SMART 2019/1083
SMART 2019/1083 – Task 6 Assessments, trials and evaluations

# Multilingualism scoring tool

## Rinalds Vīksna, 25.11.2021.

# MOTIVATION

Making European websites more multilingual is one of the targets of the Connecting Europe Facility Automated Translation (CEF AT). In order to monitor this goal, alongside other possible solutions, CEF AT needs a methodology and a tool to assess the degree of multilingualism of a web site is needed. The tool should be fully or semi-automatic and easy to use. This document covers the proposed acceptance criteria for and assessment under SMART 2019 1083 Task 6.

The emphasis is on the openness of the tool, as public administration and small business need to run it by themselves to evaluate the degree in which their website is multilingual. The tool can have proprietary components, but the main functionality of the tool should not depend on them.
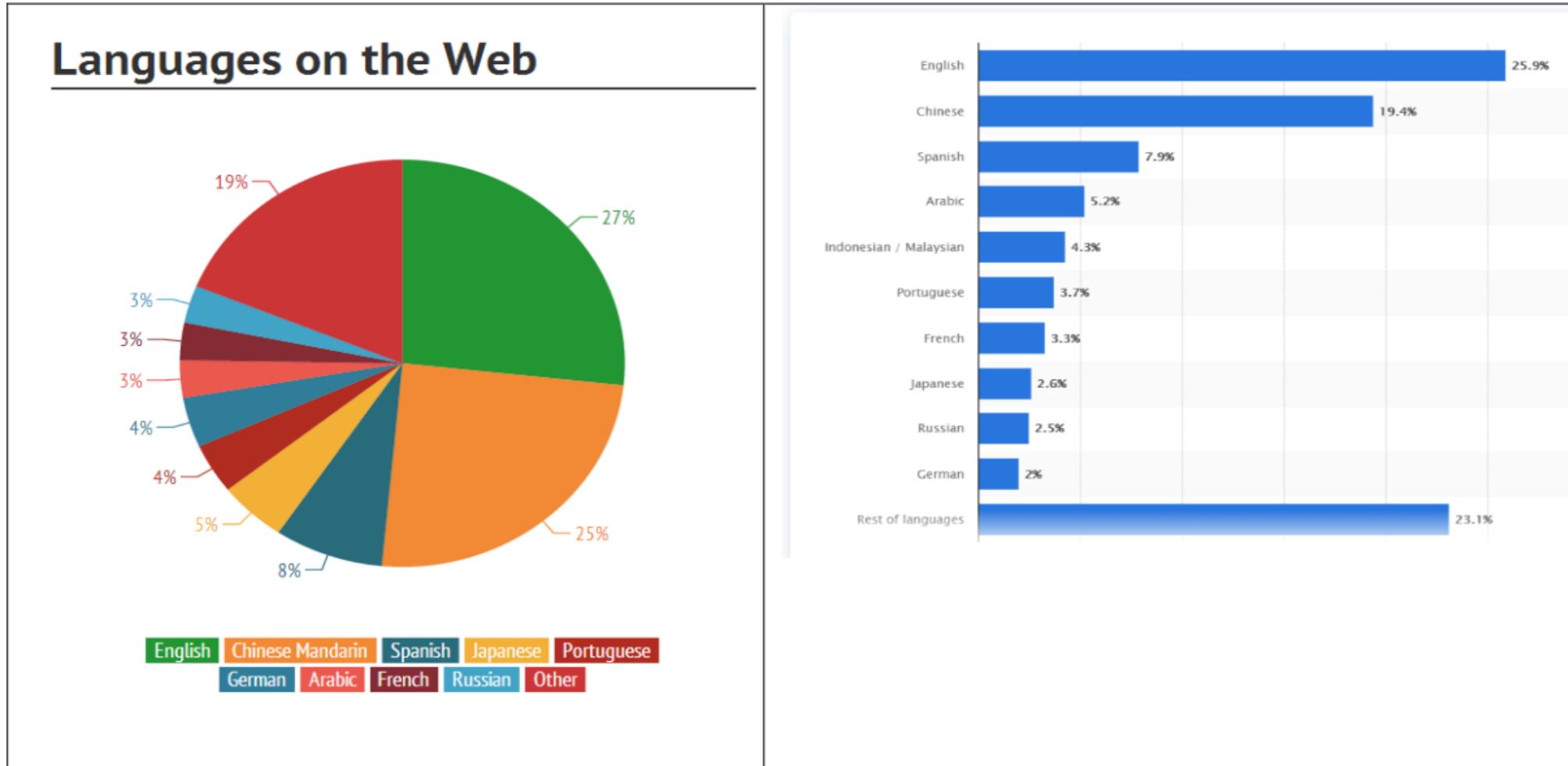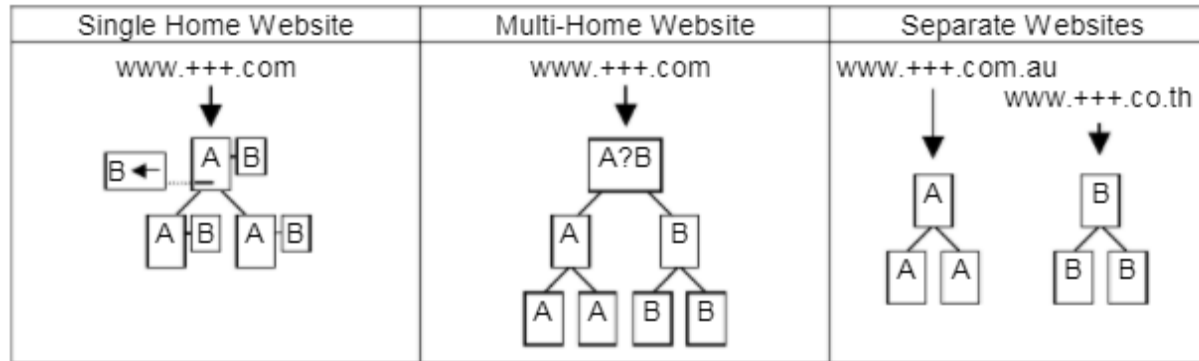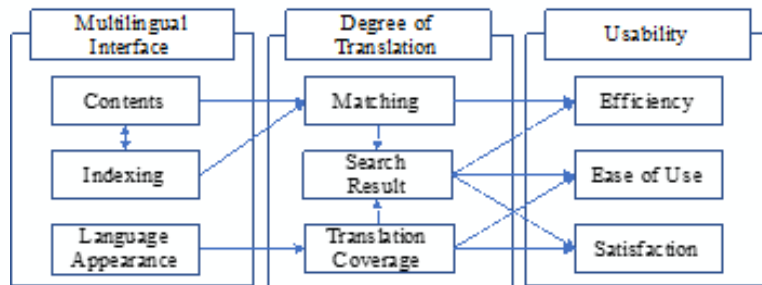
Figure 1. On left - Language of the Web (*Web Language Diversity in Numbers | Language & Technology | Globalme*); on right - Most common languages used on the internet 2020 (*• Internet: most common languages online 2020 | Statista*)

