

III. ROBOTS AND SITEMAPS

Key element of the crawler is also reading the 'robots.txt' file and sitemap. Based on the robots file, we know whether we can visit a certain link or whether it is not allowed for a crawler to visit that link. When we read a sitemap, we add all the URLs mentioned in it, into the frontier.

IV. DATABASE

The database was build inside a Docker container. We used the same schema, with one minor change, which is adding a column 'lsh_hash' into the page table. The column stores the value that LSH hash produced for the content of that website. Access to database is made by each worker when he visits a new page. The database can handle multiple workers making read and write operations by proper locking of the operations.

For viewing the results and current status in the database, we also used Adminer application, which also runs as a docker container and is then accessible on localhost.

V. LINK

VI. IMAGE AND FILES EXTRACTING

@JAKA — commnet on the issue with absolute path, very interesting:D

We stored images and files into the file system, since it gives a nicer overview of the gathered information. This was made possible by creating a folder for each new site in which we stored the information. We then stored the location of the file/image into the provided data field in the database.

VII. LSH DEDUPLICATION

We implemented LSH algorithm, which is being used to compare hashed content of different websites.

How it works

For the set of words, we hash it on, we used top 1000 Slovenian words. Or maybe triples of characters?

VIII. RESULTS - 2 SEED URLs

For the first run of the crawler, we used 2 seed URLs, which were "http://evem.gov.si" and "http://e-prostor.gov.si". We set the maximum crawl depth to None, which means it will run untill there are no more links produced. The maximum number of workers was **NUM OF WORKERS WE USED**.

The crawling process took XXXX minutes in total and retrieved XXX pages,links, images...

For the sites that are given in the instructions' seed list and also for the whole crawldb together (for both separately) report general statistics of crawldb (number of sites, number of web pages, number of duplicates, number of binary documents by type, number of images, average number of images per web page, ...). Visualize links and include images into the report. If the network is too big, take only a portion of it. Use visualization libraries such as D3js, visjs, sigma.js or gephi.

IX. RESULTS - 9 SEED URLs

The second run consisted of 9 seed urls. For the extra 5 urls of our choice, we have choosen "http://www.mz.gov.si/", "http://www.mnz.gov.si/", "http://www.up.gov.si/", "http://www.ti.gov.si/" and "http://www.mf.gov.si/".

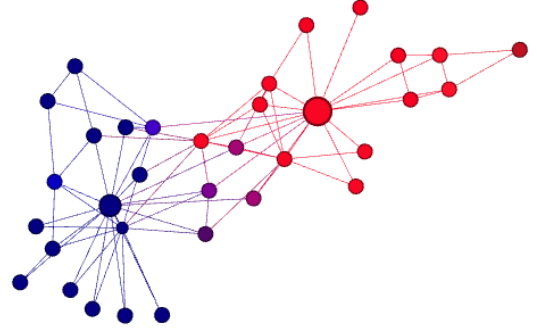


Figure 2. Example of an image.

Table I
EXAMPLE OF A TABLE.

Richard Karstark	0.013661
Jon Arryn	0.012869
Joyeuse Erenford	0.012869
Trystane Martell	0.010761
Willen Lannister	0.010684
Martyn Lannister	0.010684
Robb Stark	0.010304
Joffrey Baratheon	0.009875
Master Caleotte	0.009495
Lysa Arryn	0.009383

REFERENCES

- [1] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. Cambridge University Press, 2008.
- [2] A. Gionis, P. Indyk, and R. Motwani, "Similarity search in high dimensions via hashing," in *Proceedings of the 25th International Conference on Very Large Data Bases*, ser. VLDB '99. Morgan Kaufmann Publishers Inc., 1999, pp. 518–529.