

RDFaS

an RDFa Spider

Aim: to crawl the web extracting RDFa from pages

What is RDF and RDFa?

The Resource Description Framework (RDF) is a general method for the description or modeling of information that is implemented in web resources. RDF statements are represented as subject-predicate-object triples.

RDFa is an extension to XHTML markup which allows authors to turn existing human-visible text and links into machine-readable data without repeating the content. XHTML attributes are used to represent RDF triples.

The following HTML shows a heading with the content "Ivan Herman", which is the name of a person. However at face value, a computer would not be able to tell that this is a name.

```
<h1>Ivan Herman</h1>
```

An attribute can be set for the heading element, giving the content a property, in this case "foaf:name":

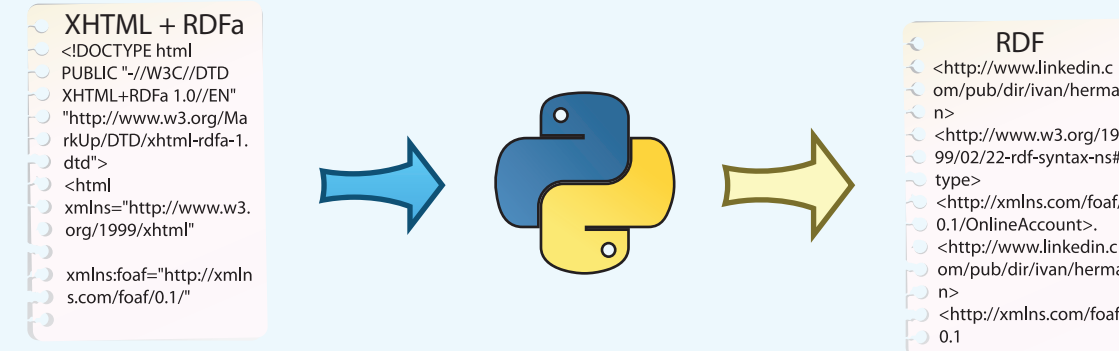
```
<h1 property="foaf:name">Ivan Herman</h1>
```

The prefix before the colon is the namespace and the suffix is the property. The foaf is linked to the foaf ontology, which is used to represent relationships between people and the name property indicates that the content represents the name of something.

Parsing Web Pages to extract RDF

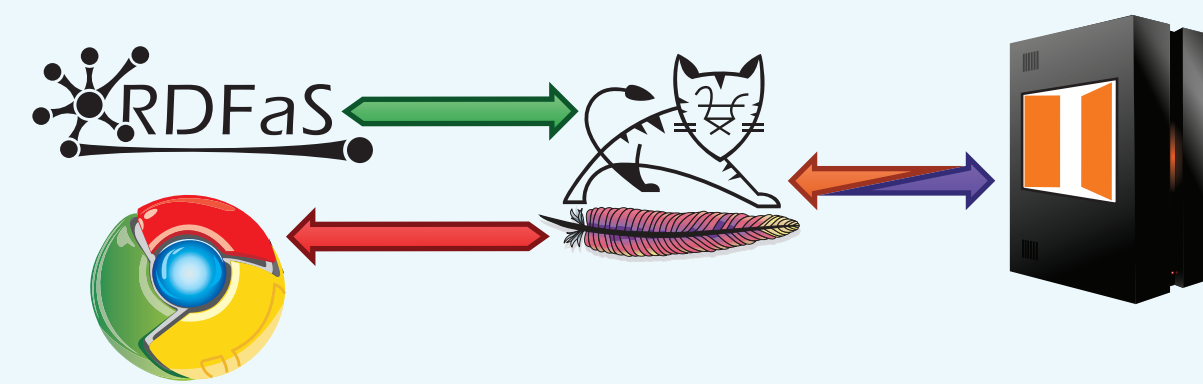
RDFa information when included in web pages is very bulky as not all of the extra XHTML from the page is needed. The ideal representation is as RDF triples which extracts the subject, predicate and object from the XHTML. Triple hierarchies can be inferred where the object of one triple matches the subject of another.

W3C's Distiller (also known as pyRDFa) is used for parsing.



