# Methodologies for big data
# Overview of Istat's activities and open challenges

Mauro Bruno

Directorate for Methodology and Statistical Process Design (DCME)

Italian National Institute of Statistics – Istat

Monica Scannapieco

Italian National Institute of Statistics – Istat (currently working at AGENZIA PER LA CYBERSICUREZZA NAZIONALE)

# Outline

- Big Data @Istat: Why, When, How

- Big Data Projects @Istat:

  - **Text Processing** Pipelines

  - **Image Processing** Pipelines

  - Improving data dissemination **timeliness**

- Conclusions

Istat

# Big Data @Istat:
# Why, When, How

# Big Data @Istat: Why & When

- **European Statistical System strategic drivers**
  - Scheveningen Memorandum "Big Data in Official Statistics" - 2013
  - Bucharest Memorandum "**Official Statistics in a datafied society - Trusted Smart Statistics**" - 2018

- **Main Objectives**
  - Enrich statistical production with **new products**
  - **Enhance timeliness** in official statistical production
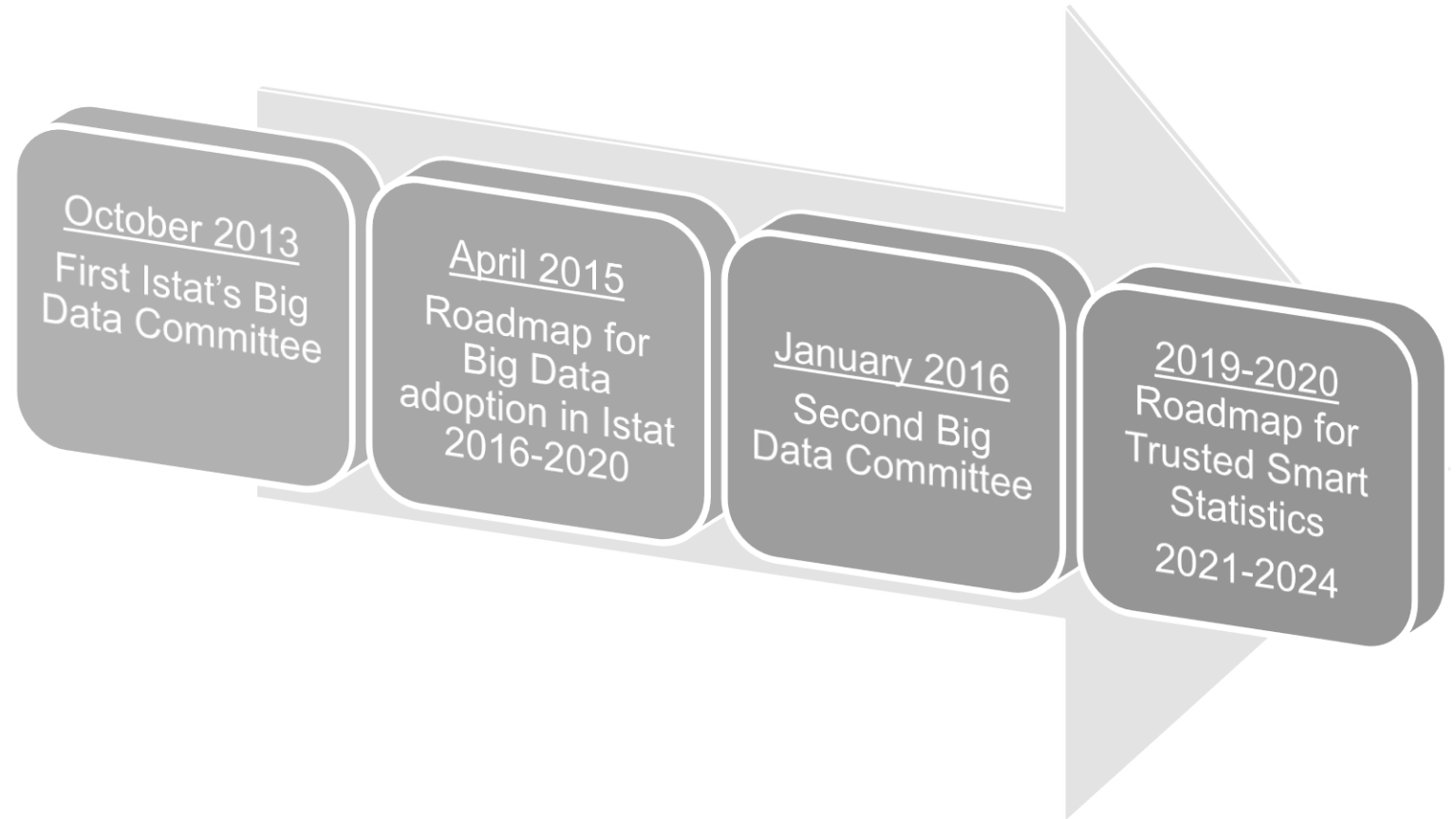  - Official statistics relevance in **new data ecosystem**



Source: Data Never Sleeps 8.0

# Big Data @Istat: How

○ **Istat's strategic context:**

The use of Big Data in Official Statistics requires **methodological**, **technological** and **organizational** investments.

Starting from 2013, Istat created a Big Data Committee responsible of the Big Data strategy…



October 2013
First Istat's Big Data Committee

April 2015
Roadmap for Big Data adoption in Istat 2016-2020

January 2016
Second Big Data Committee

2019-2020
Roadmap for Trusted Smart Statistics 2021-2024

# Big Data @Istat: How

o **Istat's strategic context:**

The use of Big Data in Official Statistics requires **methodological**, **technological** and **organizational** investments.

