# Constantly changing data landscape

Digitalization

Covid

Military conflicts

Climate change

………………

Emergent & growing

Timeliness

New domains

Cross-cutting topics

Higher granularity

……………………

New data sources

Big data

Data science

New processing possibilities

AI

ML

Increased processing power

Web Intelligence
Network

Funded by
the European Union

# What it means for official statistics

New methods

New legal arrangements

New work arrangements

New techniques

New cooperation & communication models

Space for experimentation with uncertain results

Changes in infrastructure

New skills

Re-shaping statistical production process
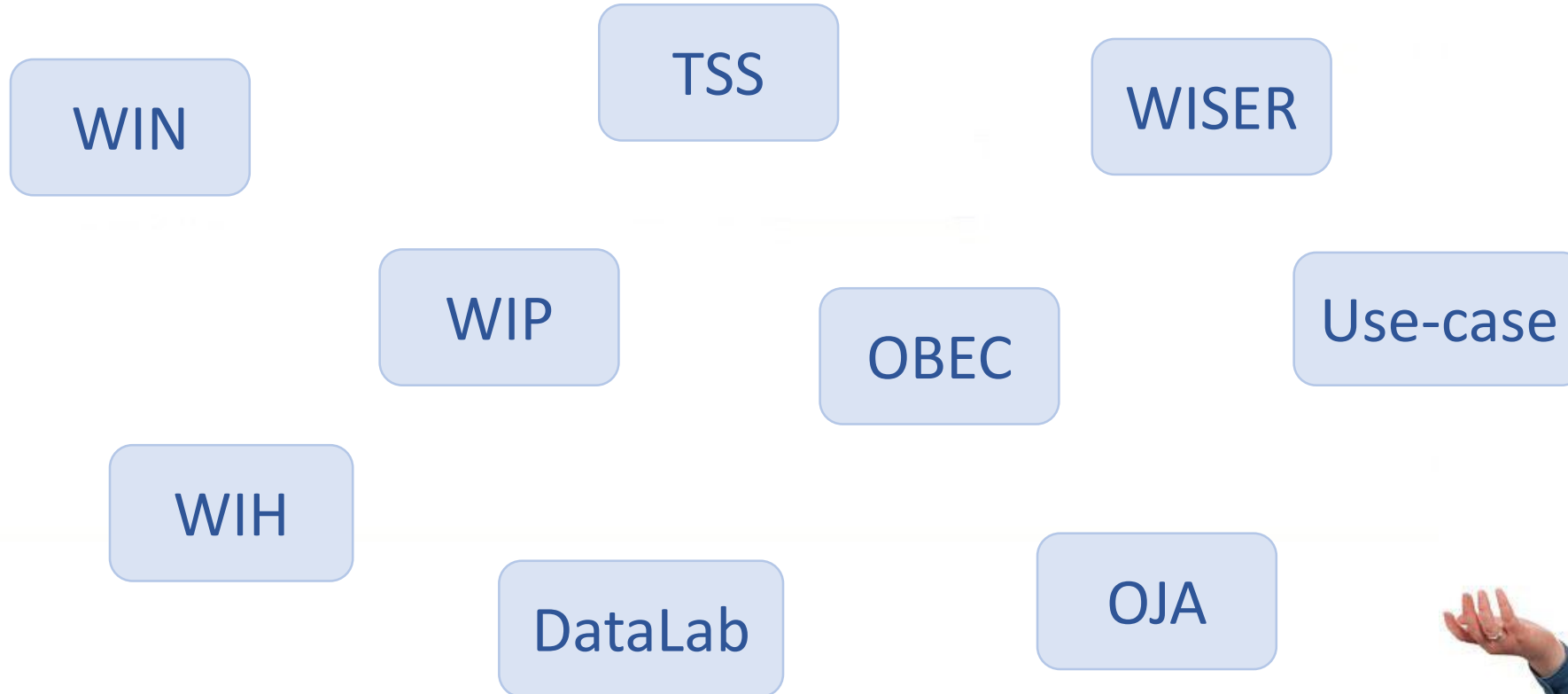
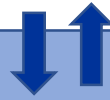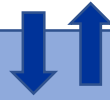Web Intelligence Network

Funded by the European Union

# WARNING!

TSS

WIN

WISER

WIP

OBEC

Use-case

WIH

DataLab

OJA

**Web Intelligence** Network

**Funded by the European Union**

# How it all started…

# What is Web Intelligence Hub (WIH)?

- Started as a concept – evolved towards tangible tools
- Web data acquisition, processing & analysis environment
- Centralized, shared system, pan-European platform



Web Intelligence Network

Funded by the European Union

# What is Web Intelligence Hub (WIH)?

# Why a shared system for web-based statistics?

- Different capacity of NSIs across Europe to use web data

- Different competency levels, scarcity of data science skills

- Infrastructure with big data capabilities required

- More efficient use of resources

**Web Intelligence** Network

**Funded by the European Union**

# What WIH concept has materialized into to date

Web Intelligence Platform (WIP)

→ Web content retrieval platform

→ DataLab (data access and analysis)