Storing RDF Data - Triple Stores

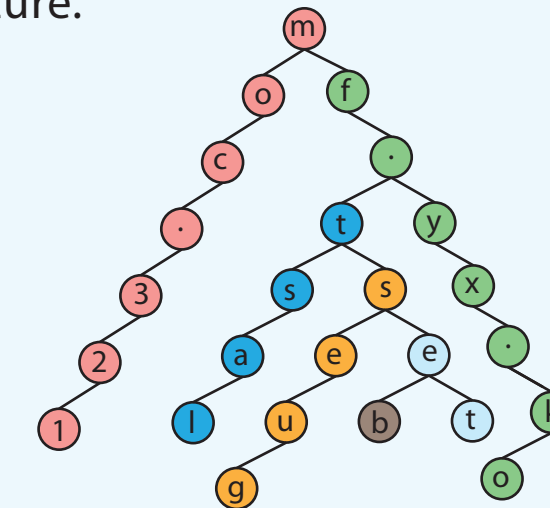
RDF data is conventionally stored in type of database system called a Triple Store. Some triple stores have a native storage backend tailored at storing RDF and performing reasoning, others use existing database engines like MySQL. Sesame, the Triple Store used for this project, provides access using Tomcat for adding and querying data over HTTP.



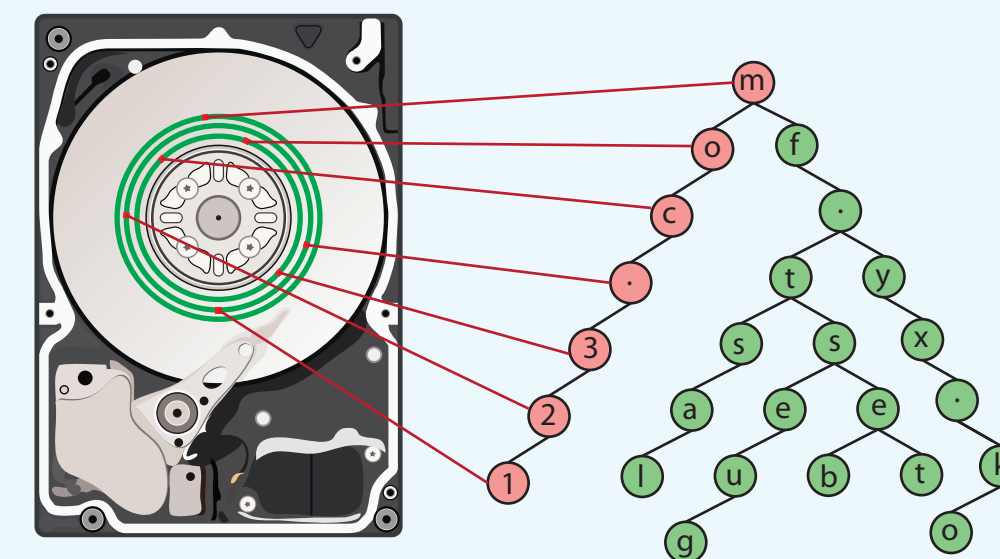
The URL Cache - Trie Implementation

The URL Cache analyzes URLs from the Crawler to determine whether or not they have been seen before and therefore whether or not they have been Crawled before or queued to be Crawled.

The cache is a large index of URLs, stored as a Trie (a form of tree structure). Duplicate prefixes are eliminated due to the hierarchical structure.



Unfortunately the Trie suffers from a performance flaw. Trie nodes are added sequentially as they occur from requests but often read in the order of the hierarchy meaning lots of drive head jumps (or seeks) will be required adding varying levels of overhead. On average a seek takes 10ms meaning a 100 character URL will take 1000ms (1 second) to check.



The URL Queue - Paged Queue

The URL Queue is a crucial part of the system, it tells the Crawler Threads the next URL to crawl but also determines the order in which URLs are crawled to avoid over-crawling any particular site.

URLs could come in to the queue in any order, it is likely many from the same site will be discovered in a short space of time.