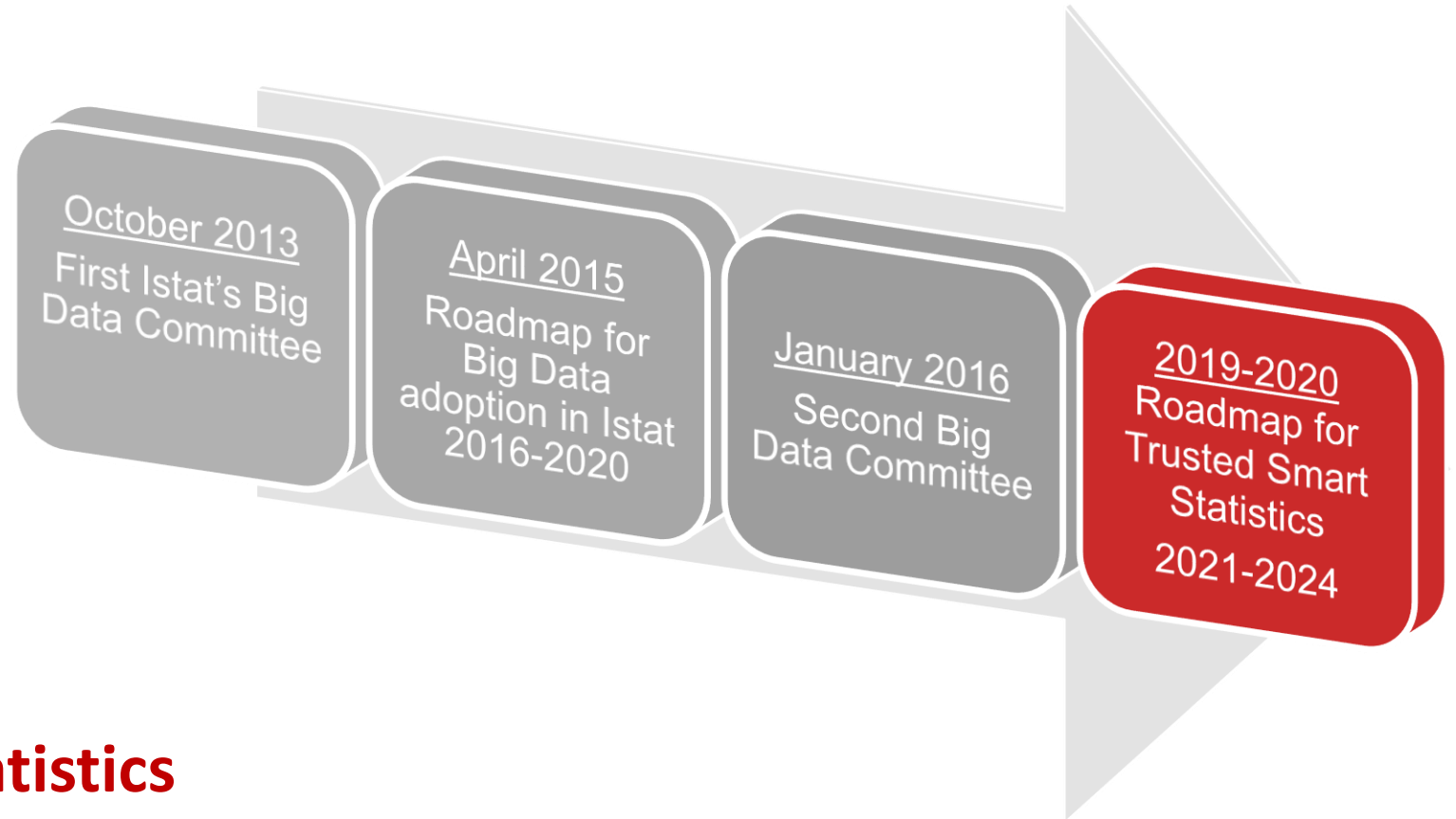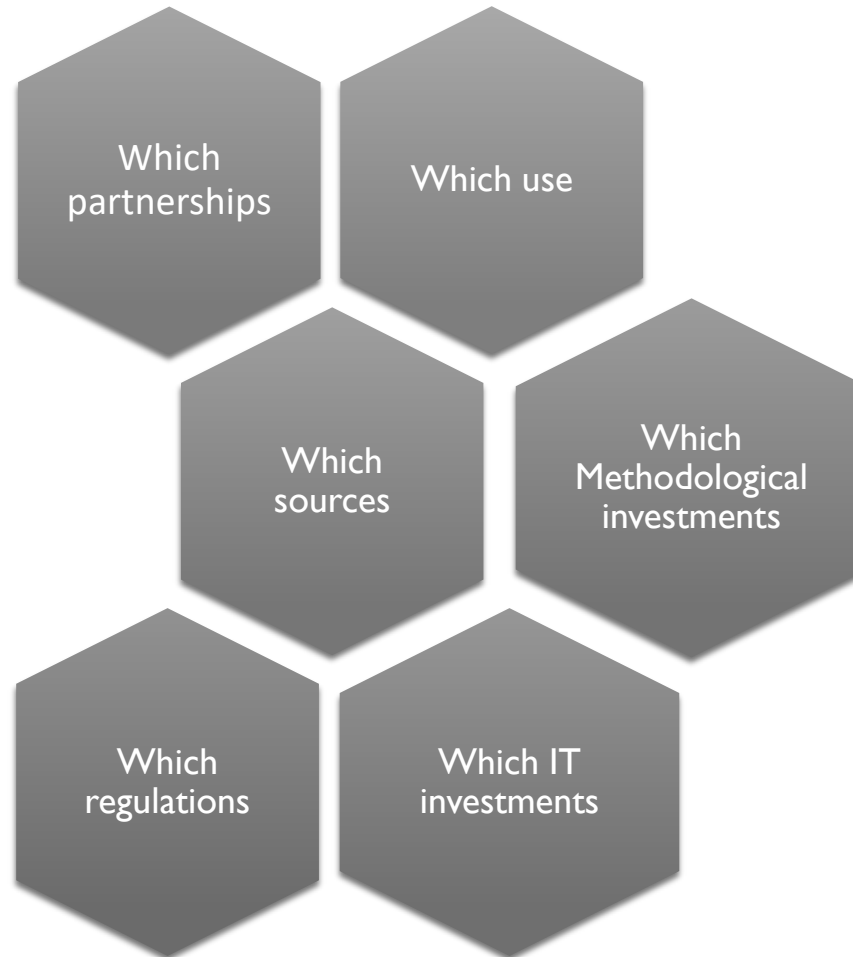Starting from 2013, Istat created a Big Data Committee responsible of the Big Data strategy...

October 2013
First Istat's Big Data Committee

April 2015
Roadmap for Big Data adoption in Istat 2016-2020

January 2016
Second Big Data Committee

2019-2020
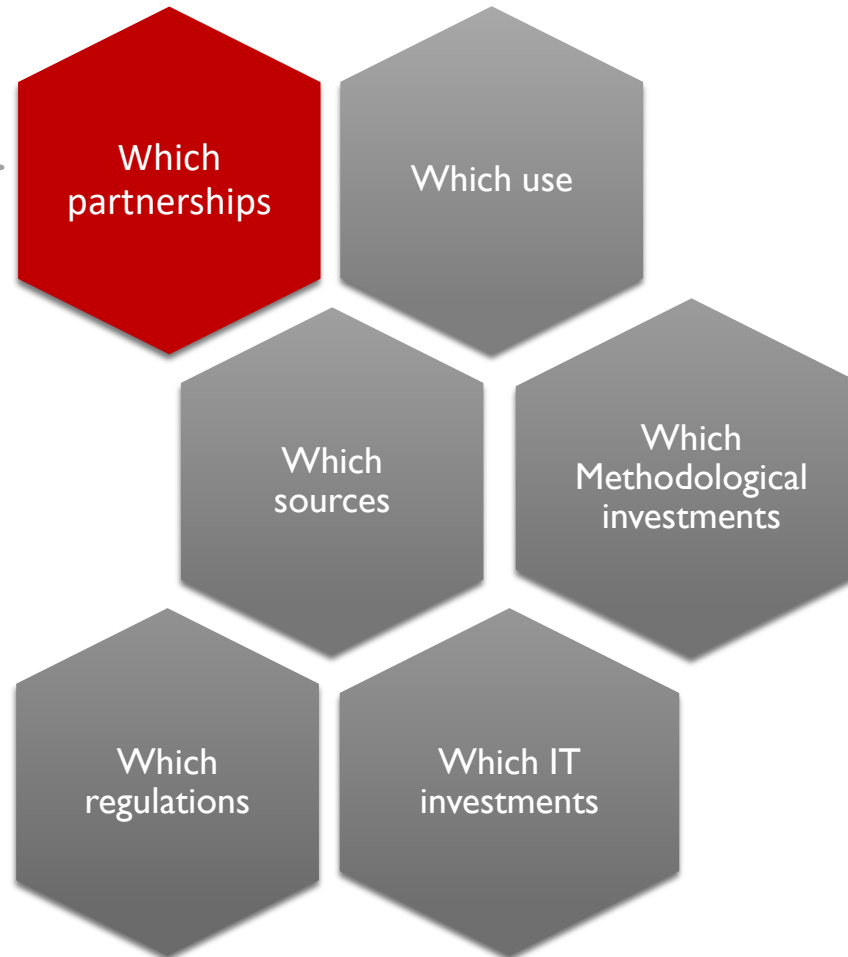Roadmap for Trusted Smart Statistics 2021-2024

