# Business register quality enhancement

WP3-UC5 members:
CBS, HSL, SCB, SF, STATA

WIH-CON, June 12 2023, Brussels

**Trusted Smart Statistics – Web Intelligence Network**
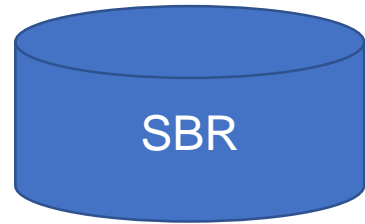Grant Agreement: 101035829

Web Intelligence
Network

**Funded by
the European Union**

# Contents

- The main concept
- Search, scrape, link, train, predict, derive
- Other data sources
- Putting it all together
- Zooming out: web data and survey design
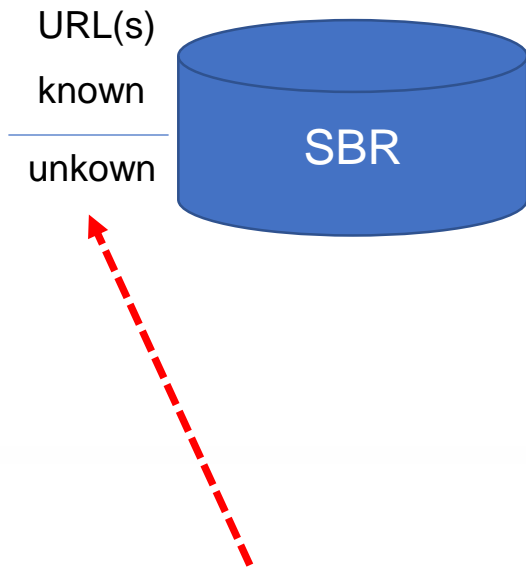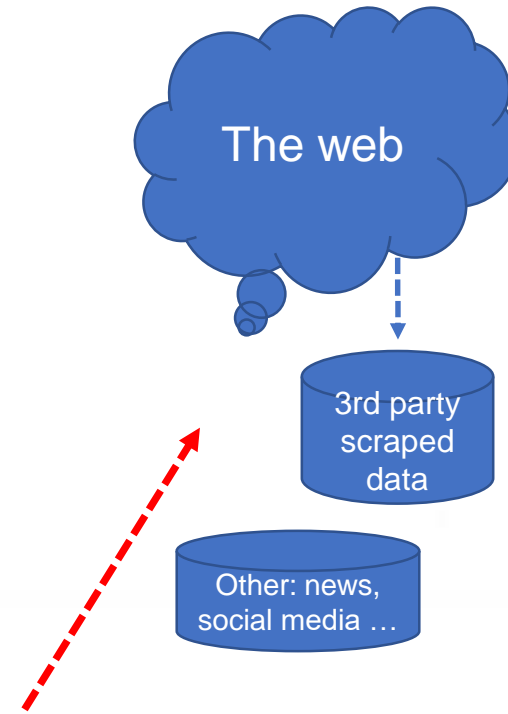- Future work
- Wrap up

Web Intelligence
Network

Funded by
the European Union

# The main concept (1)

SBR

Statistical Business Register

The web

# The main concept (2)

URL(s)