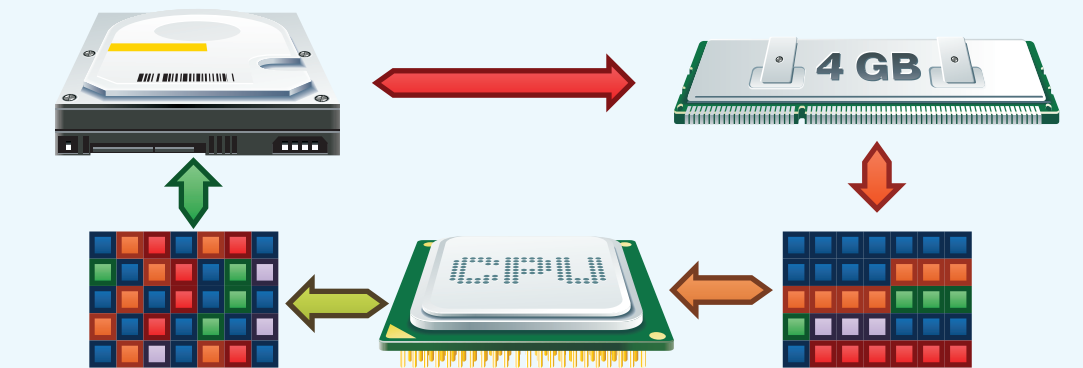


As URLs enter the queue they are stored in a Temporary Page until the page becomes full. Once it is full they are rearranged into a polite, non-sequential order (based on the domain) and saved in a Queue Page. Any domains that occur too frequently are put back in a Temporary Page to be thinned out during the next processing cycle.