...towards **Trusted Smart Statistics**

Istat

# Big Data @Istat: How – Operational Dimensions

# Big Data @Istat: How – Operational Dimensions

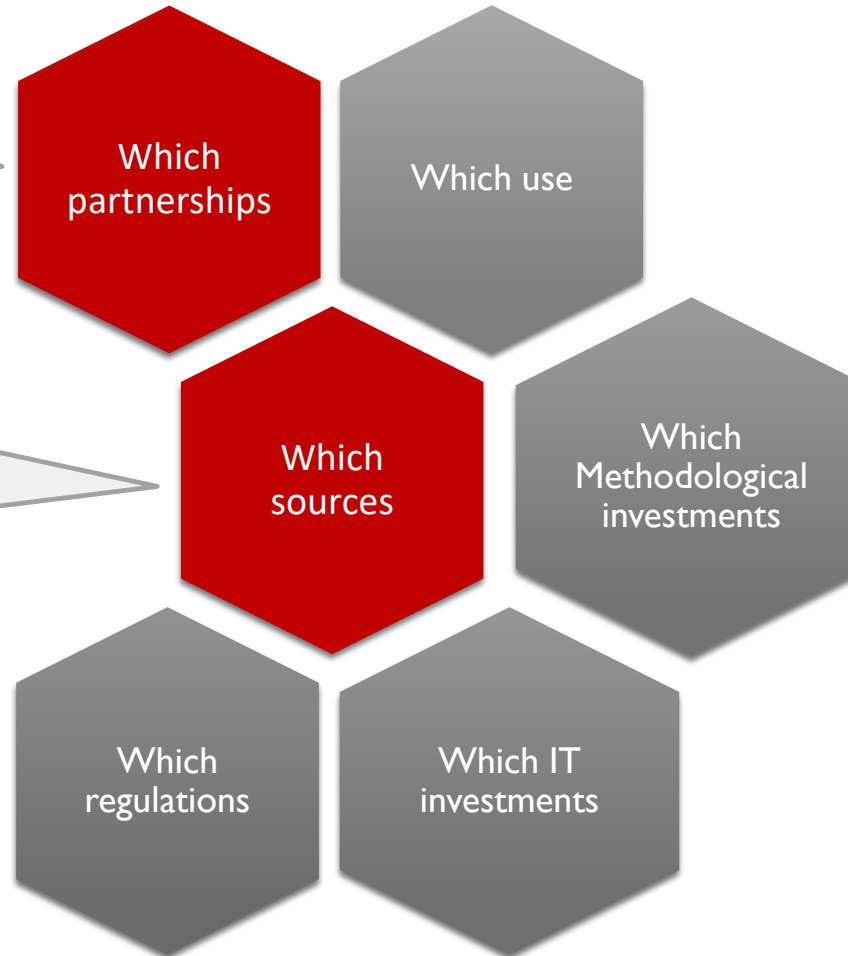- Mobile Network Operators
- Academic partnerships
- Private Companies (e.g. Nielsen)
- ...

Which partnerships

Which use

Which sources

Which Methodological investments

Which regulations
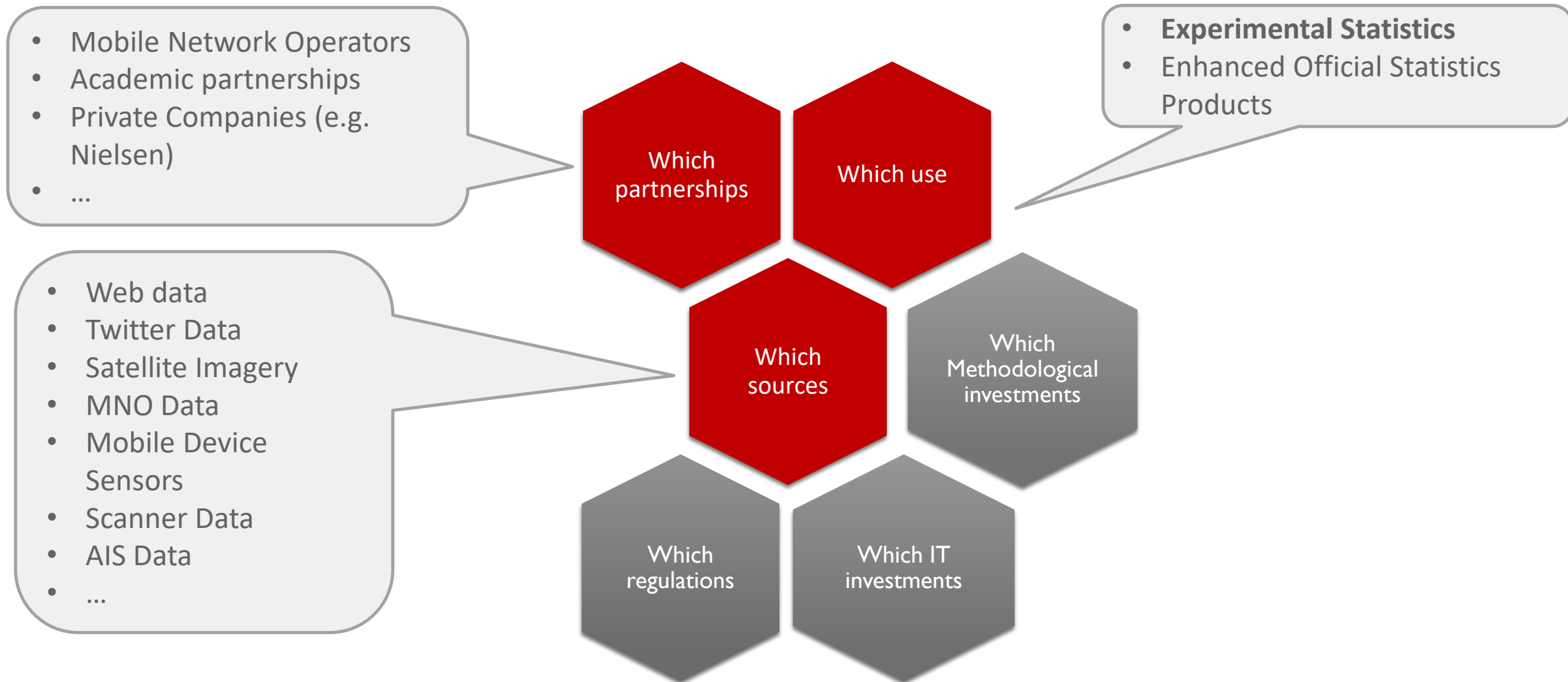
Which IT investments

Istat

# Big Data @Istat: How – Operational Dimensions

- Mobile Network Operators
- Academic partnerships
- Private Companies (e.g. Nielsen)
- ...