**Web Intelligence** Network

Funded by the European Union

# Web Intelligence Network (WIN)

**14 countries, 17 organizations, 100 members**

Contribute to the development of the WIH

Reach out to **all ESS countries**

Use web data, use the WIH

Web Intelligence
Network

Funded by
the European Union

# How these two work together

# What topics WIN is looking to

**Most mature use-cases**

Online job advertisements
OJA

Enterprise characteristics
OBEC

Online prices

Construction activities

BR quality enhancement

**Experimental research**

Real estate

Tourism

Traffic cameras

Web Intelligence Network

Funded by the European Union

# What topics WIN is looking to

**Most mature use-cases**

Online job advertisements
OJA

Enterprise characteristics
OBEC

Construction activities

Online prices

BR quality enhancement

**Experimental research**

Real estate

Tourism

Traffic cameras

Web Intelligence
Network

# Online Job Advertisements (OJA)



skills

localization

salary €

# Online job advertisements (OJA)

- Data acquisition and analysis system already in place – data centrally scraped by Eurostat/ Cedefop, accessible by NSIs via DataLab

- Job portals for each country selected, based on landscaping process with the involvement of national experts – coordinated by Eurostat/Cedefop

- Goal:
– augment labor market statistics
– provide information on skills

# What WIH concept has materialized into to date

Web Intelligence Platform (WIP) → Web content retrieval platform
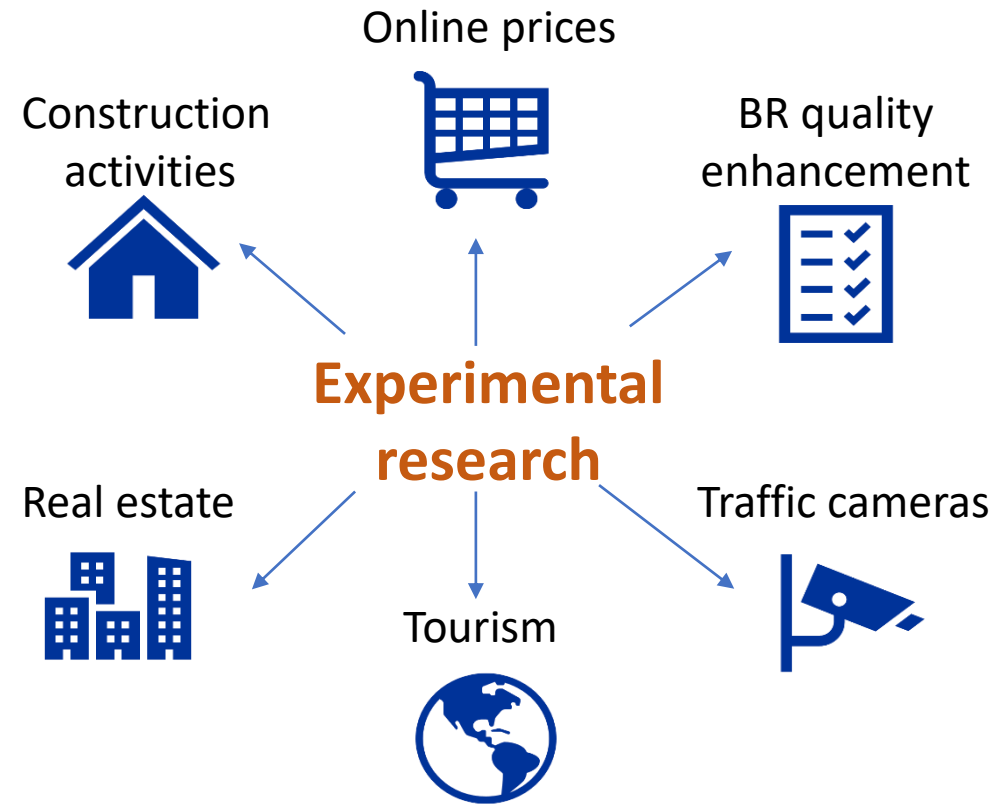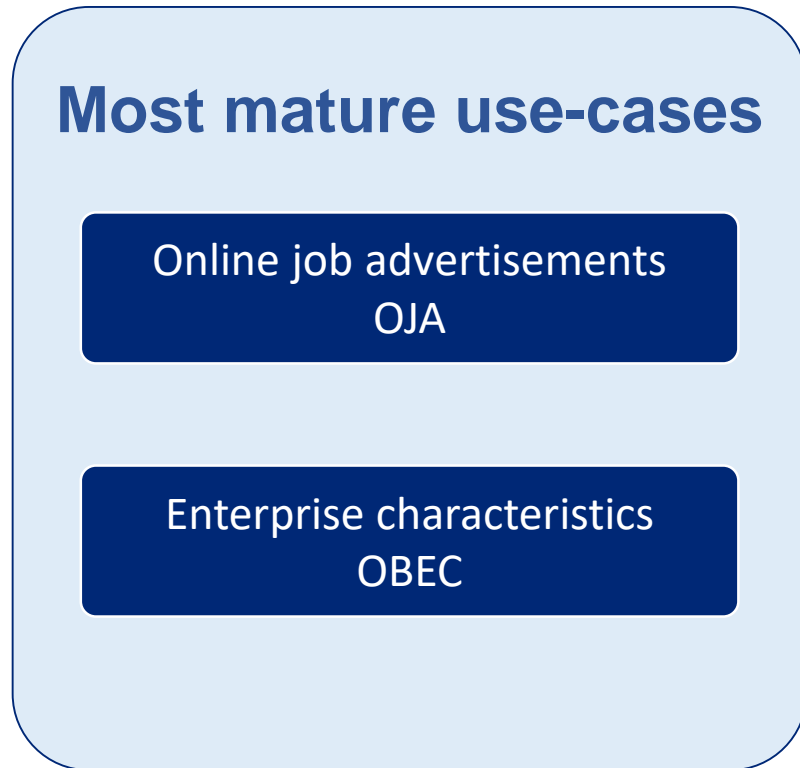
Web Intelligence Platform (WIP) → DataLab (data access and analysis)