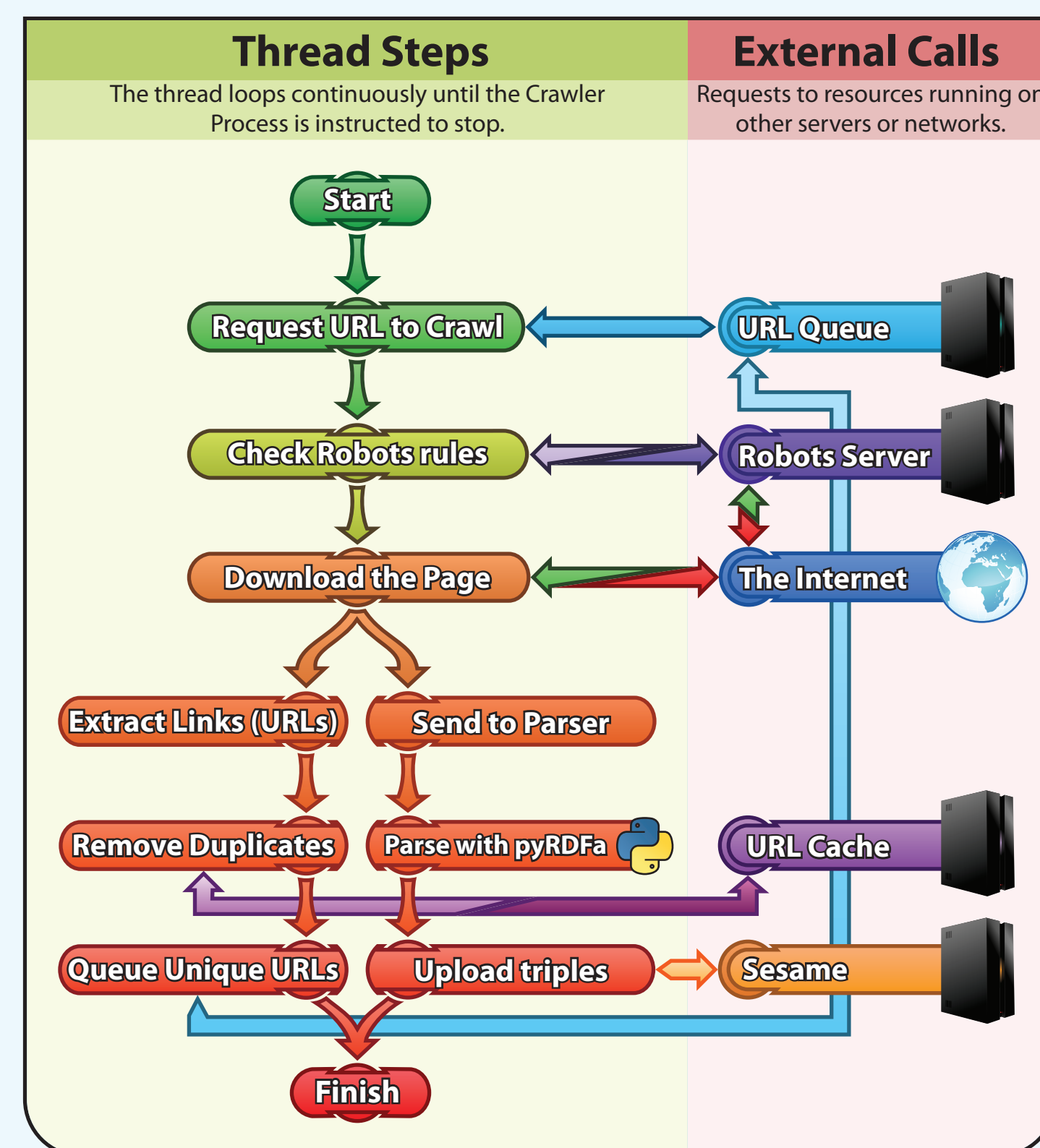


The ideal diversity of a Queue Page (the limit on the most frequent domains) is automatically tuned based on recent crawling and speed statistics.

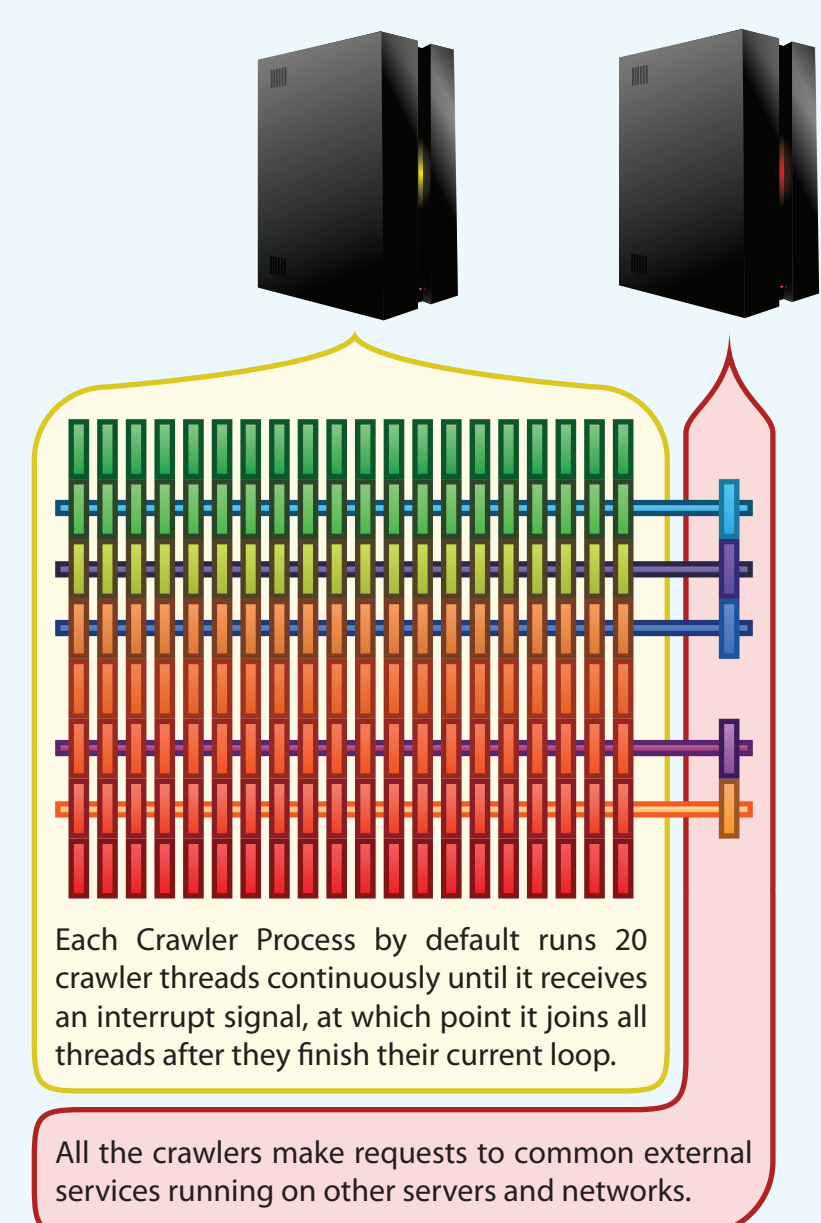
The queue is broken down in to pages so that data can be read and written in sequential chunks and each page is always small enough so that it can be processed in memory.