- Web data
- Twitter Data
- Satellite Imagery
- MNO Data
- Mobile Device Sensors
- Scanner Data
- AIS Data
- ...

Which partnerships

Which use

Which sources

Which Methodological investments

Which regulations

Which IT investments

Istat

# Big Data @Istat: How – Operational Dimensions

- Mobile Network Operators
- Academic partnerships
- Private Companies (e.g. Nielsen)
- ...

- **Experimental Statistics**
- Enhanced Official Statistics Products

**Which partnerships**

**Which use**

- Web data
- Twitter Data
- Satellite Imagery
- MNO Data
- Mobile Device Sensors
- Scanner Data
- AIS Data
- ...

**Which sources**

**Which Methodological investments**

**Which regulations**

**Which IT investments**

Istat

# Big Data @Istat: How – Operational Dimensions

- Mobile Network Operators
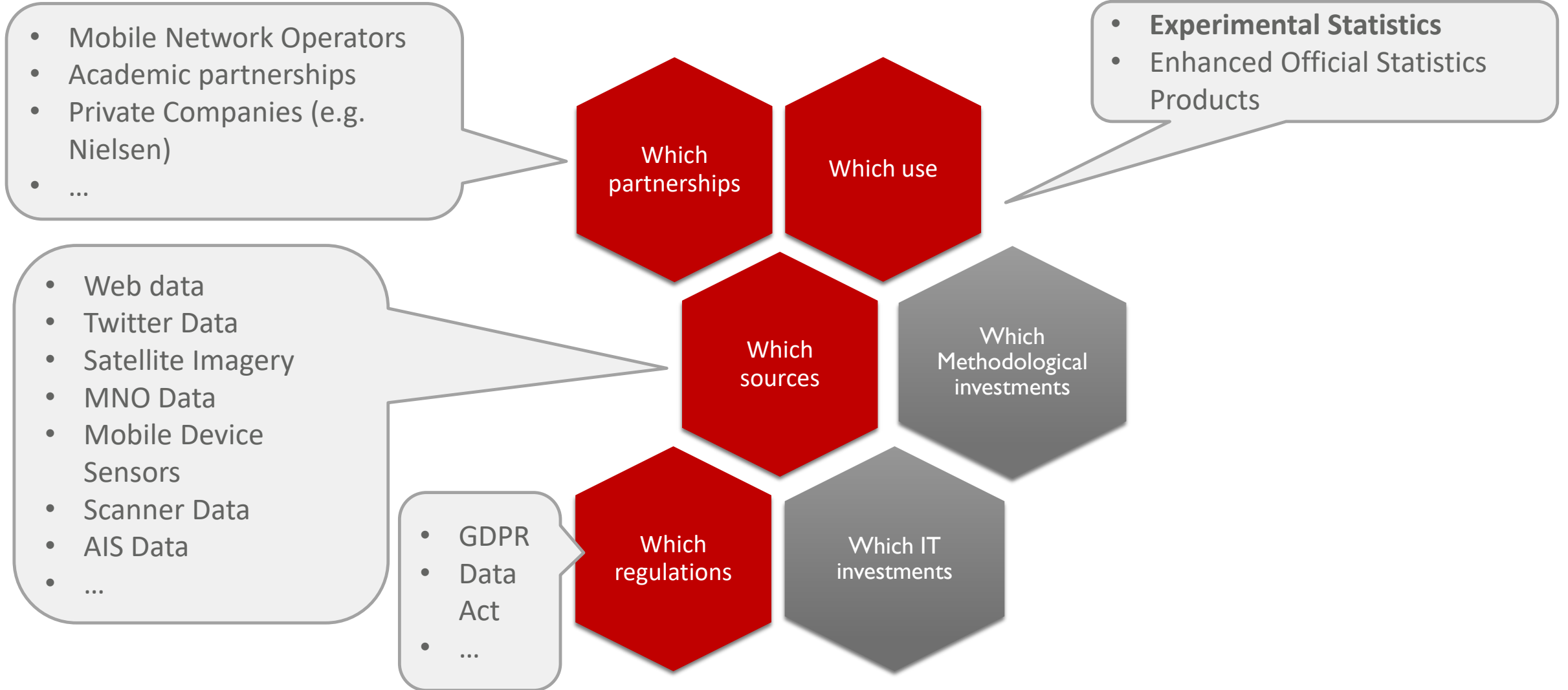- Academic partnerships
- Private Companies (e.g. Nielsen)
- ...

- **Experimental Statistics**
- Enhanced Official Statistics Products

**Which partnerships**

**Which use**

- Web data
- Twitter Data
- Satellite Imagery
- MNO Data
- Mobile Device Sensors
- Scanner Data
- AIS Data
- ...