# MULTILINGUALITY

- Multilingual - "(of people or groups) able to use more than two languages for communication, or (of a thing) written or spoken in more than two different languages"

- The Bilingualism Versus Multilingualism Dimension
  - For purposes of this task, we will define "Multilingualism" as written or spoken in two or more different languages. Bilingualism or trilingualism are instances of multilingualism.

- W3C: "a "multilingual" web site refers to a web site that uses more than one language."
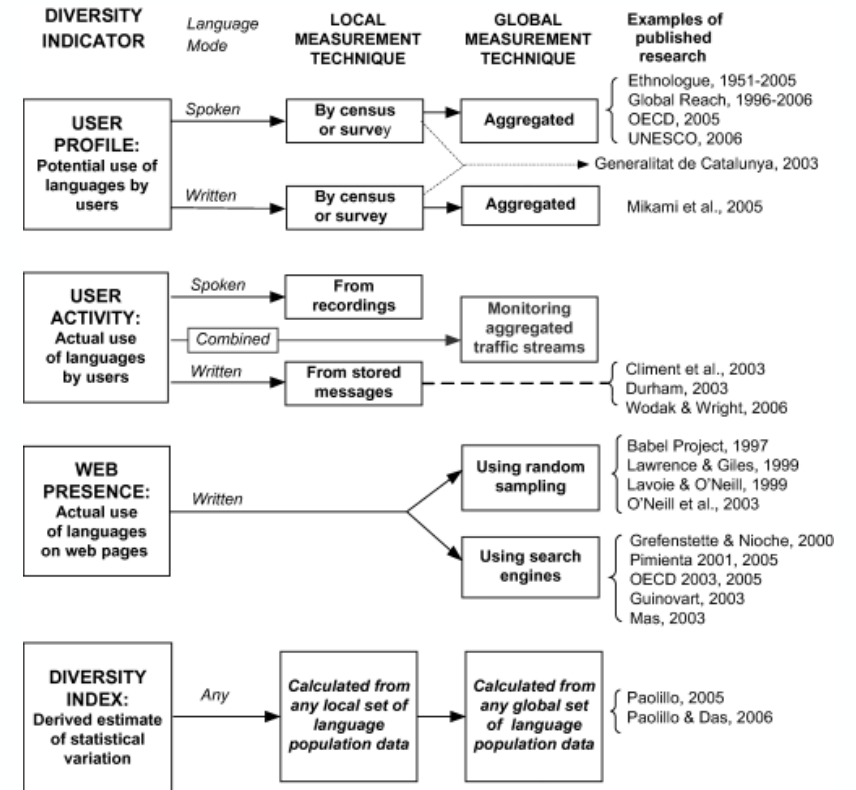
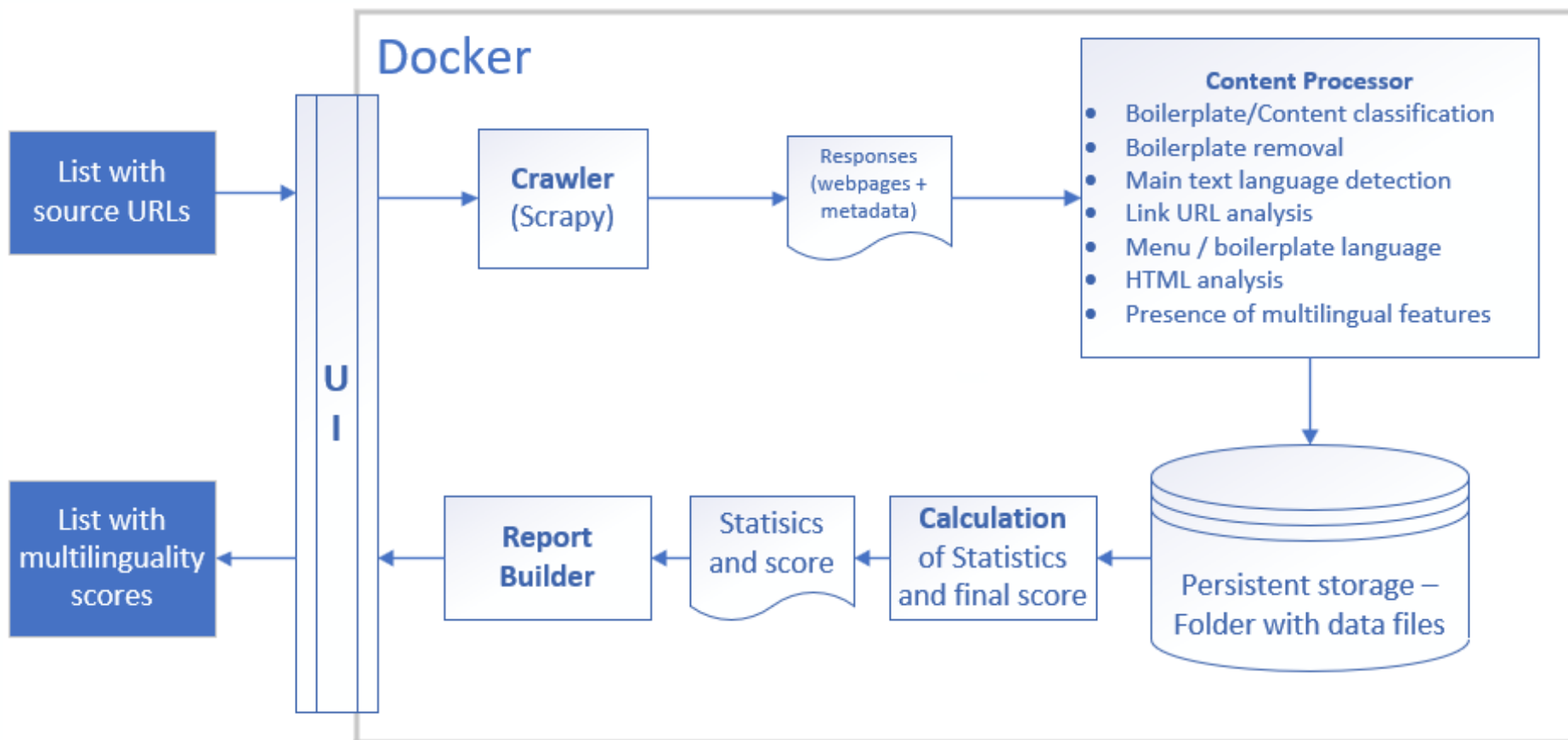Three Broad Types of Multilingual Websites



