# Providing Online Based Enterprise Characteristics with the Web Intelligence Hub.
## Experience of the Web Intelligence Network activities

**Jacek Maślankowski** (j.maslankowski@stat.gov.pl)
Statistics Poland
University of Gdańsk

2025/02/04

**Trusted Smart Statistics – Web Intelligence Network**

Grant Agreement numer: 101035829 (2020-PL-SmartStat)

Web Intelligence Network

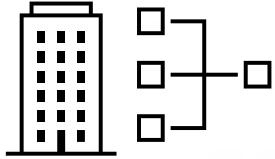Funded by the European Union

Funded by the European Union

# Agenda

Background information

Methodology and data processing

Main findings and conclusions

Web Intelligence Network

Funded by the European Union

# What is Online Based Enterprise Characteristics?

OBEC – term have been used since 2018

The use of a website by the enterprise to present its 'business', with extension to social media perspective.

It includes not only the existence of a website which is located on servers belonging to the enterprise or at one of the enterprise's sites, but also third party's websites.
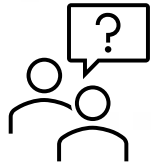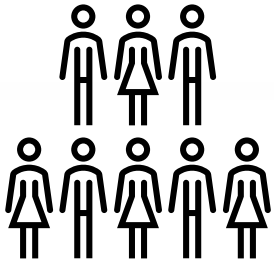
# How can we define OBEC population?

Methodological manual for data compilers and users of the ICT survey:
A4. Does your enterprise have a website?

The population of OBEC use case consists of enterprises having website and employing 10 or more employees.

**Web Intelligence** Network

Funded by the European Union

# OBEC population = URL database for WIH

## It consists of two datasets:

- For the Web Intelligence Platform (web scraping tool with data storage Amazon OpenSearch):
  - Anonymized ID
  - Enterprise URL
  - Group to which it belongs (e.g. /OBEC)
- For further processing in Datalab (JupyterLab with Python, RStudio):
  - Business register ID
  - Anonymized ID
  - Other attributes

## To link with business register

# What indicators OBEC can provide?

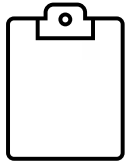| Main indicators | | Other indicators | | |
|---|---|---|---|---|
| Social Media Presence | E-commerce | User friendliness | Climate neutrality | Innovations |
| Type of social media / Purpose of social media | E-shops / Product list | GDPR policy / Accessibility / Sustainability | Green economy / SDG support | Patents / Collaboration with Universities / Address - location - technical incubator / Innovative companies and regions |

# List of indicators

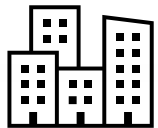| INDICATOR | BREAKDOWN1 | BREAKDOWN 2 | … | MEASURE |
|---|---|---|---|---|
| Social media presence | Instagram/X/Facebook/LinkedIn/Xing/Tiktok/<br><br>Pinterest/?Viadeo/?Yammer/Flickr/Instagram/Snapchat/<br><br>Slideshare/Threads/Whatsapp/Reddit/Telegram/Discord/WeChat/ | NACE | Company size, Region | Present / Not in any; in specific |
| E-commerce | - | NACE | Company size, Region | Yes / No |
| Chatbot | - | NACE | Company size, Region | Yes / No |
| Innovation | With extension to other keywords | NACE | Company size, Region | Yes / No |
| Multilanguage support | - | NACE | Company size, Region | Yes / No |
| Extracting contact information | Phone number, e-mail address, zip code | NACE | Company size, Region | Yes / No |

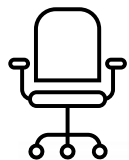# What is the reason to include enterprises having 10+ employees?

Traditional questionnaire – Survey on ICT Usage and E-commerce in Enterprises

URL database (Uniform Resource Locator – website address)

Significant percentage of small enterprises does not have a website compared to larger enterprises
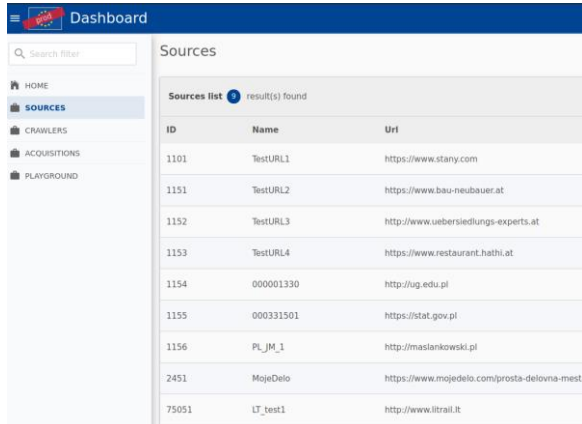
Names of small companies are usually not unique