Web Intelligence Network

Funded by the European Union

# Online job advertisements – key issues for the WIN

- Data source stability

- Quality of the data (overcoverage, undercoverage)

- Quality of data classification (e.g. ISCO, NACE, NUTS)

- Relevance of existing classifications (e.g. ISCO, ESCO)

# Online Job Advertisements – what WIN works on

| Data quality analysis | Improvement of ontologies | Calculation of experimental indicators |
|---|---|---|
| ↓ | ↓ | ↓ |
| Data annotation - company and economic activity | Verification of the NLP algorithm - detecting job occupation (ISCO classification) | Development of harmonized methodology |

Web Intelligence Network

Funded by the European Union

# Enterprise Characteristics (OBEC) – what WIN works on

**Definition of the OBEC population**

- Selection of potential data sources
- Methodology of URLs retrieval

**Selection of core and additional indicators**

- social media presence (C)
- e-commerce (C)
- user friendliness
- climate neutrality
- innovation

**Testing the WIP (web content retrieval platform)**

- Functional and non-functional requirements for the WIP
- Web scraping tests

Web Intelligence Network

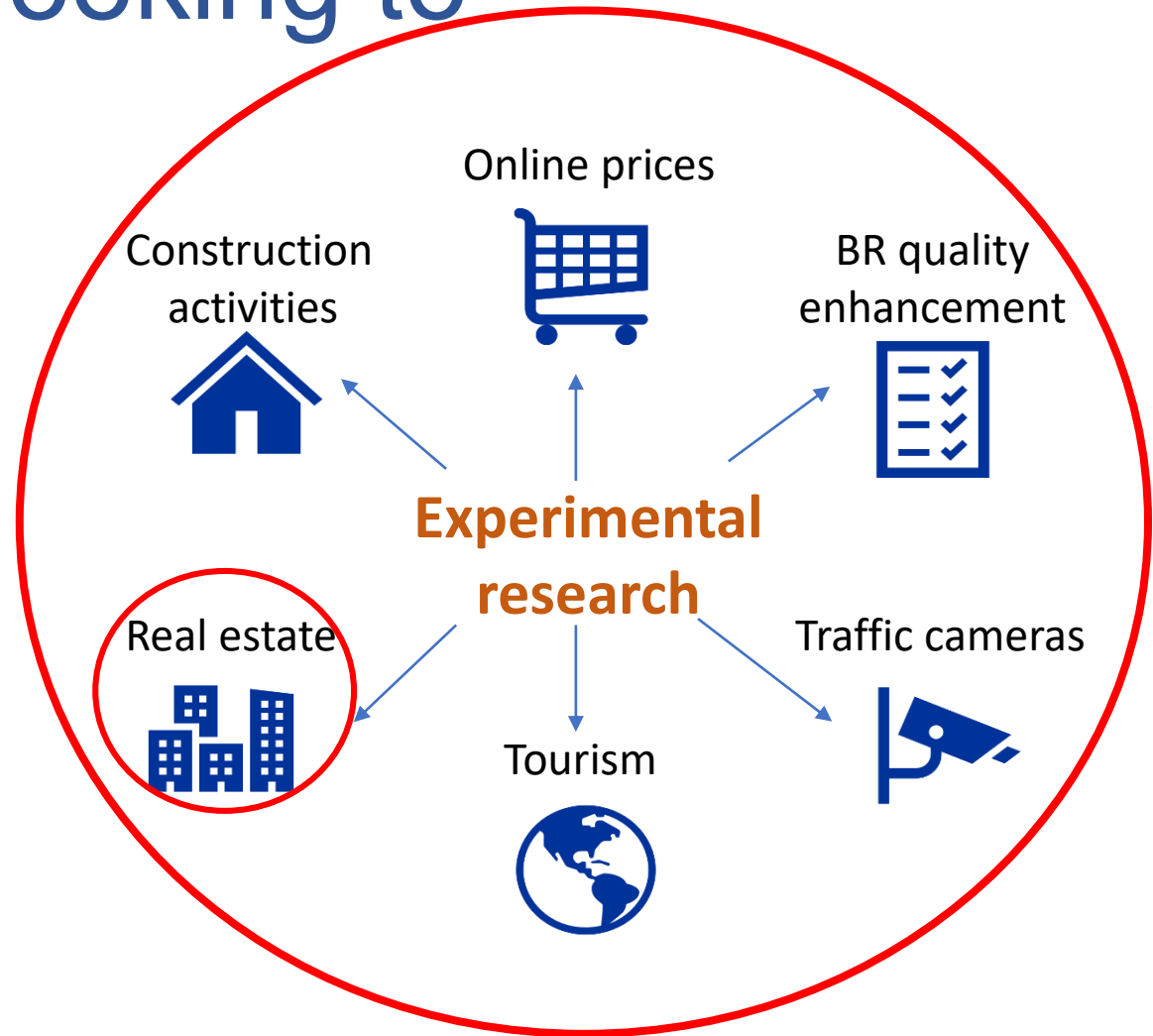Funded by the European Union

# What WIH concept has materialized into to date

# What topics WIN is looking to

**Most mature use-cases**

Online job advertisements
OJA

Enterprise characteristics
OBEC

Online prices

Construction activities

BR quality enhancement

Real estate

**Experimental research**

Traffic cameras

Tourism

**Web Intelligence** Network

**Funded by the European Union**

# Experimental research – how WIN works

- New data sources exploration & landscaping

- Programming, production of software

- Data acquisition and recording

- Data processing (e.g. de-duplication)

- Modelling and interpretation (i.e. data analysis and quality assessment)