**Which sources**

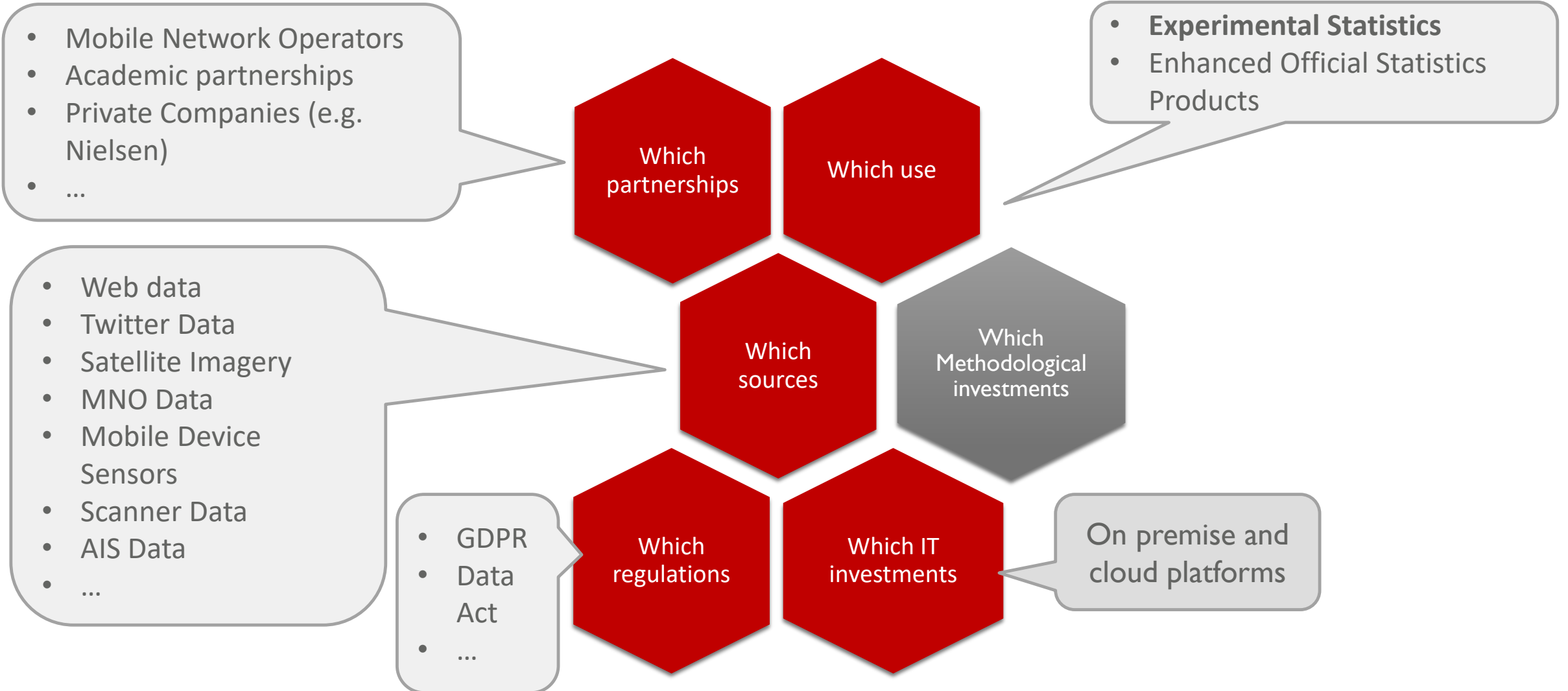**Which Methodological investments**

- GDPR
- Data Act
- ...

**Which regulations**

**Which IT investments**

Istat

# Big Data @Istat: How – Operational Dimensions

- Mobile Network Operators
- Academic partnerships
- Private Companies (e.g. Nielsen)
- ...

- **Experimental Statistics**
- Enhanced Official Statistics Products

**Which partnerships**

**Which use**

- Web data
- Twitter Data
- Satellite Imagery
- MNO Data
- Mobile Device Sensors
- Scanner Data
- AIS Data
- ...

**Which sources**

**Which Methodological investments**

- GDPR
- Data Act
- ...

**Which regulations**

**Which IT investments**

On premise and cloud platforms

Istat

# Big Data @Istat: How – Operational Dimensions

- Mobile Network Operators
- Academic partnerships
- Private Companies (e.g. Nielsen)
- ...

- **Experimental Statistics**
- Enhanced Official Statistics Products

**Which partnerships**

**Which use**

- Web data
- Twitter Data
- Satellite Imagery
- MNO Data
- Mobile Device Sensors
- Scanner Data
- AIS Data
- ...

**Which sources**

**Which Methodological investments**

- GDPR
- Data Act
- ...

**Which regulations**

**Which IT investments**

On premise and cloud platforms

Istat

# Big Data @Istat: What – Methodology

o **Dealing with Big Data heterogeneity**

  - New **data preparation** pipelines

    - Text

    - Images

  - New **inference** paradigm

    - **Machine learning**

Which Methodological investments

o **Dealing with access to external (private) sources**

  - Privacy preserving methods
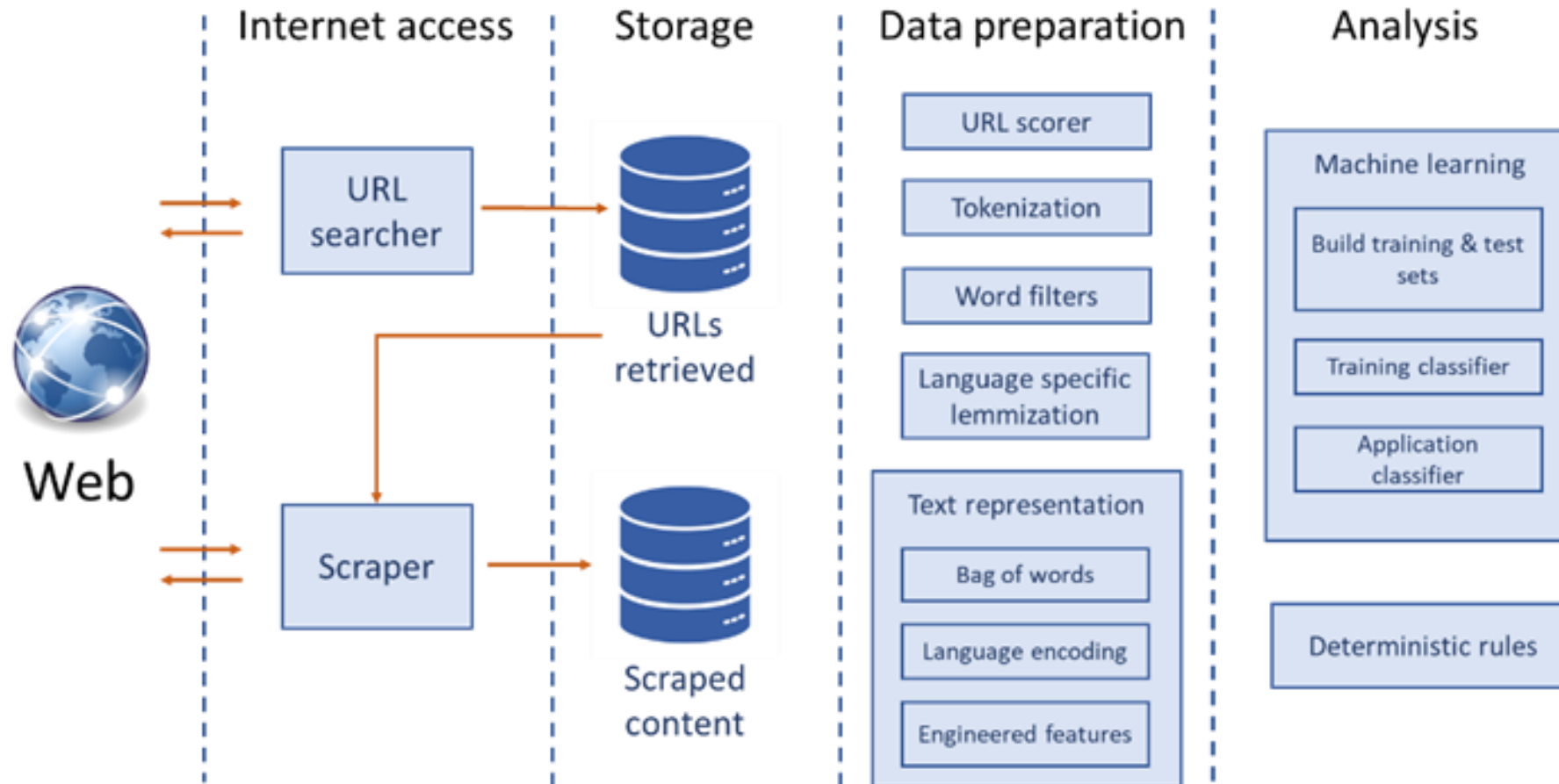
  - Web scraping

Istat

# Big Data Projects @Istat

# Text processing pipelines: Enterprise Characteristics

**GOALS:**

o The main goal of the project is the estimation of enterprises characteristics starting from the **web scraping of enterprises websites** (e.g., web ordering, job vacancy advertisement, link to social media)

o We implemented an algorithm that allows to predict enterprises characteristics (supervised machine learning model). The model is trained with survey data serving as training set for the machine learning task

o Simulations have demonstrated that these **new estimates are comparable to survey-based estimates**