*Multilingual Evaluation Guideleni model proposed by Lee and Choi (2019)*



*A taxonomy of different methodologies used to estimate language diversity on the Internet (Gerrand, 2007)*

https://www.researchgate.net/publication/2939572_Multilingual_Website_Usability_Cultural_Context

# MULTILINGUALITY MEASURE

- ✓ Language coverage
- ✓ Language balance

- ▪ Linguistic quality
- ▪ Technical quality – i18n attributes and others
- ▪ Content parallelism and Multilingual functionality (Navigation)

# Architecture

- Crawler – Scrapy

- Boilerplate removal – jusText

- Language detection – LangDetect

- Scoring – 2 formulas:
  - Normalized language balance = sum(Share1, Share2, …, ShareN)/N
  - Lieberson's diversity index $LDI = 1 - \sum P_i^2$
    - where Pi represents the share of i-th language

Multilingualism Scoring Tool is available for download either as a code and Docker container on the github:
https://github.com/tilde-nlp/Multilingualism-scoring-tool

- Crawled bottom 230 websites from https://moz.com/top500 with depth 1 in two days

- Found 34 empty results:
  - Redirects 14
  - No content 6
  - Bad address 10
  - Javascript 1
  - Restrictive robots.txt 2
  - Forbidden 1
  - OK 196