- Dissemination of the experimental statistics and results

Joint work of 2-6 NSIs on each use-case ➡ Solutions applicable at the ESS level
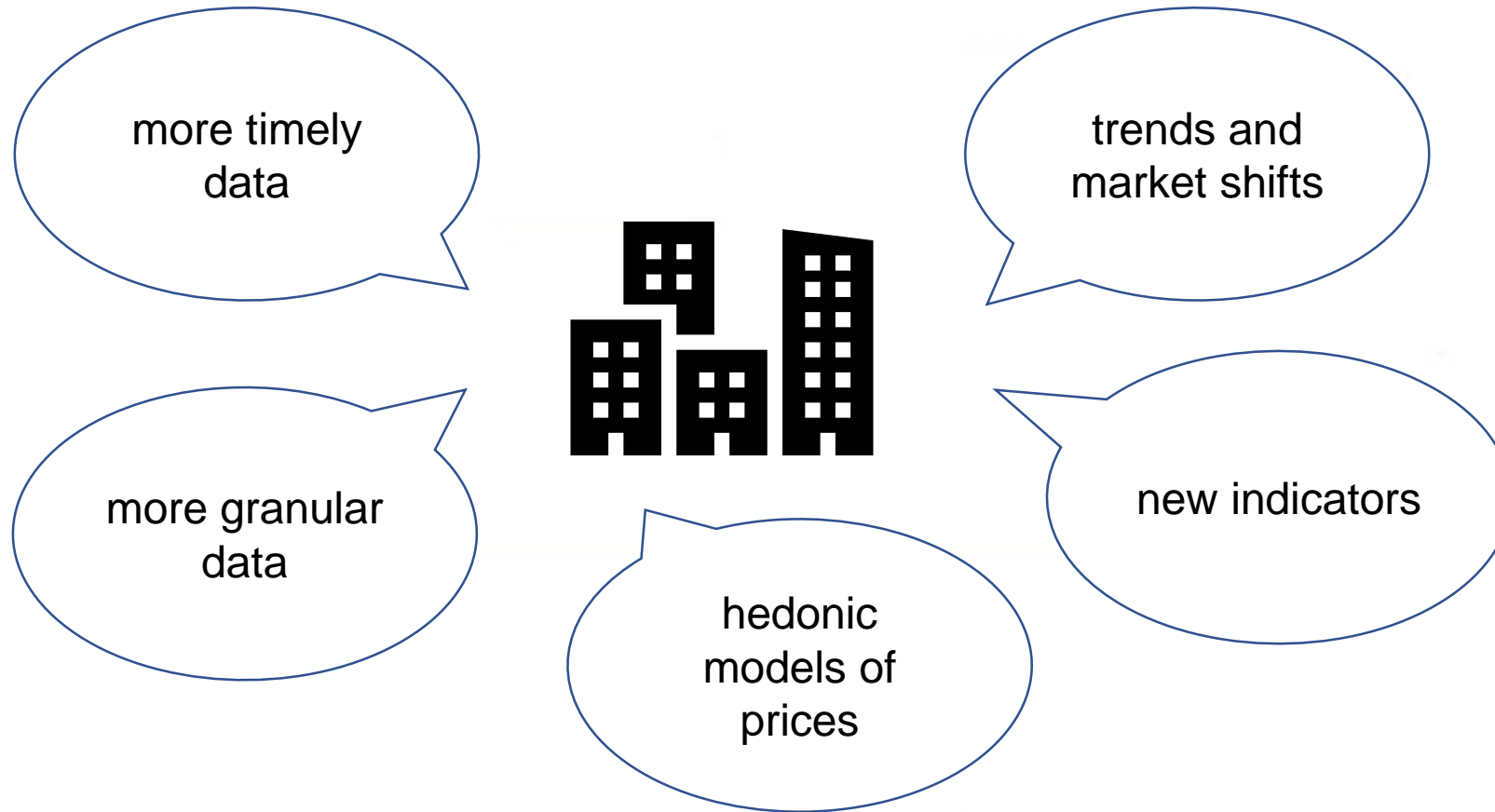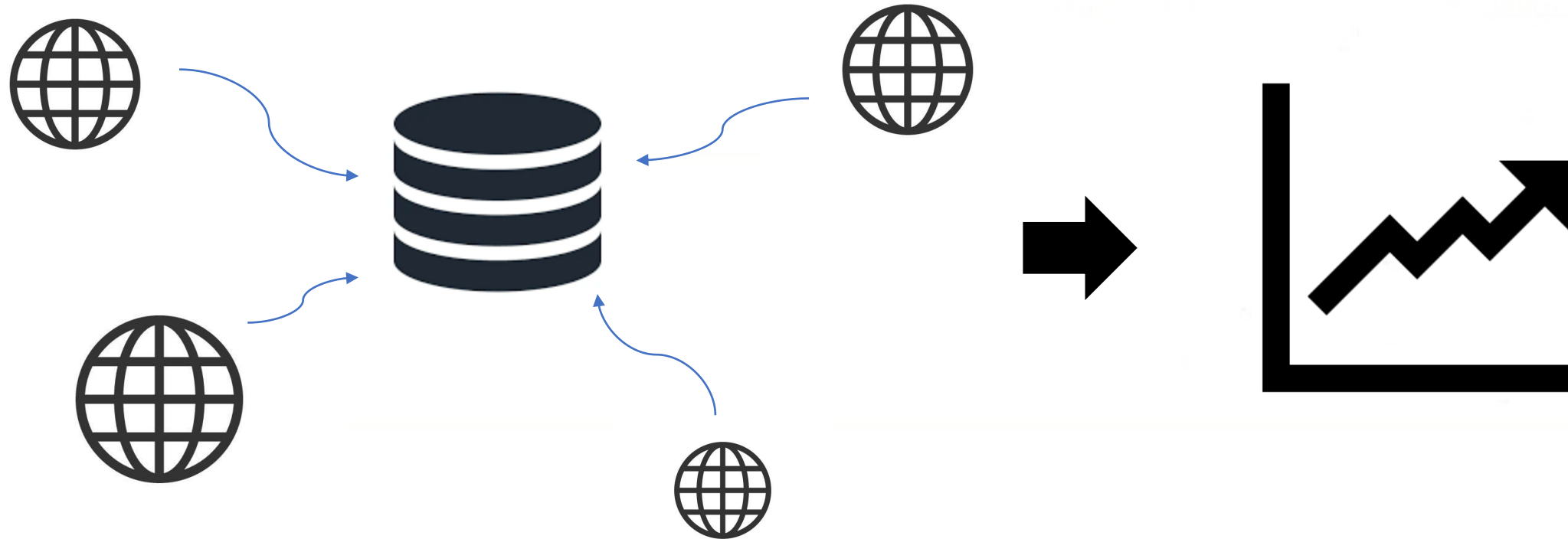
Web Intelligence
Network

