# 52north
## exploring horizons

# WEBCRAWLER FOR HYDROLOGICAL DATA

## ECMWF Summer of Weather Code

Norwin Roosen @ 52° North
ESoWC Webinars, 20.09.2018

# THE PROJECT

Web Crawler for hydrological data

- goal: *"Development of a tool to search the web systematically, identifying data sources for observed environmental data"*

solution:

- multilingual search engine, but with thematic focus

- classification of web pages (contains/links to data?)

- data set metadata extraction

# CHALLENGES: UNSTRUCTURED DATA

- only some data on the web is listed in catalogs

    → easily accessible, not our focus

- data without machine readable annotation hard to discover

    – links to data without proper annotation

    – embedded into HTML pages

    – access on request only

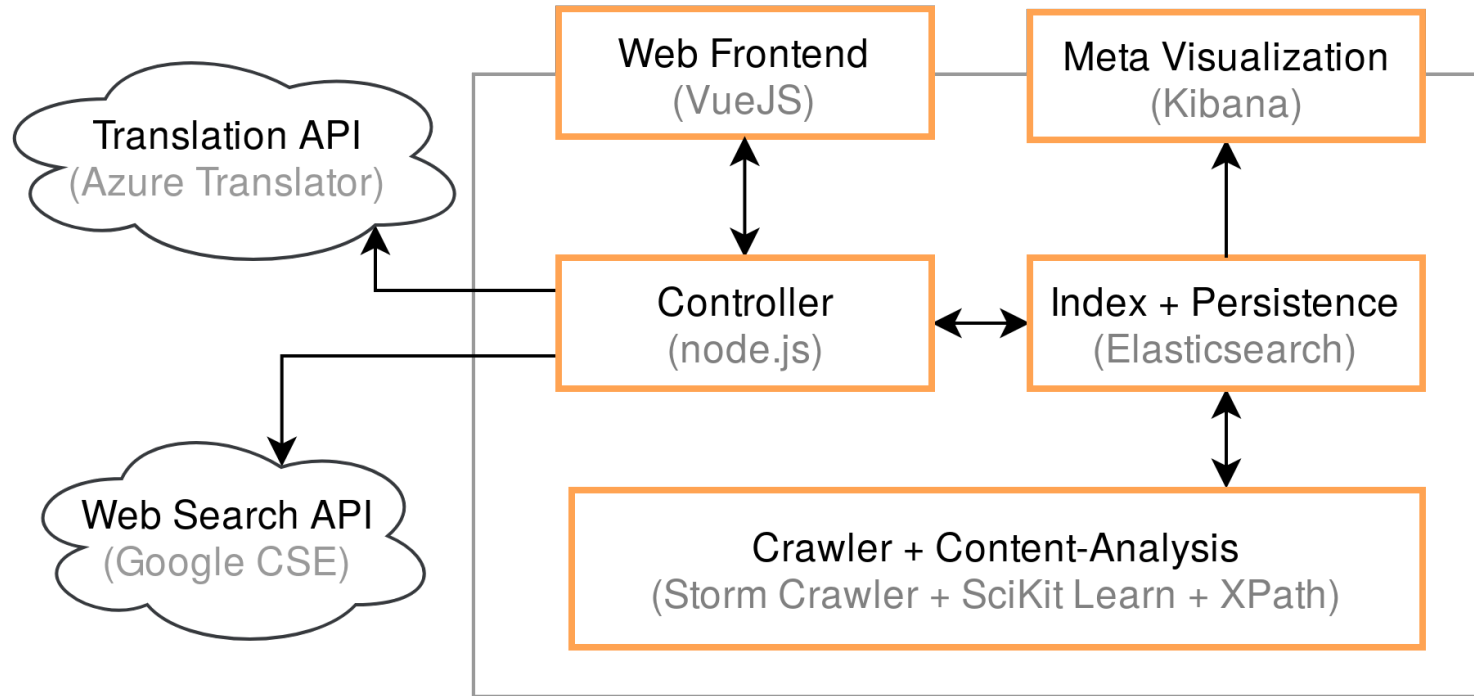    – web (map) applications → undocumented APIs

52n

# CHALLENGES: MULTILINGUAL SEARCH

- global data coverage is helpful for predictive models
  - effective search in local languages only

  → need to reduce language barriers via translation of
  - search keywords
  - page content

  → language independent content analysis

52n

# FEATURES

- web search
  - automated translation of search terms by country
  - crawling of outlinks for broader spectrum
- content analysis
  - classification via custom machine learning model
  - metadata extraction
- result list & crawl management via web-UI
  - result list with extracted metadata & hotlinks to translated pages
  - result labeling interface for classifier training
  - performance metric visualization

# ARCHITECTURE



- job based: minimize crawl effort by crawling on specific topics only
- stream based: results accessible as they come in

# DEMO

# DEMO

# REVIEW

Were the initial goals met?

✔ effective discovery of websites & data sets

✔ language-abstracting interface

✗ classification not operational for every language

52n

# REVIEW

Future work

- train classifier for more languages

- check more dataset annotations (→ schema.org)

- performance optimisation of crawler

52n

# THANKS!

- Code: Apache 2 licensed on GitHub
  github.com/52north/ecmwf-dataset-crawler

- Read more in the..

  - project wiki:  github.com/52north/ecmwf-dataset-crawler/wiki

  - 52N blog:     blog.52north.org/2018/06/26/ecmwf/