Istat

# Text processing pipelines: Enterprise Characteristics

**Generic pipeline for processing textual data from enterprise websites:**

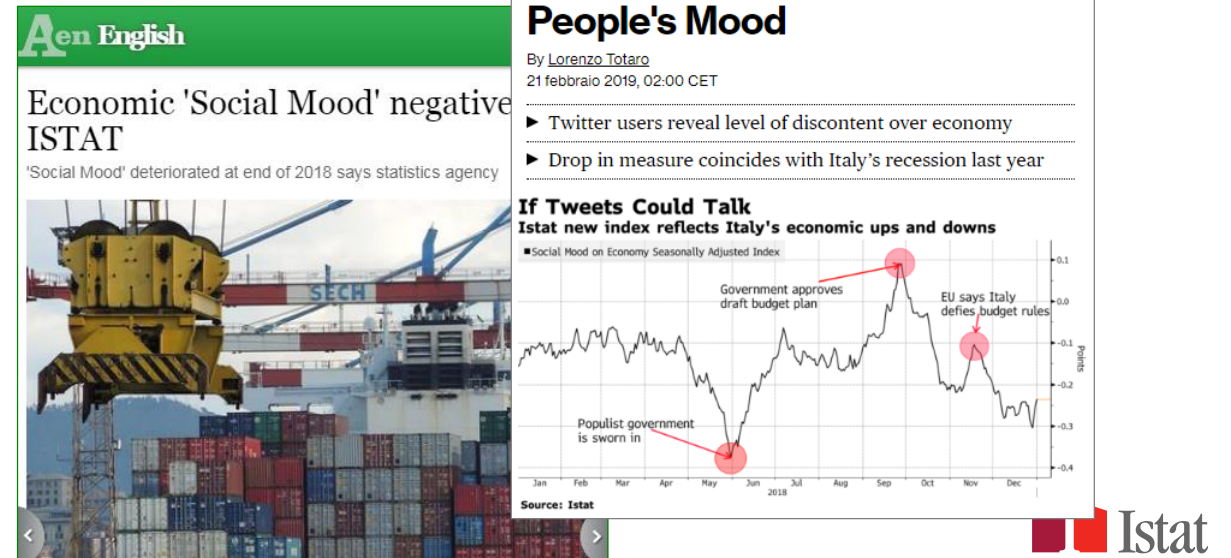# Text processing pipelines: Enterprise Characteristics

## Open challenges:

o The described project faces (some aspects of) the important issue of integrating a Big Data source with survey data. In the project, survey data are used as a training set of a Machine Learning classifier executed on Web extracted data

o Recently, in (Pratesi et al., 2022), more complex data integration methods are used to reduce the bias by combining a probability and a non-probability sample through a vector of common auxiliary variables, as an extension of (Kim & Wang, 2019).
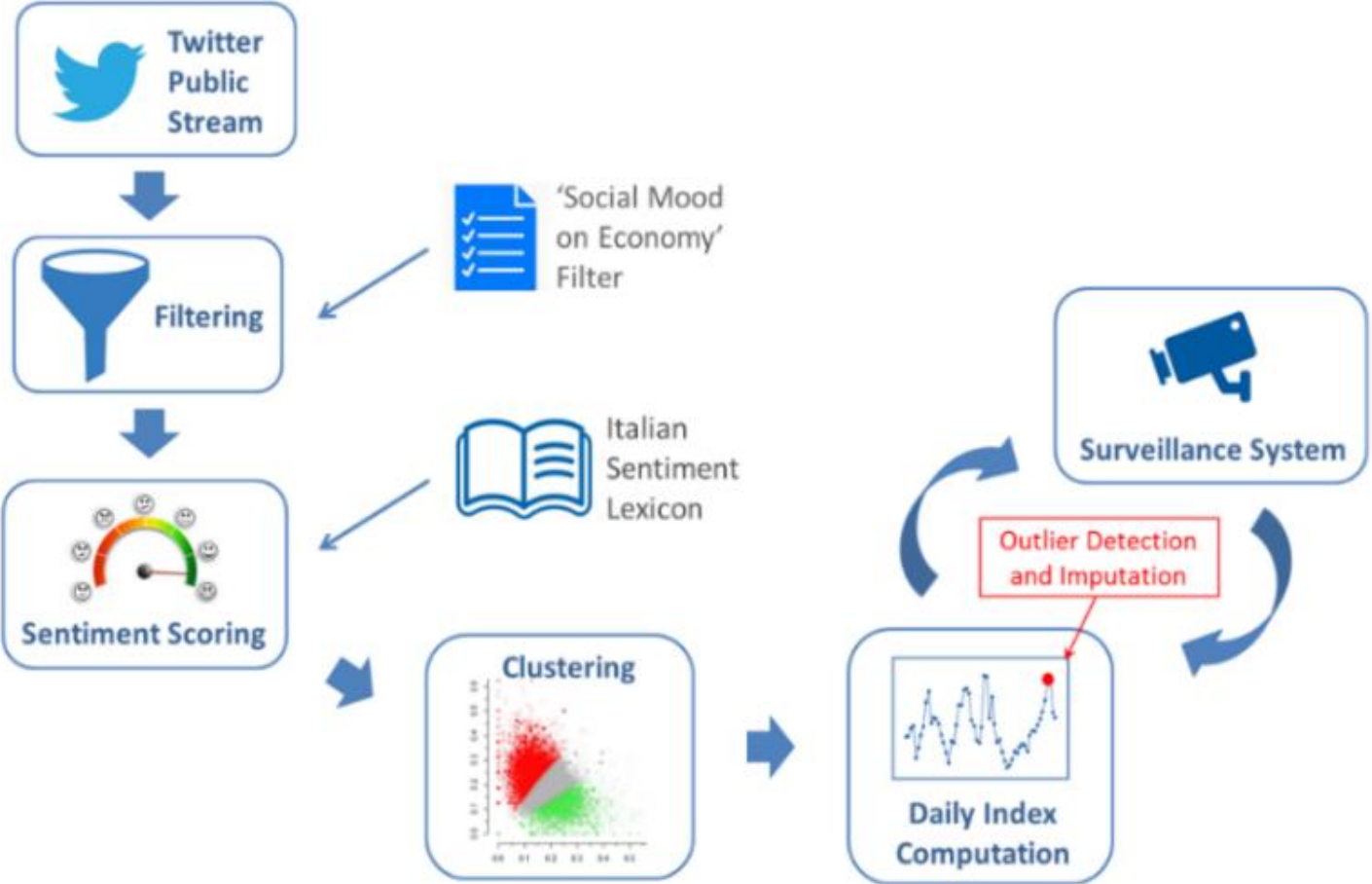
Istat

# Text processing pipelines: Social Mood

## GOALS:

○ The Social Mood on Economy Index (SMEI) is an **experimental statistic** published by Istat since 2018. It provides daily measures of the Italian sentiment on the economy, these measures derived from samples of public tweets in Italian language captured in real time

○ Data collection started in February 2016 and has been active since then almost without interruptions

• The dissemination of the new index attracted significant interest from the media (both traditional and online)

• Receptions were predominantly positive (the overwhelming majority praising Istat's openness to innovation)

• But few skeptical comments too!