Funded by
the European Union

# Experimental research – example: real estate

more timely data

trends and market shifts

more granular data

new indicators

hedonic models of prices

Web Intelligence Network

Funded by the European Union
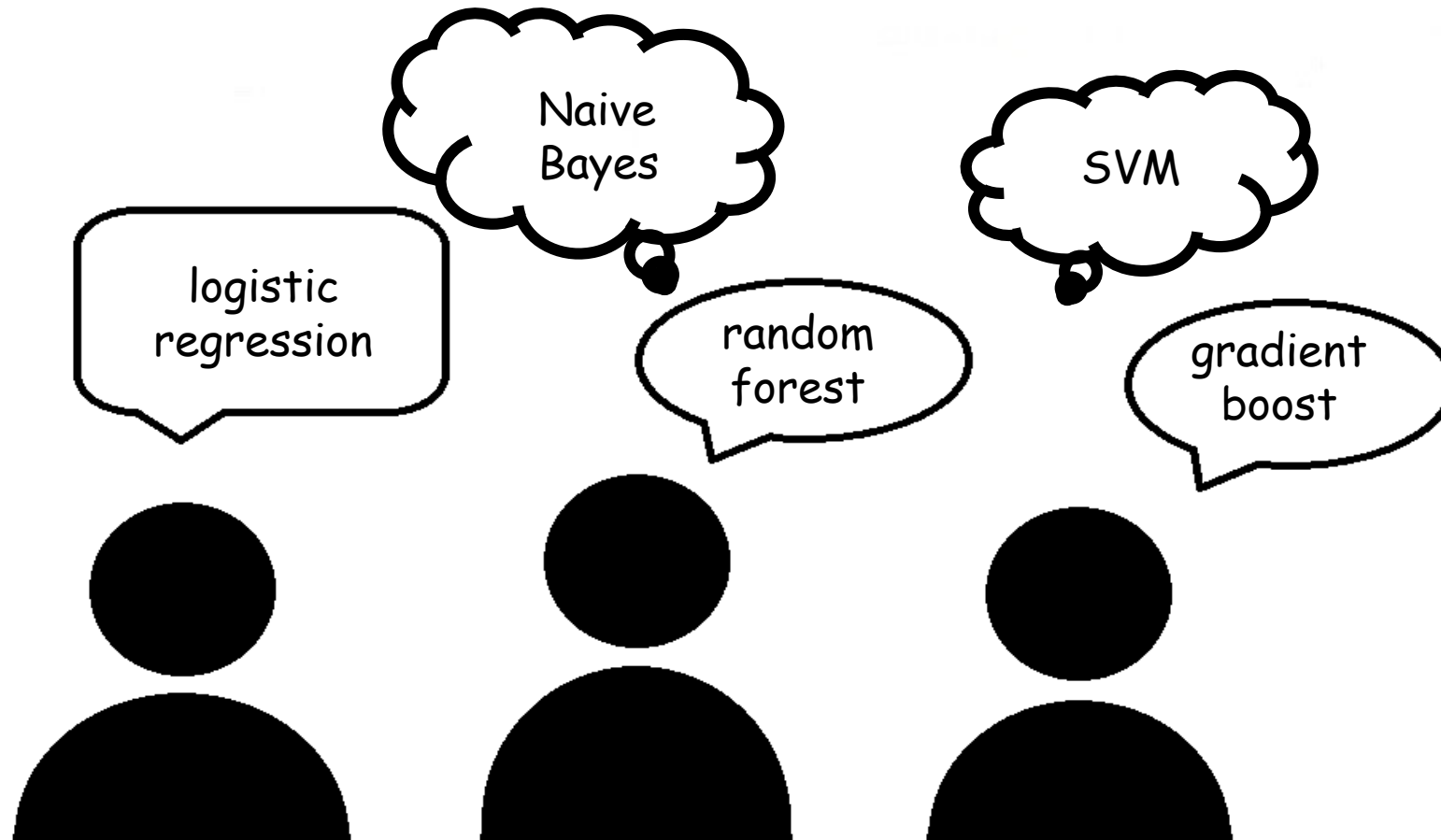
# Experimental research – rental market in Poland