known

───────

unkown

SBR

The web

3rd party scraped data

Other: news, social media …

The ratio known/unknown is country-specific

Web data is more than websites only

Web Intelligence
Network

Funded by
the European Union

# The main concept (3)

URL(s)
known
___
unkown

SBR

A. Find legal units
(URL finding / linking)

The web

B. Derive / improve
statistical variables
(NACE codes)

3rd party
scraped
data

Other: news,
social media …
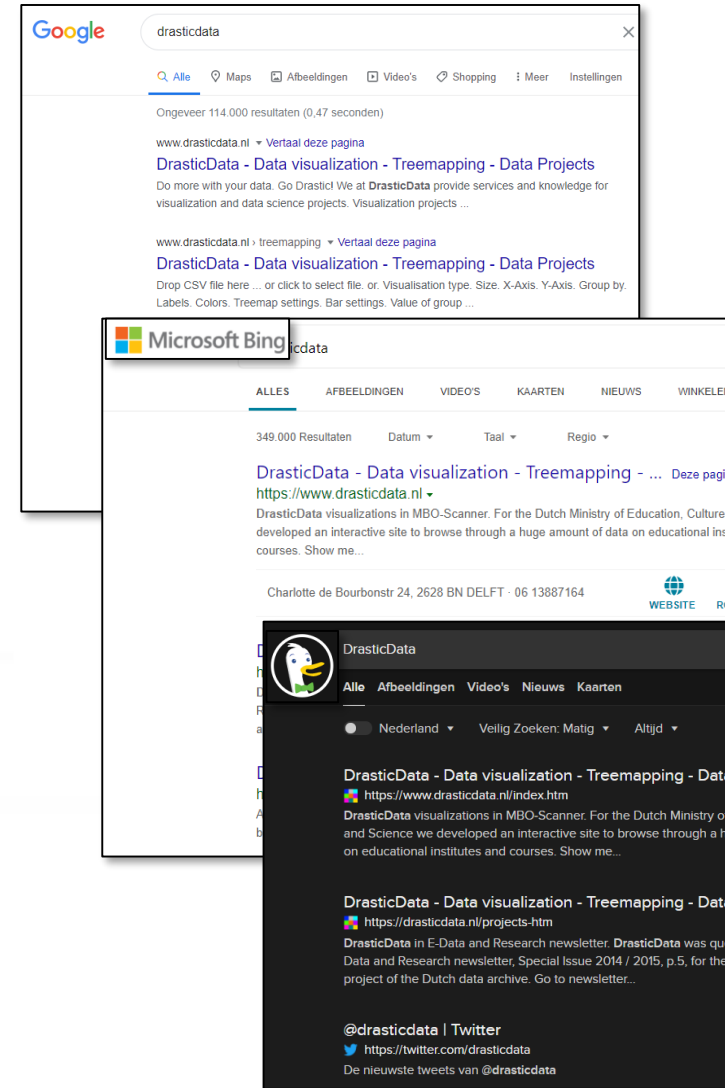
**Web Intelligence**
Network

**Funded by**
**the European Union**

# Search (1)

- To **find** or **verify** URLs for legal units (LUs)
- **Automatically** search on what we know from SBR
  - E.g. Name, address, municipality, id, and/or contact info
- Using a free or paid **API**
- Search engine leakage manageable:
  - Use paid/trusted search engines
  - Use search phrase wisely
  - Spread across search engines and in time

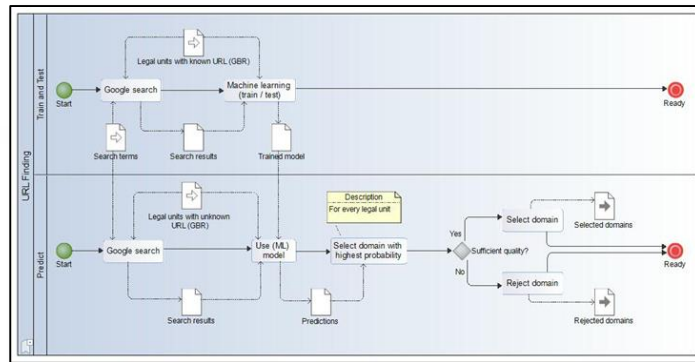https://ec.europa.eu/eurostat/cros/content/url-finding-methodology_en

# Search (2)

How to select the right URL from a list of search results?

- Using *snippet* and/or *extra scraping* step
- Use an *ML* model capturing the *search engine behaviour*:
  - Train and test on set of LUs with known URL
  - Predict URL from list of search results