Æn English

### Economic 'Social Mood' negative ISTAT

'Social Mood' deteriorated at end of 2018 says statistics agency



**Bloomberg** Subscribe

Economics

### Move Over GDP: Italy's Statistics Office Checks People's Mood

By Lorenzo Totaro
21 febbraio 2019, 02:00 CET

▶ Twitter users reveal level of discontent over economy

▶ Drop in measure coincides with Italy's recession last year

**If Tweets Could Talk**
Istat new index reflects Italy's economic ups and downs
■ Social Mood on Economy Seasonally Adjusted Index

Government approves
draft budget plan

EU says Italy
defies budget rules

Populist government
is sworn in

Jan Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec
2018

Source: Istat

# Text processing pipelines: Social Mood

## Pipeline to produce the Social Mood on Economy Index
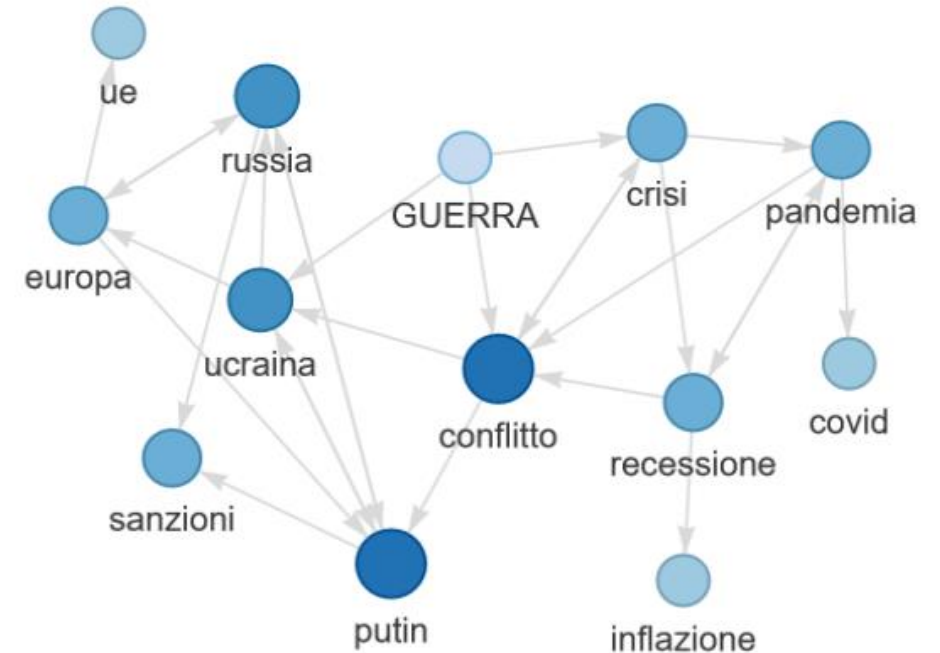
# Text processing pipelines: Social Mood

## Open challenges:

o In relation to SMEI, there are two major research directions that we would like to explore, namely:

  (i)   evaluation of the quality of Twitter's filters

  (ii)  improvement of the index interpretability

## WordEmBox

To evaluate the quality of the filter keywords we have exploited Word Embeddings (WE) methods. To this aim we used WordEmBox, an ad-hoc tool developed by Istat aiming at exploring WE spaces



Graph analysis for word "GUERRA" with the WordEmBox

# Image processing pipelines: Land Cover

**GOALS:**

Land Cover (LC) statistics and maps are a very important statistical product. As they require a big effort to be created, the idea is to build a set of algorithms to process satellite images in order to generate:

- Automatic Land Cover **Estimates**
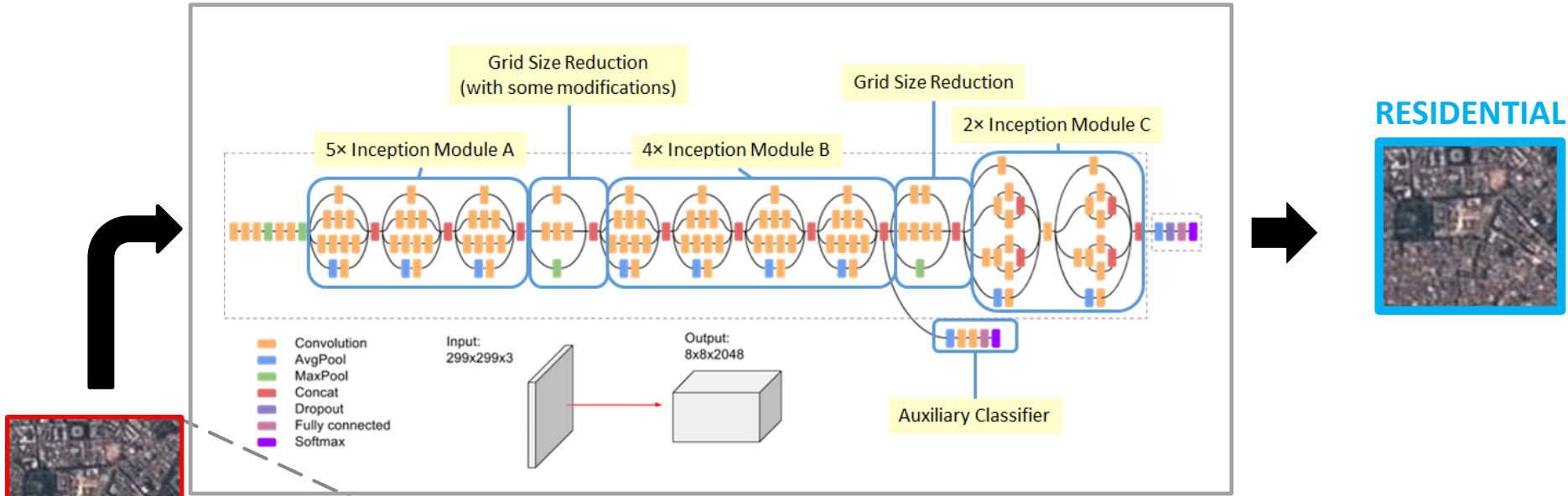- Automatic Land Cover **Maps**

**HOW:**

- **Standard approach**: Spectral Signature
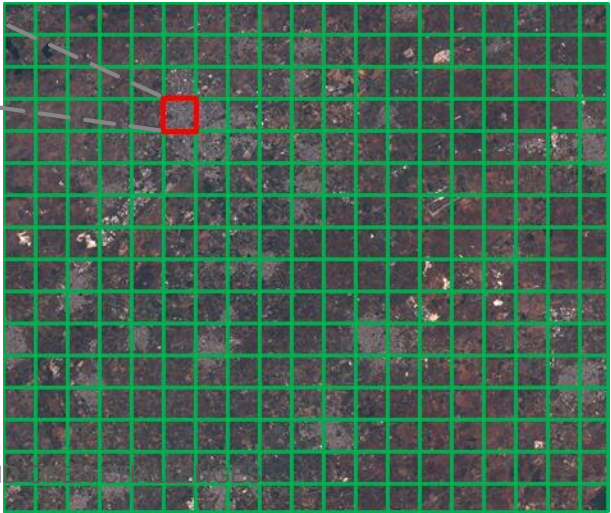- **New approach**: Using Deep Learning (**CNN** for classify-and-count and **U-Net** for segmentation)