# Apartment for rent, Krakow, Poland

We are pleased to present a property with an area of less than 70m2, located in close proximity to the Main Market Square in Krakow. The apartment is located in a historic tenement house, from which the view spreads over the tower of St. Mary's Basilica.
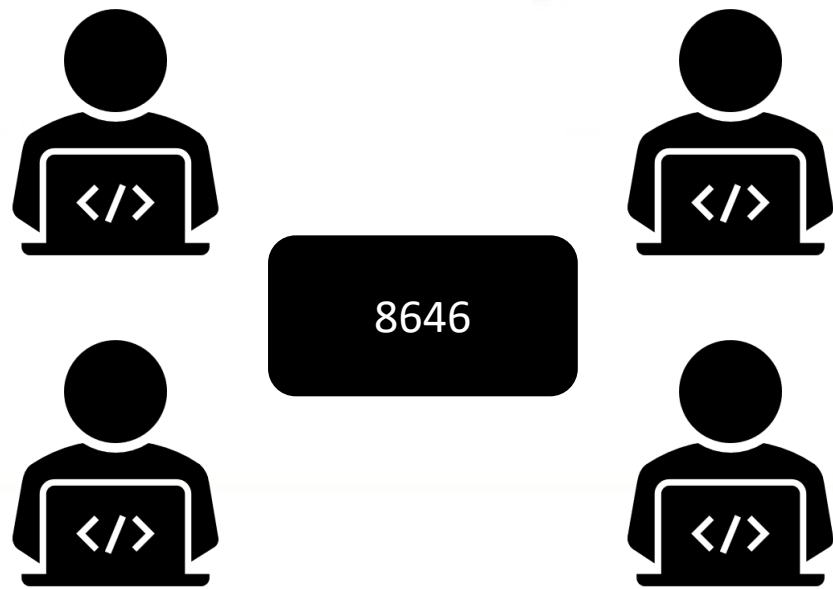
The apartment is located on the 4th floor of a historic tenement house located at the Main Market Square in Krakow. The apartment with a cadastrial area of 67.34 m2 consists of a large living room, 2 bathrooms with a shower and a bathtub, a kitchenette, two separate bedrooms and a spacious mezzanine with a second bedroom. The apartment is fully equipped and fully furnished. From the windows of the premises we can see the charming roofs of the Old Town and the tower of St. Mary's Basilica. The premises is suitable for short-term rental activities and is leased until 2028. With the current layout, it will comfortably provide accommodation for 8 people. The apartment is equipped with air conditioning and video security system.
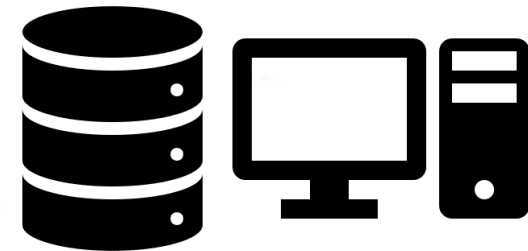
# Experimental research – method selection

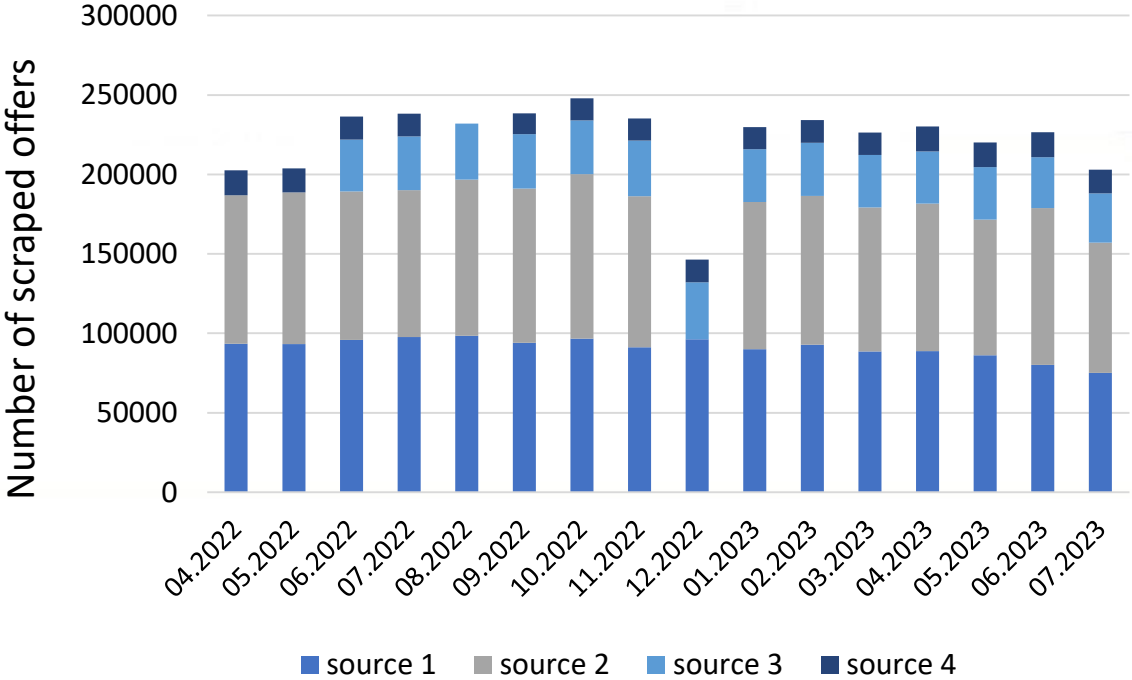# Experimental research – automatic classification of the apartment type
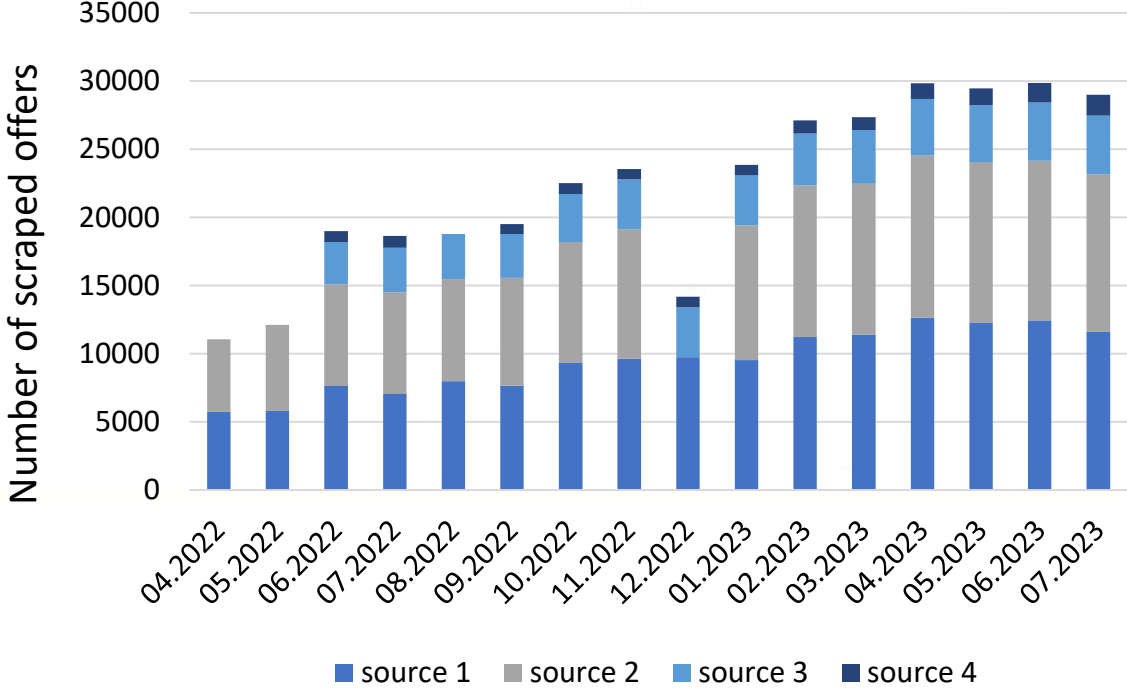
# Experimental research – sales /rental offers in Poland