# WIP for the OBEC

General rules:
- one unified data storage for all countries
- each country is contributing to the list of URLs by managing and maintaining
- websites are downloaded for further data processing
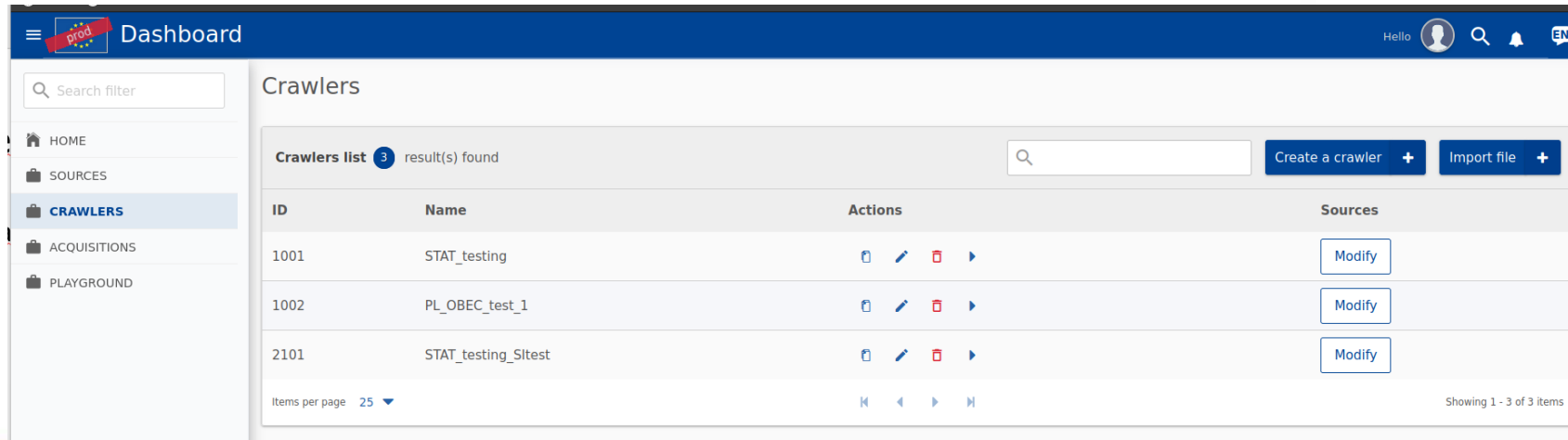- scalable – data stored in NoSQL database

Web Intelligence Network

Funded by the European Union

# Prerequisites for OBEC.
## What we need to create statistics with WIH?

- Access
  - WIN WIP (https://prod.wihp.ecdp.tech.ec.europa.eu/screen/home)
  - WIN Datalab (https://prod.wihp.ecdp.tech.ec.europa.eu/screen/home)

- Scripts
  - Latest release on GitHub
  - JSON file with URLs for web scraping

- Skills (basics)
  - Python
  - Amazon OpenSearch
  - Linux

# Collecting and processing data in 6 steps



1. URL enterprises with anonymized Business Register numbers.

2. Uploading the dataset of URLs into WIP

3. Modifying necessary parameters and starting the Crawler.

4. Opening WIP – cloning the Github open code on SMP/E-commerce/Multilanguage

5. Executing the software in Python and accessing the data from Crawler already scrapped.

6. Waiting for the results and open CSV file already prepared.

Web Intelligence Network

Funded by the European Union

# Methodological obstacles identified

## OBEC

- Target population may be unknown
- Providing the list of URLs may be difficult for selected ESS countries
- Different rules apply to websites in different countries, e.g. for some countries it is obligatory to provide tax number on the website while for others not

Web Intelligence
Network

Funded by
the European Union

# Other issues

Legal aspects

Web scraping policy

Changes in rules of sharing URL population and data by NSIs

# Main findings

- Varying legal regulations at national level in terms of the transfer of data from registers

- The inconsistency of the algorithms for data collection and processing

- The discrepancy between the data obtained with the use of different methods are among the issues under particular consideration within the TSS-WIN project

Web Intelligence Network

Funded by the European Union

# Conclusions

Adaptation of ESSnet Big Data I, II, Eurostat work

Strict co-operation with Eurostat

Key success factor – the data on OBEC disseminated in the Eurostat Database

**Web Intelligence** Network

**Funded by the European Union**

# Thank you!

Jacek Maślankowski, j.maslankowski@stat.gov.pl (Statistics Poland)