# Harvesting Tables in the Wild

Jack Henschel, Rohit Raj, Eelis Kostiainen
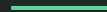
12/02/2021

# The Knowledge on the Internet

- World Wide Web is vast and has lots of "data"

- Mostly focused on human consumption

- Turning this "data" into knowledge is a semantic challenge!

- *TSLA, 2020-04-06, 420.2*

Our Goal:

**Build a scalable and reliable pipeline to extract HTML tables from web pages**

# Related Work

# Literature Review

- Comprehensive review by S. Zhang and K. Balog[1] about:
    a. Table Extraction
    b. Table Interpretation
    c. Table Search
    d. Table Question Answering
    e. Knowledge Base Augmentation
    f. Table Augmentation

[1] S. Zhang, K. Balog: "Web Table Extraction, Retrieval, and Augmentation: A Survey" ACM Transactions on Intelligent Systems and Technologies (2020).
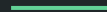
# Using reference knowledge bases

- Work by Muñoz et. al.[2] to match tables from Wikipedia
  - Extracts data from wikipedia as RDF
  - Uses the knowledge gained to both augment the knowledge DB

- D. Ritze et. al.[3] proposed using DBpedia
  - Evaluates performance of HTML to a KB matching system
  - T2K Match: an iterative matching method combining schema and instance matching for matching common HTML tables against cross-domain knowledge bases

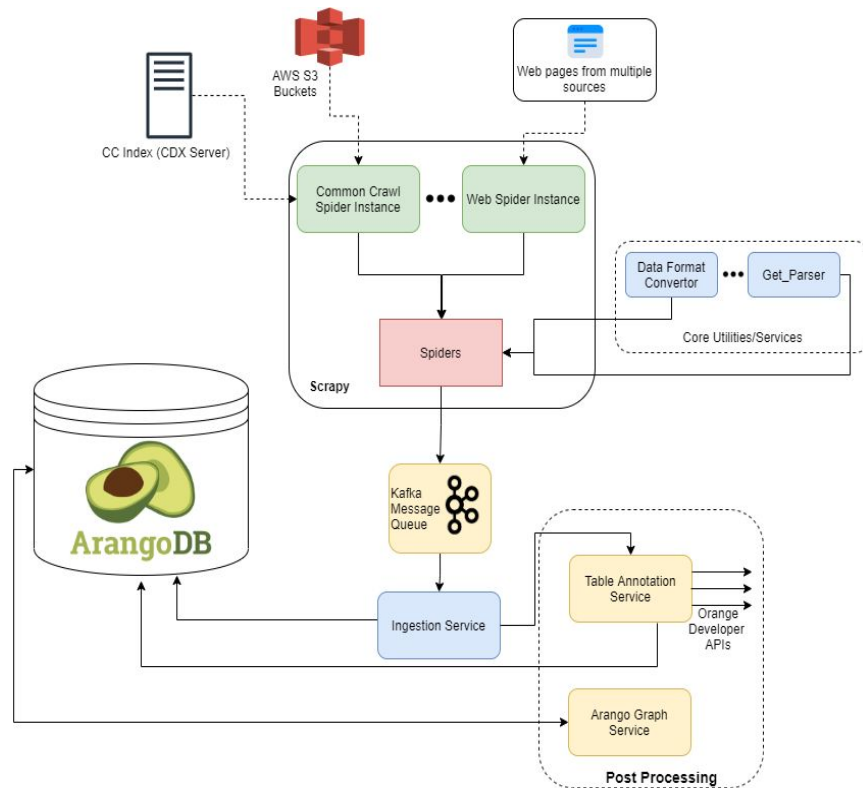[2] E. Muñoz, A. Hogan, A. Mileo: "Using linked data to mine RDF from wikipedia's tables" WSDM (2014).
[3] O. Lehmberg, D. Ritze, R. Meusel, C. Bizer: "A Large Public Corpus of Web Tables containing Time and Context Metadata", WWW 2016.
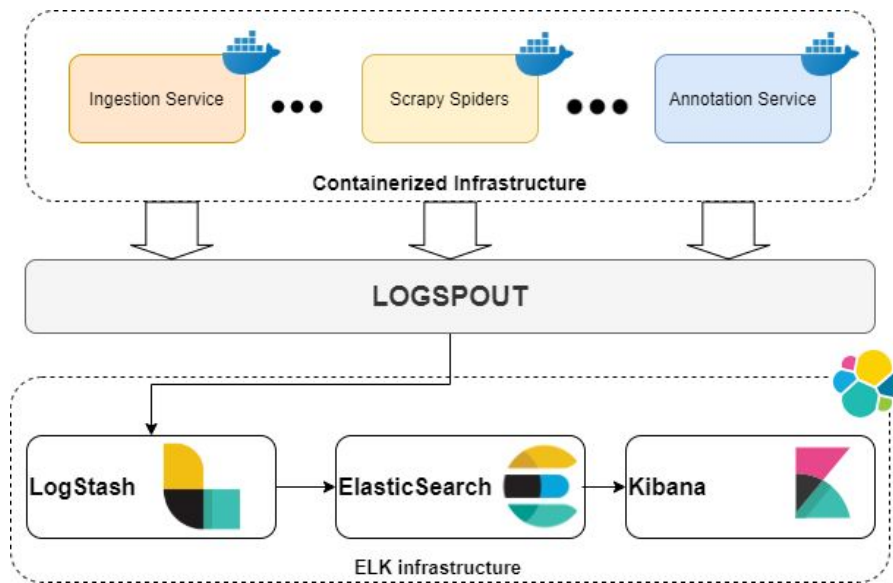
# System requirements & design