230 - All websites

34 Empty websites

Redirects    No content    Bad address

Pure javascript    robots.txt    Forbidden

OK

# Observations while using tool

- Number of prepared requests in each depth (before filtering)
  - 'request_depth_count/0': 14,
  - 'request_depth_count/1': 1000,
  - 'request_depth_count/2': 72004,
  - 'request_depth_count/3': 1259091,
- Stopped after ~110K requests
  - 'exception_count': 17784 detailed breakdown:
    'ValueError':2, 'IgnoreRequest':12784, 'CancelledError':1, 'ConnectionRefusedError':4929, 'DNSLookupError':54, 'TimeoutError':11, 'ResponseNeverReceived': 3,
  - Responses 92927:
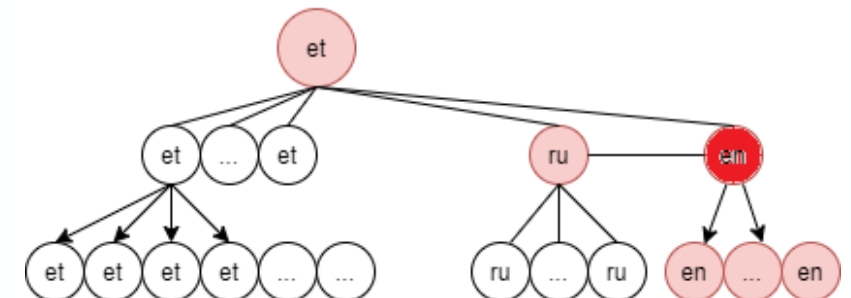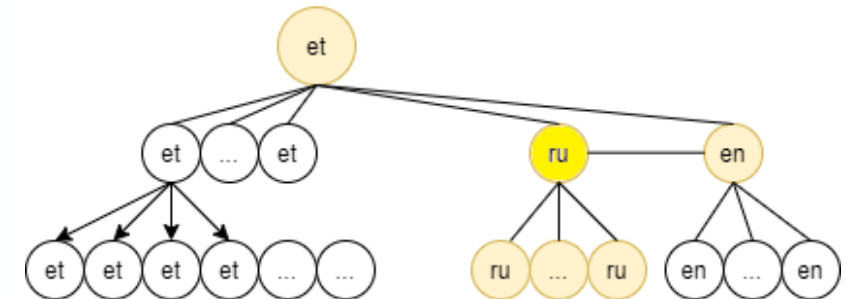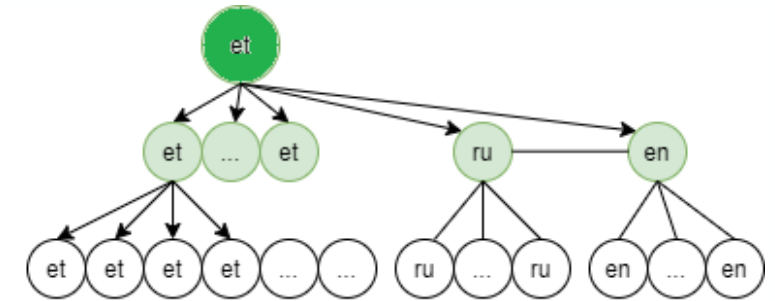    200 – 73096; 301 – 7534; 302 – 6338; 307 – 598; 404 – 228; 500 – 77, others <15 each

# Observations while using tool

https://president.ee/et/
- 1 Hops: Pages: 45; w/o 8, et 35, ru 1, en 1
- 2 Hops: Pages: 2892; w/o 1274, et 1497, ru 39, en 72, fi 2, lv 1, uk 1, lt 1, el 1, ….
- 3 Hops: Pages: 8272; w/o 3670, et 3164, ru 463, en 910, fi 10, lv 4, uk 5, lt 3, …
- 4 Hops: Pages: 13179; w/o 6566, et 3938, ru 857, en 1694, fi 24, lv 6, uk 8, lt 9, …