### Offers for sale
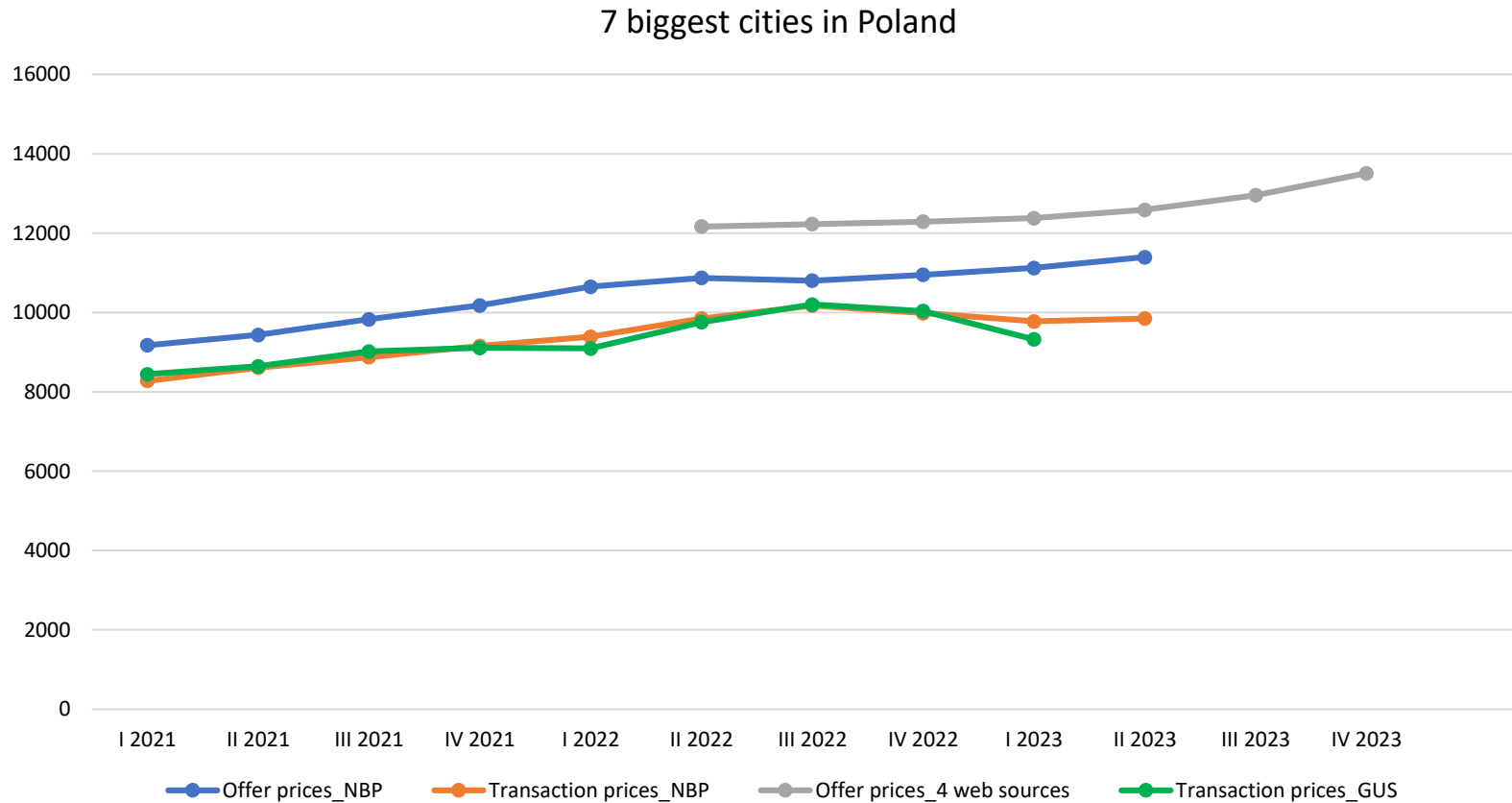


### Offers for rent

# Experimental research - comparison between web data and official statistics
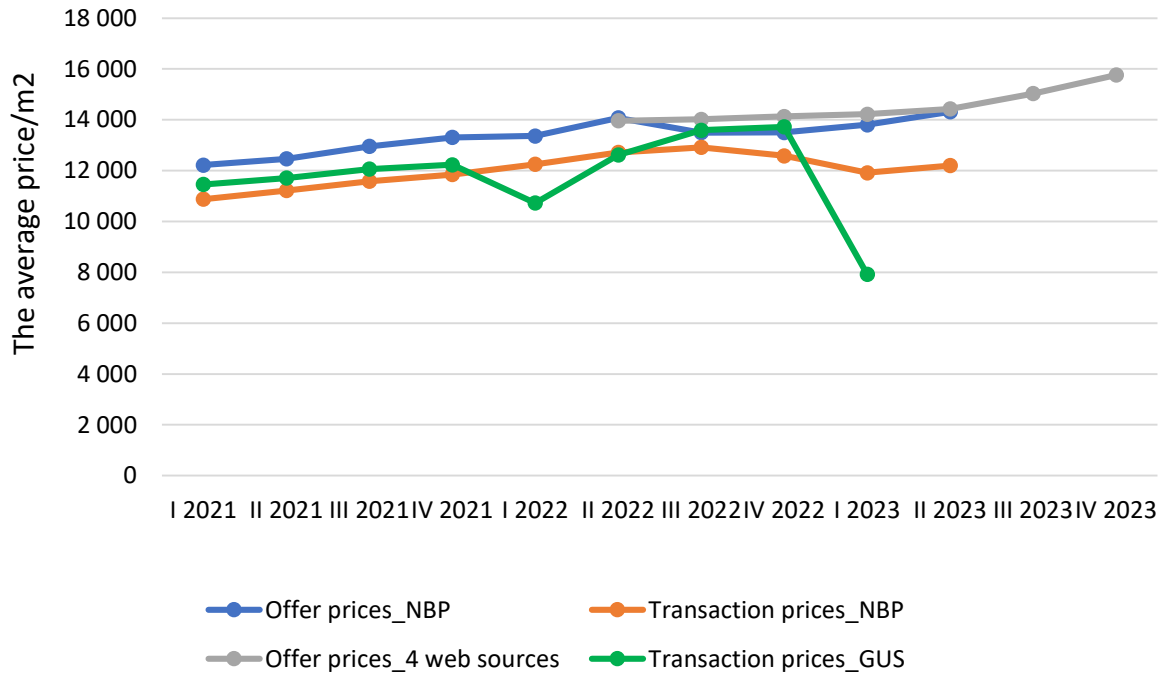
## 7 biggest cities in Poland



Legend: Offer prices_NBP · Transaction prices_NBP · Offer prices_4 web sources · Transaction prices_GUS

# Experimental research - comparison between web data and official statistics



Warsaw

Krakow

# Experimental research - what is ahead of us

- Quality assessment and improvement (de-duplication of offers, dealing with missing data etc.)

- Tests of different ML models

- Cooperation with domain experts and validation of results

- Calculation of experimental statistics

# Experimental research – quality and methodology issues

- Difference between offer and transaction prices

- Undefined population (lack of reference source for rental market)

- Duplication of offers (within a single portal/ across different portals)

- Multi-offers (apartments in new constructions)

- Missing values (e.g. price of apartments in new constructions)

# Challenges - use of web data & WIH by NSIs

Where does the innovation begin?

Who is the „owner" of the development?

Where are domain experts in this process?

Is web data of sufficient quality for official statistics?

Are we ready to sacrifice quality for timeliness?

Do we still need different methods to embrace new data sources for official statistics?

Is our approach to web data acquisition optimal?

Are we facing a shift of the paradigm of official statistics?

Web Intelligence Network

Funded by the European Union

# We need you to discuss it all!

**Meet us at conferences**

NTTS conferences

EUROPEAN CONFERENCE ON QUALITY IN OFFICIAL STATISTICS 2024 ESTORIL - PORTUGAL

IAOS INTERNATIONAL ASSOCIATION FOR OFFICIAL STATISTICS

**Visit us on-line**

**Look for our training, webinars, tutorials**

| | |
|---|---|
| Statistical Business Registers | Construction activity |
| Tourism data | Web Intelligence in Practice - OBEC |
| Online real estate market | OJA Training for WIN and WISER |
| Architecture, methodology and quality | And more… |

**Join WISER!!! Web Intelligence uSers Group**

Web Intelligence Network

Funded by the European Union

# Thank you…
# (…and do not hesitate to contact us!)

do.nowak@stat.gov.pl – Dominika Nowak, ESSnet WIN project coordinator
righi@istat.it - Alessandra Righi, WISER group coordinator

**Web Intelligence**
Network

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