# Crawling Pipeline

- High-level requirements
  - Continuous crawling of data from different sources
  - Single-node
- Data sources
  - World Wide Web
  - CommonCrawl (CC) web corpus
- Technologies
  - Crawlers: long-running tasks
  - Storage: represent semantic relationships
  - Ingestion: post-processing scalable to high volumes
  - Monitoring: provide insights to crawling pipeline



8

# Monitoring Infrastructure

- **Problem**: pipeline insights

  unknown

- **Solution**: continuous monitoring

- Elasticsearch, Logstash, Kibana

  (ELK stack)

  - Logspout for docker log redirection

  - Explicit redirection

# Implementation & Milestones

# Getting started

Minimal implementation with Scrapy framework:

1.  Download HTML from given URL
2.  Extract *<table>* elements
3.  Store them in a JSON file

# Advanced table parsing

- HTML tables are used for many different purposes

  ➜ need to **select the relevant ones**

- HTML may be **erroneous**

- Some websites have specific patterns for their tables

  ➜ special focus on Wikipedia's ***wikitables***

# Data format

- Retain and store as much information about the web page as possible

- Extensive study about data formats other people have used

- **Final choice**: JSON schema based on DWTC format with additional fields

- Validation through JSON Schema

```
1 ▾ {
2     "hasHeader": true,
3     "pageTitle": "COVID-19 situation update for the EU
4     "url": "https://www.ecdc.europa.eu/en/cases-2019-r
5     "headerPosition": "FIRST_ROW",
6     "tableType": "RELATION",
7 ▾   "relation": [
8 ▾     [
9         "EU/EEA",
10        "Sum of Cases",
11        "Sum of Deaths",
12        "14-day case notification rate per 100 000 inh
13        "14-day death notification rate per 1 000 000
14        "Reporting period YYYY-WW"
15      ],
16 ▾    [
17        "France",
18        "3053617",
19        "73049",
20        "403.45",
21        "79.07",
22        "2021-02 and 2021-03"
23      ],
24 ▾    [
25        "Spain",
26        "2593382",
27        "56208",
28        "1026.05",
29        "83.79",
30        "2021-02 and 2021-03"
31      ],
32 ▾    [
33        "Italy",
34        "2466813",
35        "85461",
36        "315 31"
```

13

# Common Crawl

- Extremely large corpus of web pages (12 years)