https://github.com/SNStatComp/urlfinding

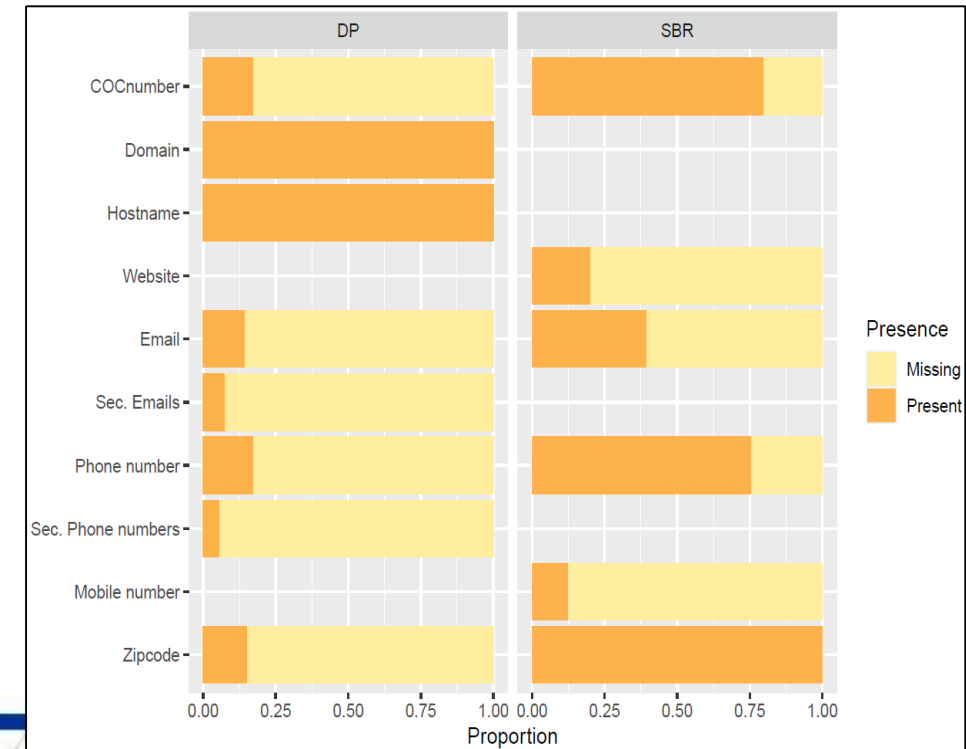https://github.com/EnterpriseCharacteristicsESSnetBigData/StarterKit

# Scrape

- Two types:
  - *Generic*: no prior knowledge of site structure
  - *Specific*: scraper is designed for specifics of website
- *Focused* scraper:
  - Gives priority to those parts of websites that are expected to contain valuable info, for example "about us" or "vacancies"
- National legislation might enforce *identifying information* on websites:
  - tax-id or COC-id
- Be aware of *n-to-m relationships* LU <-> website
  - LU might have multiple websites
  - Register the main website (if identified) or all?
  - Small business might be present only on business services portal listing many different small companies

# Link 3rd party web data

- Use web data collected by 3rd parties if added value is considerable
- NL > 2 yrs experience with DataProvider (DP) data
- Monthly datasets, linking to SBR
- Not always easy, gaps in data
- For ~15% of LUs a URL could be deduced

# Train, predict, derive



```
## [1] "enterprise"   "company"      "unternehmen"  "home"
## [5] "welcome"      "ueber"        "über uns"     "über"
## [9] "geschichte"   "about us"     "uber uns"     "about"
## [13] "unsere"      "willkommen"   "produkt"      "product"
## [17] "artikel"     "article"      "organisation" "dienstleistung"
## [21] "angebot"     "leistung"     "offer"
```

NACE detection

- AT: Word-driven NACE-1 prediction (XGBoost)

- NL: Predict whether a registered NACE is incorrect

- SE: NACE detection experiments with KB-BERT method adapted and extended for Swedish language

Correcting or complement administrative information:

- HE: contact information discovery from websites
  - emails, classified into functional/high/medium/low
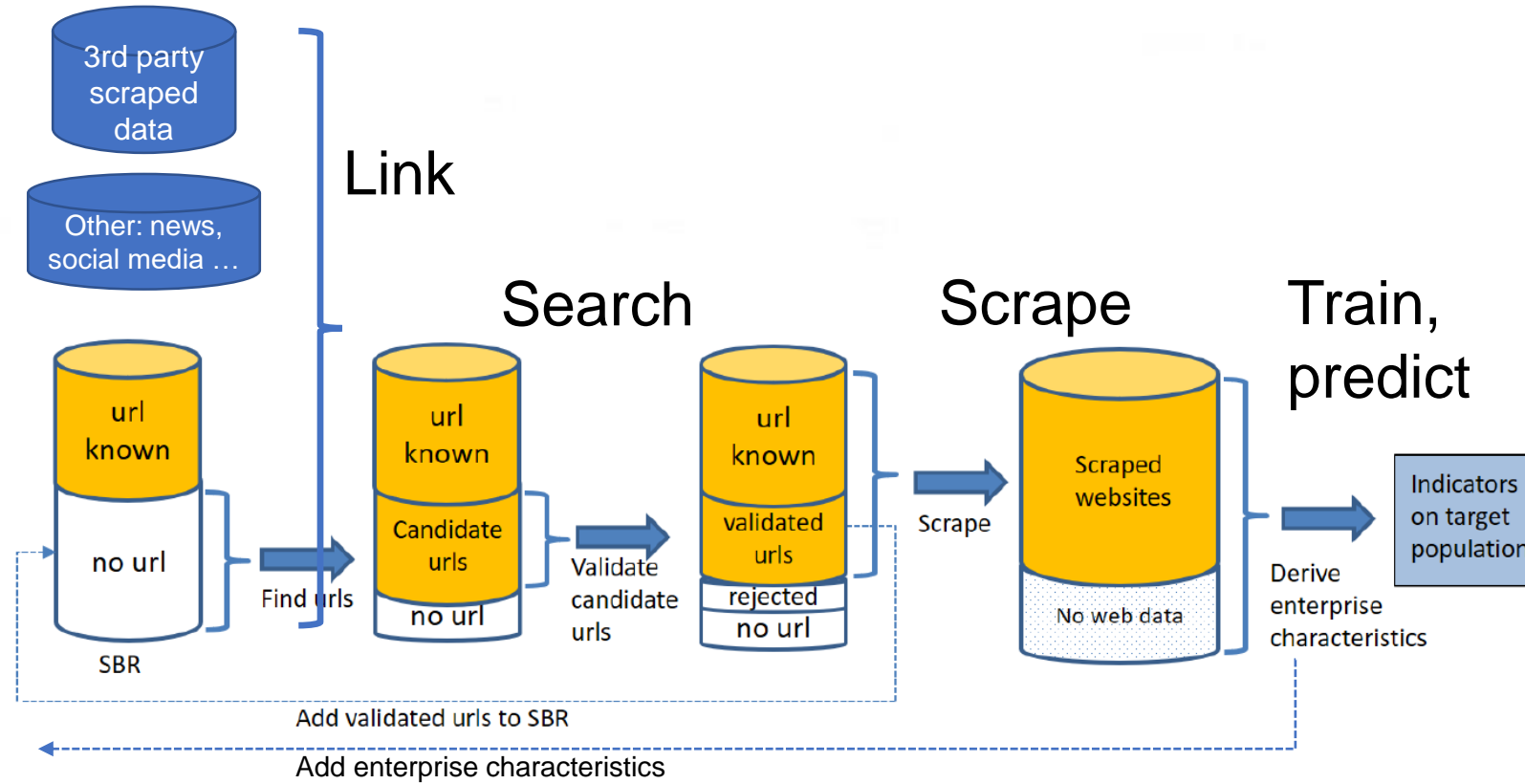
Web Intelligence
Network

# Other data sources

- Use domain registry for (additional) URL finding:
  - Degree of openness varies per country and domain
  - .nl domain is not open; .fi domain has public API
- Wikipedia / dbpedia / business collection portals
- News releases, social media
- Financial / annual reports
- OJAs

# All together

# All together: work in year 2 of project:

| Main topic | Detailed work | Update / New |
|---|---|---|
| URL finding | URL finding: updates on experiences from Statistics Hesse (HE) | Update |
| | Update on linkage process at (NL) | Update |
| | Using domain registry data (FI) | New |
| | URL finding first experiences: finding domain registry data (SE) | New |
| Business register enhancement | Update on NACE classification at Statistics Austria (AT) | Update |
| | Update on on detection of NACE misclassifications and on NACE prediction (NL) | Update |
| | First experiences NACE detection (SE) | New |
| | Contact information discovery from enterprise websites (HE) | New |