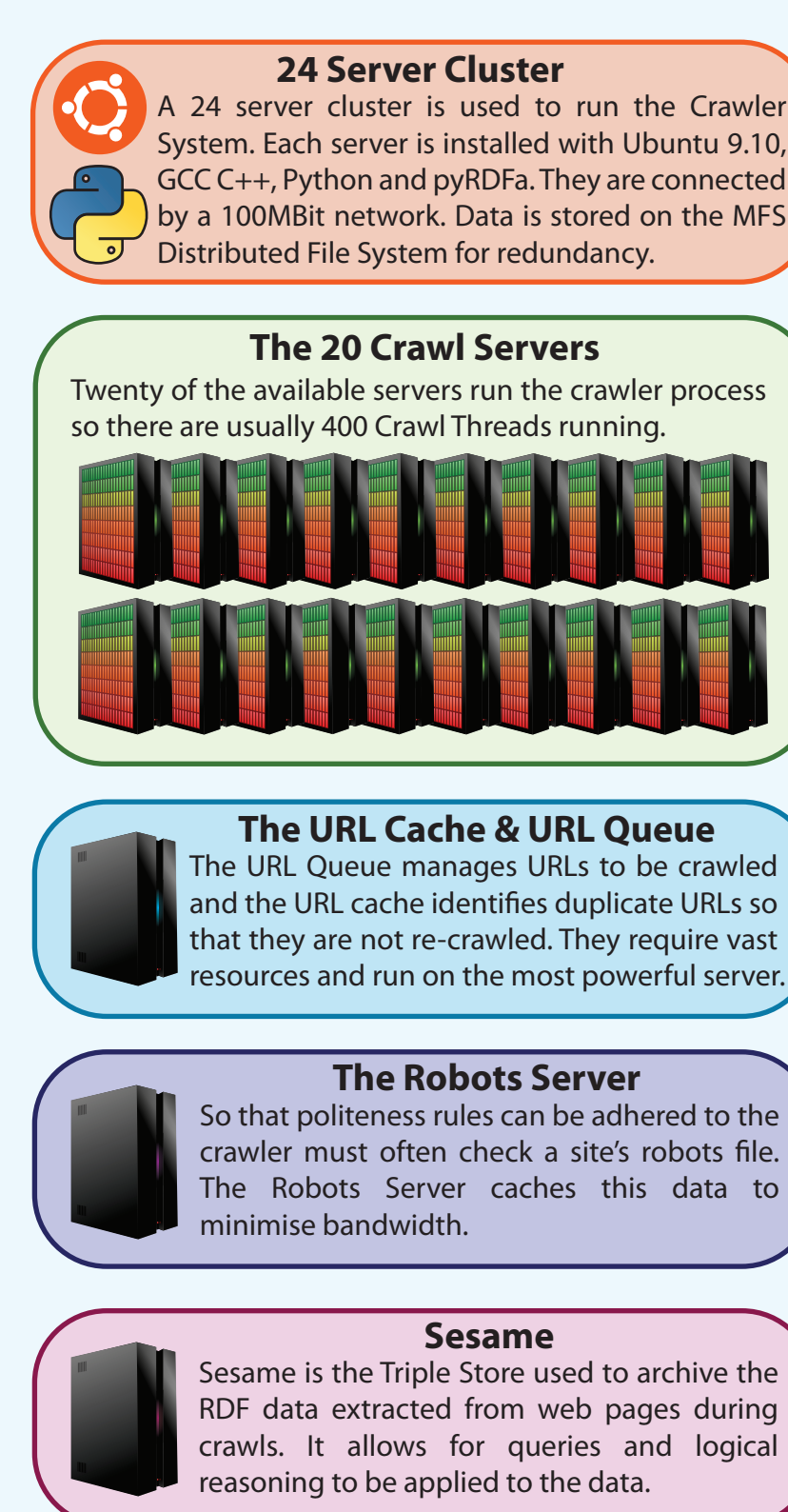
A Crawler Thread



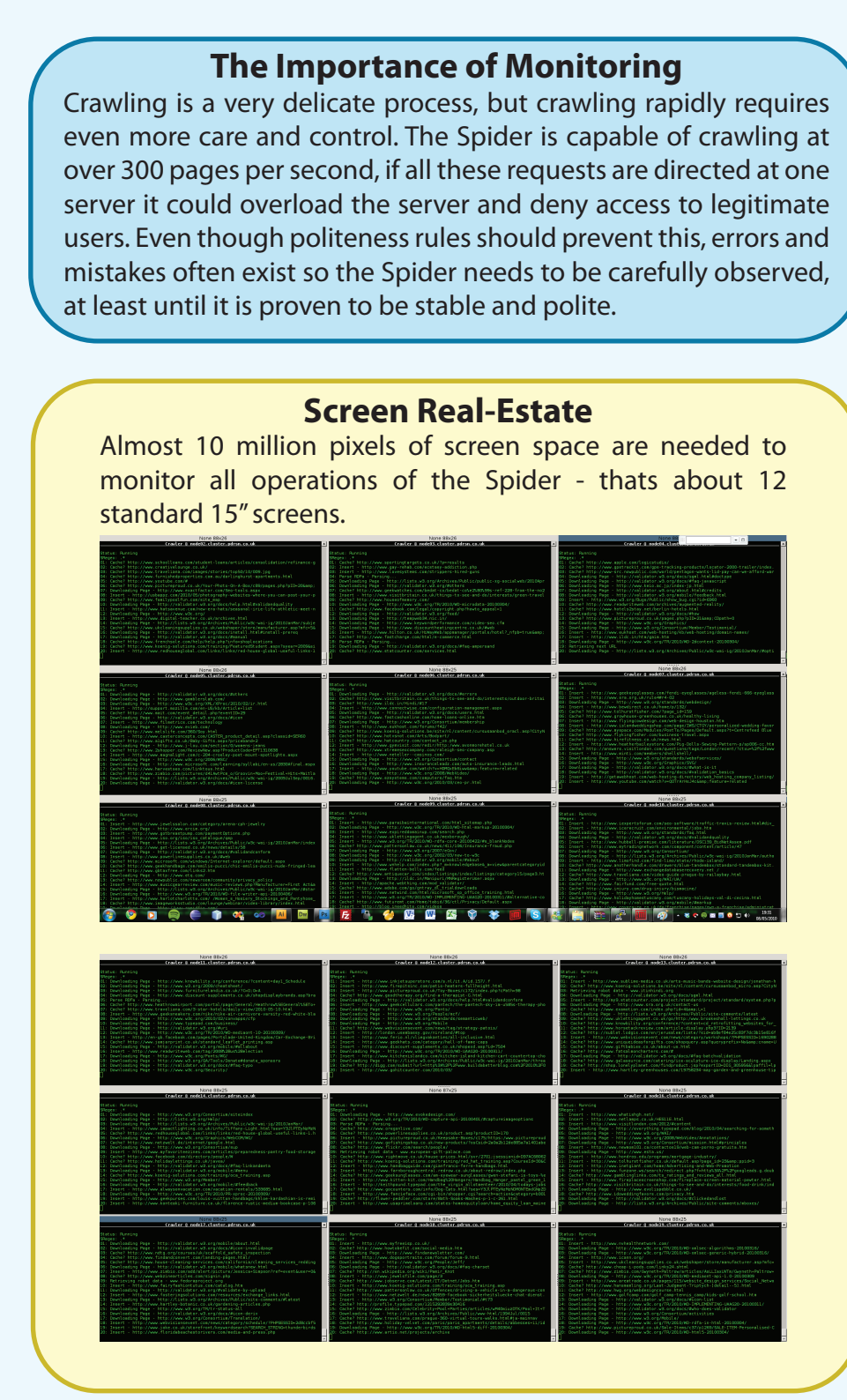
The Crawler Process



The Crawl Cluster



Monitoring the Crawl



Querying the Triple Store