- Publicly available index: Data hosted on Amazon S3

- Additional data source for our pipeline

- Data saved in **WARC format**

  - Search the index server

  - Download the WARC file from S3 bucket

# Crawling Strategy

- *Which websites should we visit? Where should we start?*

- Crawler follows all links on a web page

  - Two tier regex approach

- **Whitelisted** and **blacklisted domains**

- Custom selection of **Seed URLs**

- Based on Alexa ranking and CC Top 1000

# Visualization

Visited web pages (and extracted tables) form a relationship
➜ **Graph**

ArangoDB natively supports interacting with graph structures

**Vertices**: *visited_pages* and *parsed_tables*

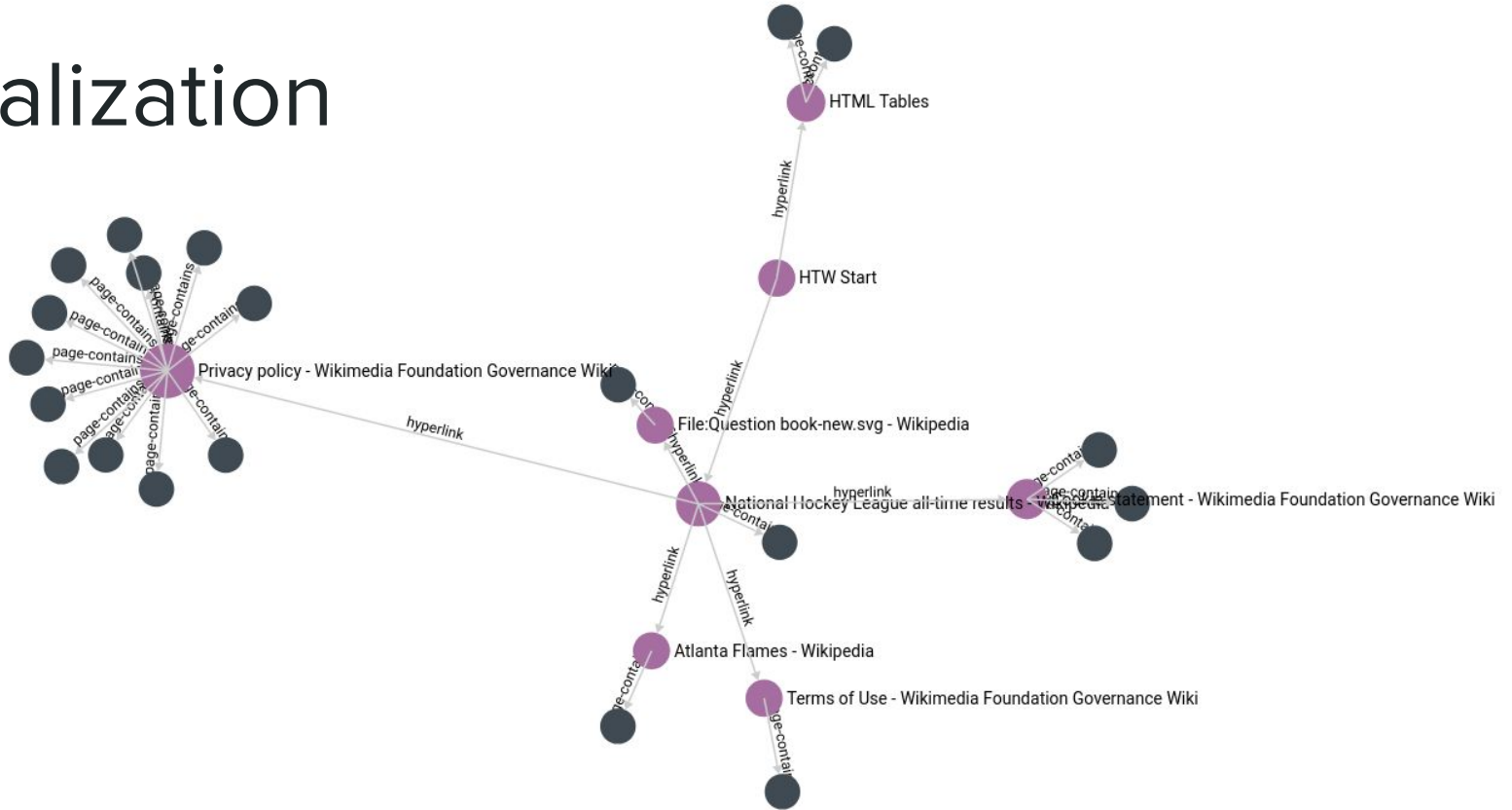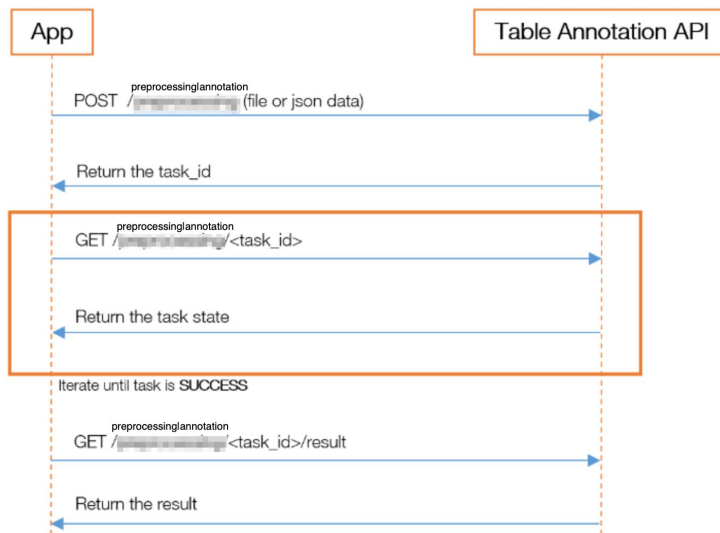**Edges**: *hyperlink* and *page_contains*

# Visualization

# Table Processing

- Goal: extracting meaningful relationships from crawled data

- Solution: integrate pipeline with Orange Table Annotation API [4]

- Steps:
  - Preprocessing
  - Annotation



[4] Orange, *Table Annotation: Semantic Annotation Toolkit for Tabular Data*, Retrieved 10.02.2021, From:
https://developer.orange.com/apis/table-annotation

# Results
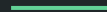
# Partial Graph of final collection

# Results - Some Stats

- Final run during last week of January

- 3469 tables collected

- 2002 websites visited

- 135 unique domains visited

- Low count due to a constrained whitelist strategy and

  suboptimal Seed URL selection

# Results - Qualitative

- Built a robust crawling pipeline from the ground up

- Used modern technologies

- Based on relevant and recent scientific literature

- Followed software development best practices to keep code extensible and maintainable

# Future Work

# Pipeline Improvements

- Refresh previously crawled data

- Parse pages with client-side rendering

- Improvements to seed-pages and whitelisted domains

- Optimizations in extraction and pre-processing
  algorithms

# Table Annotation

- Integrating pipeline with annotation service once stable release available
    - Current beta release too buggy for a stable integration
    - `{'code': 1, 'message': 'Internal error', 'description': "apigee - CatchCommonErrors exception: TypeError: Can't use 'in' on a non-object."}`

- Provide more valuable data insights

# Thanks for your attention!
## Any questions?

*Full report available at:* https://**u9k.de/htw-report**