https://president.ee/ru/index.html
- 1 Hops: Pages: 37; w/o 7, ru 27, en 1, et 2
- 2 Hops: Pages: 673; w/o 158, ru 434, en 38, et 42, uk 1
- 3 Hops: Pages: 5373; w/o 2147, et 1540, ru 824, en 812, uk 5, fi 5, lv 3, lt 2, …
- 4 Hops: Pages: 12031; w/o 5980, et 3367, ru 898, ne 1670, uk 8, fi 19, lv 6, …
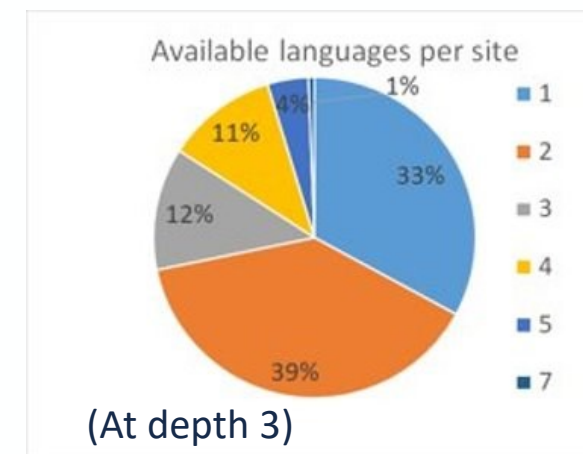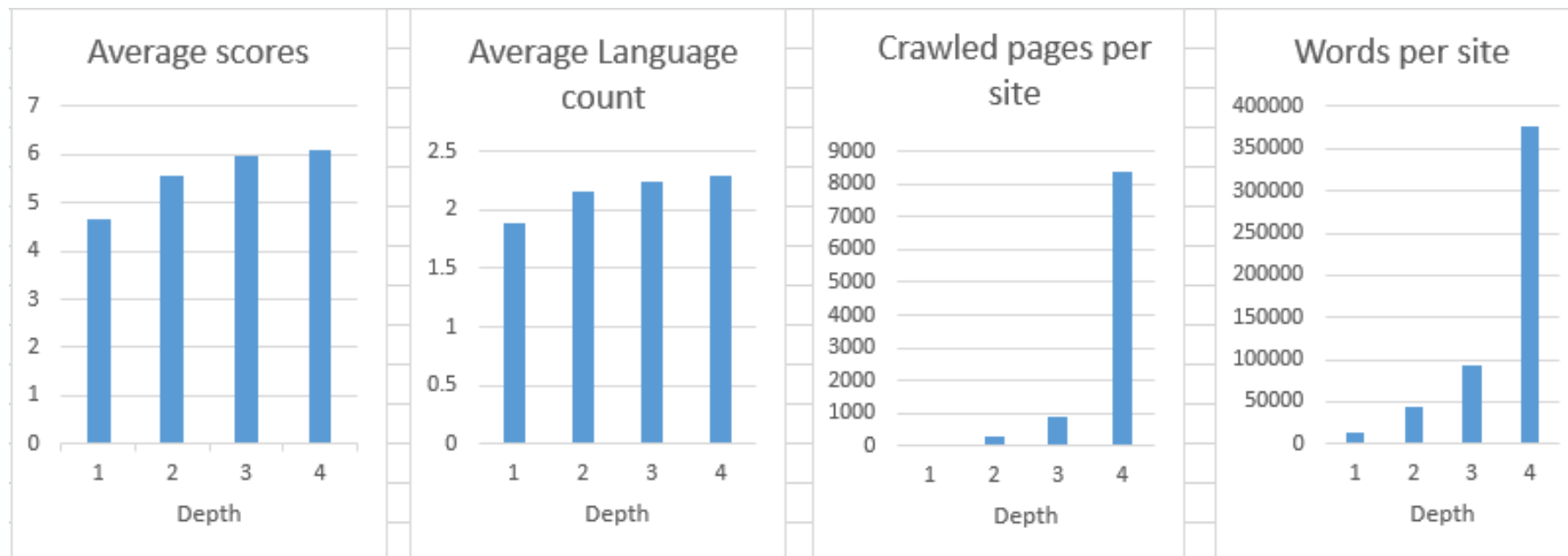
https://president.ee/en/index.html
- 1 Hops: Pages: 39; w/o 6, ru 1, en 31, et 1
- 2 Hops: Pages: 1133; w/o 293, ru 28, en 733, et 43, uk 2, fi 3, de 6, sv 2, lv 2, …
- 3 Hops: Pages: 6162; w/o 2445, et 1625, ru 456, en 1537, fi 14, de 15, sv 4, lv 5, …
- 4 Hops: Pages: 12471; w/o 6149, et 3421, en 1881, ru 879, fi 28, de 21, sv 5, lv 6, …

# RESULTS

| URLs | Crawling depth | Crawl time | Memory usage max | Average score | Average coverage of EU languages | Average number of pages/site | Average number of words/site |
|------|------|------|------|------|------|------|------|
| 198 | 1 | 0:50h | 160MB | 4.65 | 1.89 | 33 | 14516 |
| 198 | 2 | 12:54h | 1030MB | 5.57 | 2.15 | 262 | 44592 |
| 198 | 3 | 48h | 1238MB | 5.97 | 2.25 | 885 | 93242 |
| 198 | 4 | >250h | 6672MB | 6.08 | 2.29 | 8374 | 376787 |
| 600 | 2 | 52h | 1023MB | 5.56 | 2.02 | 227 | 53663 |



(At depth 3)

# CONCLUSIONS

- The tool can crawl a large number of given websites and give some results almost immediately, updating scores as more pages are crawled;

- Time necessary for "polite" crawling increases quickly as we try to exhaustively crawl a website;

- "Multilinguality score" we are currently using represents coverage of translated content in respect to the most represented language on the website, as well as the number of languages the content is translated in;

- European websites currently are not very multilingual – on average content is presented only in 2-3 languages;

# THANK YOU FOR YOUR ATTENTION!