Deliverable 3.2: WP3 2nd Interim technical report

Web Intelligence Network

Funded by the European Union
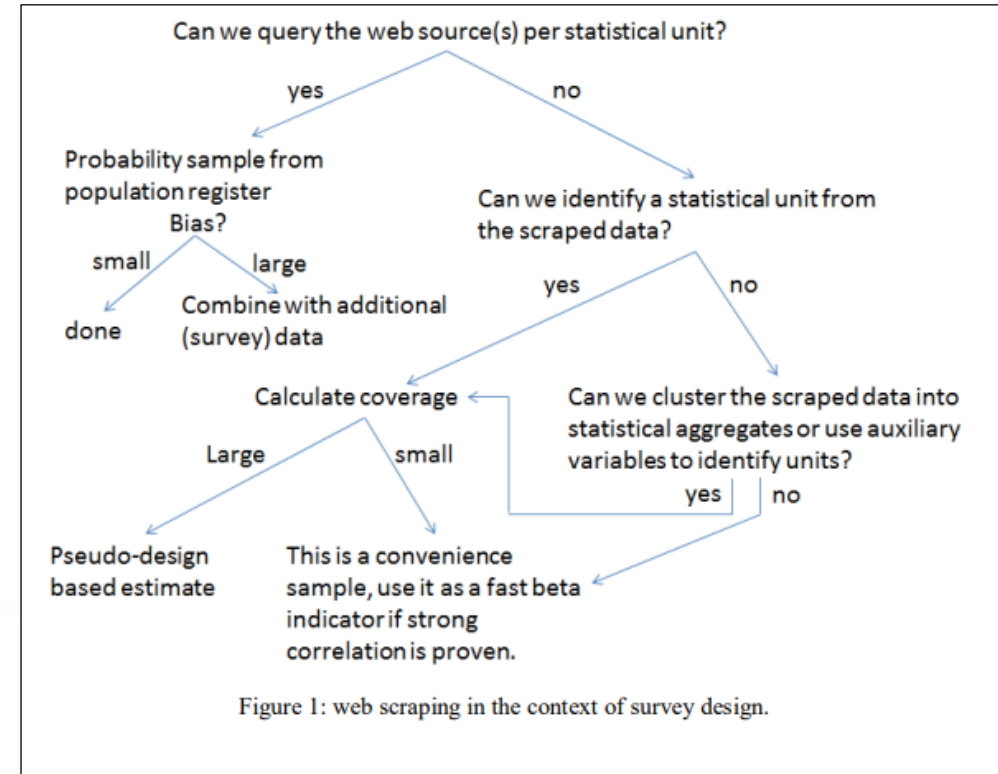
# Zooming out: web data and survey design (1)

#### Abstract

*Web scraping – the automatic collection of data on the Internet – has been used increasingly by national statistical institutes (NSIs) to reduce the response burden, to speed up statistics, to derive new indicators, to explore background variables or to characterise (sub) populations. These days it is heavily used in the production of price statistics. In other domains it has proven to be a valuable way to study the dynamics of a phenomenon before designing a new costly statistical production chain or to supplement administrative sources and metadata systems. Technical and legal aspects of web scraping are crucial but also manageable. The main challenge in using web scraped data for official statistics is of a methodological nature. Where survey variables are designed by an NSI and administrative sources are generally well-defined and well-structured, data extraction from the web is neither under NSI control nor well-defined or well-structured. A promising approach however is to combine high-quality data from traditional sources with web data that are more volatile, that are usually unstructured and badly-defined but in many cases also richer and more frequently updated. In this paper we reflect on the increasing use of web scraping in official statistics and report on our experiences and the lessons we learned. We identify the successes and challenges and we philosophise how to combine survey methodology with big data web scraping practices.*

Figure 1: web scraping in the context of survey design.

General workflow for any web source

BigSurv2018

https://www.researchgate.net/publication/327385487_Web_scraping_meets_survey_design_combining_forces

# Future

- What can web data tell us about *enterprise networks*?
- Can *AI*, such as Chat-GPT help us finding websites or deriving variables?

# Wrap up

- The web is a rich source on enterprise information for official statistics
- Starting from a Statistical Business Register we use a mix of searching, scraping, linking 3rd party web data and machine learning to enhance the business register
- Be aware of n-to-m relationships LU <-> web data
- Other sources can be domain registry, news, social media, OJAs
- This fits into a more general view: web data and survey design
- Future topics: enterprise networks, AI

# Questions / ideas welcome

- Olav ten Bosch (CBS, Netherlands), coordinator, [o.tenbosch@cbs.nl](mailto:o.tenbosch@cbs.nl)
- Alexandra Ils (HSL, Germany)
- Arnout van Delden (CBS, Netherlands)
- Heidi Kuhnemann (HSL, Germany)
- Johannes Gussenbauer (STATA, Austria)
- Katja Löytynoja (SF, Finland)
- Nick de Wolf (CBS, Netherlands)
- Pieter Vlag (SCB, Sweden)
- and others