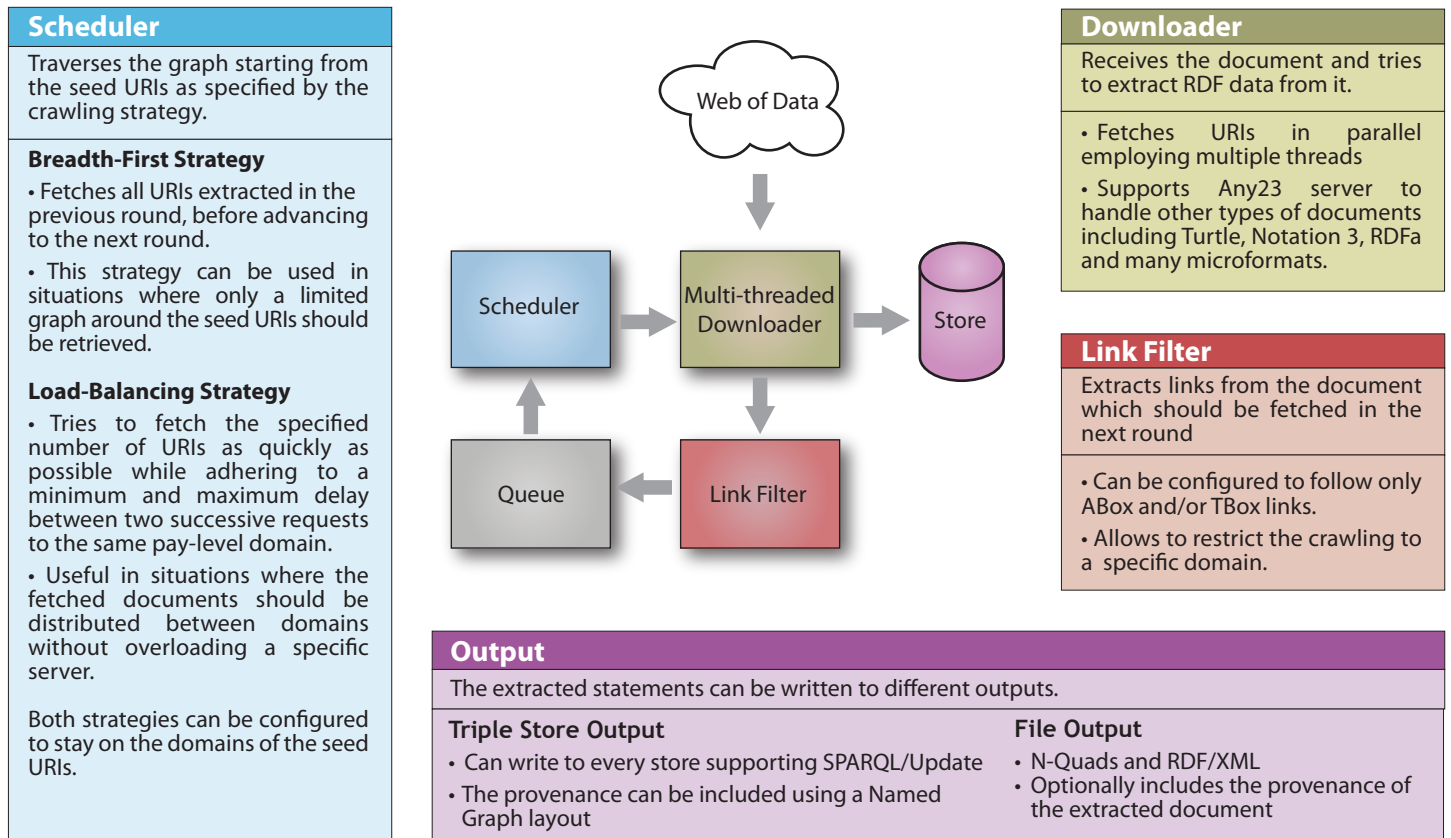
LDSpider - An open-source Crawling Framework for the Web of Linked Data

LDSpider is an extensible Linked Data crawling framework, which enables client applications to traverse and to consume the Web of Linked Data. LDSpider provides an easy-to-use command line interface as well as a Java API which allows applications to configure and control the details of the crawling process.

Features

- LDSpider can process a variety of Web data formats including RDF/XML, Turtle, Notation 3, RDFa and many microformats.
- Crawled data can be stored together with provenance meta-information either in a file or via SPARQL/Update in an RDF store.
- LDSpider offers different crawling strategies, such as breadth-first traversal and load-balancing, for following RDF links between data items.
- Besides of being usable as a command line application, LDSpider also offers a simple API which allows applications to configure and control the details of the crawling process.
- The framework is delivered as a small and compact jar with a minimum of external dependencies.
- LDSpider is high-performing by employing a multi-threaded architecture.
- Can be used to collect small to medium-sized datasets up to tens of millions of triples.

Architecture



Usage Examples

Crawling FOAF profiles

Crawling of interlinked FOAF profiles starting with single seed profile.

round	1	2	3	4	5
profiles	1	10	101	507	6730

Crawling Twitter profiles

Crawling of interlinked Twitter profiles, which expose structured data using RDFa, starting with a single seed profile.

round	1	2	3
profiles	1	38	1160