**Standard approach**: Spectral Signature

# Image processing pipelines: Land Cover



RESIDENTIAL

Input Satellite Image

| LAND COVER CLASS | AREA SHARE |
|---|---|
| ... | ... |
| RESIDENTIAL | $\frac{45}{16 * 19} \cong 15\%$ |
| ... | ... |

# Image processing pipelines: Land Cover

**Results:**

The integrated architecture (CNN + U-Net) works very well for all LC classes:

o The U-Net takes care of LC classes "**River**" and "**Highway**"

o The CNN copes with all the other LC classes

o Partial LC maps produced by 1) and 2) are merged to yield a final complete LC map

**Open challenges:**

o One of the major problems in automated Land Cover (LC) estimation project is the **lack of a benchmark to validate the algorithm**, according to the chosen resolution and type of classification

o Suitable training dataset, compatible with the requested resolution for output. **Integration of input data with administrative sources** (e.g., data from regional technical charts, cadastral maps, and agricultural census)

Istat

# Improving data dissemination pipeline: TERRA

Istat's experience at European Big Data Hackathon, where we used different data sources traditional (**Comext Data**) and produced by sensors (**Google Mobility**) to analyze the impact of mobility restrictions on import / export



**TERRA - imporT ExpoRt netwoRk Analysis**
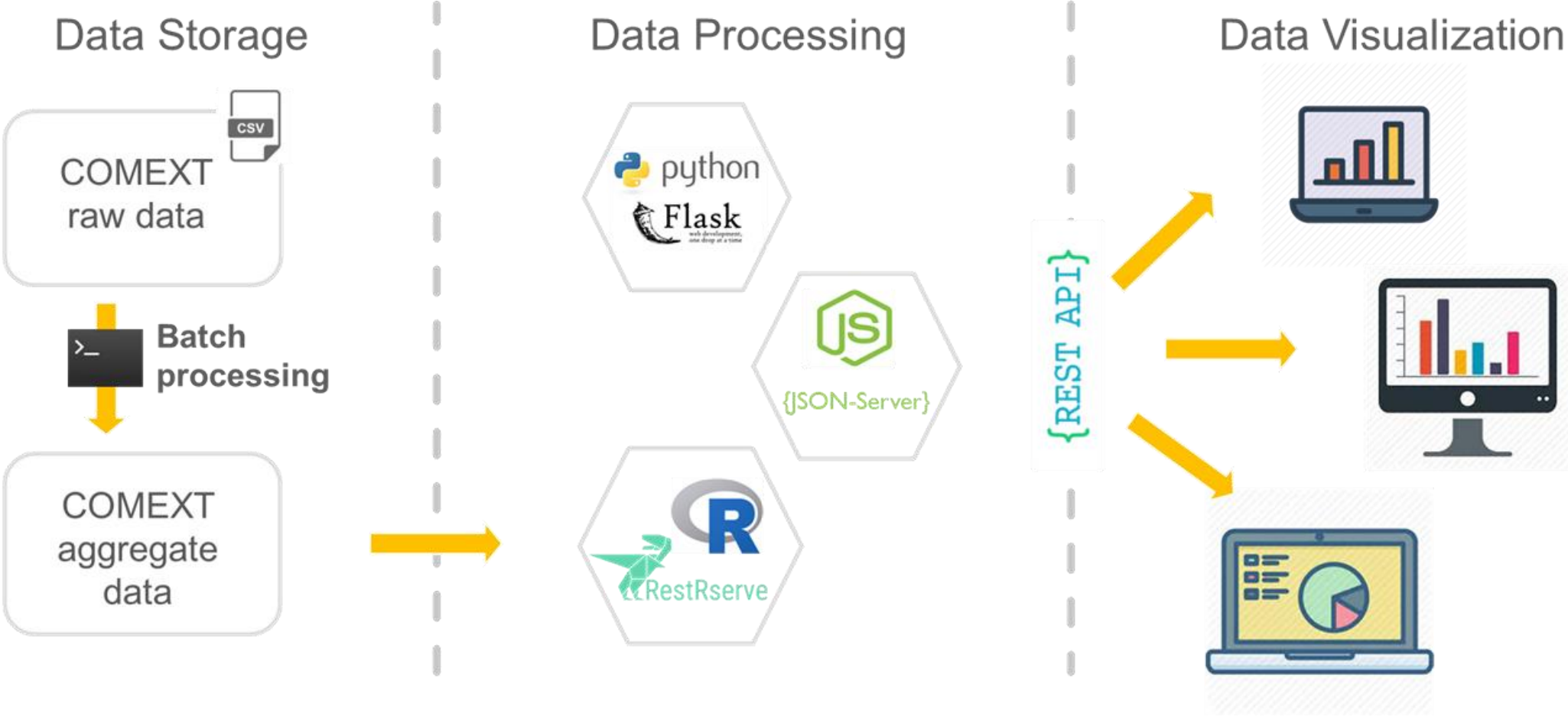
# Improving data dissemination pipeline: TERRA

## Big Data Hackathon 2021:

o Create an interactive dashboard to predict and simulate international trade relations in a high-resolution network by product and time

- o **Scenario analysis** and support for international trade policies
- o Ability to represent international global exchange networks
- o **Visualization of relationships** for partners, products and means of transportation
- o Analysis at product level as disaggregated as possible
- o **Scenario simulation for transport interruptions**
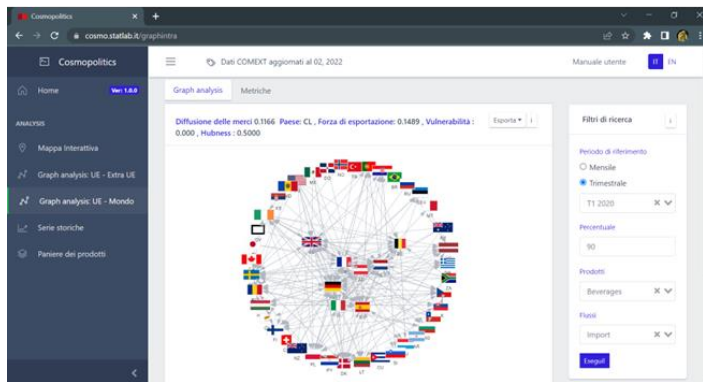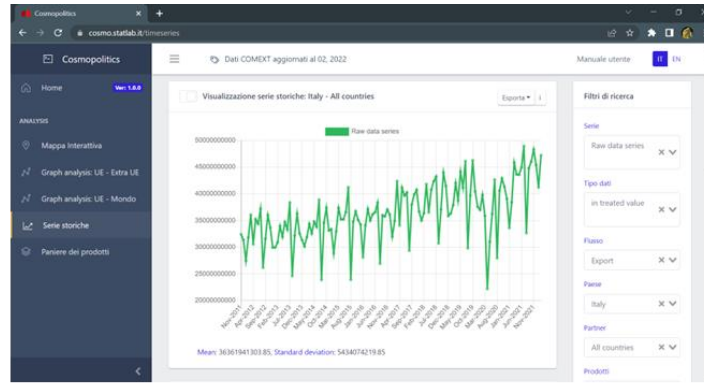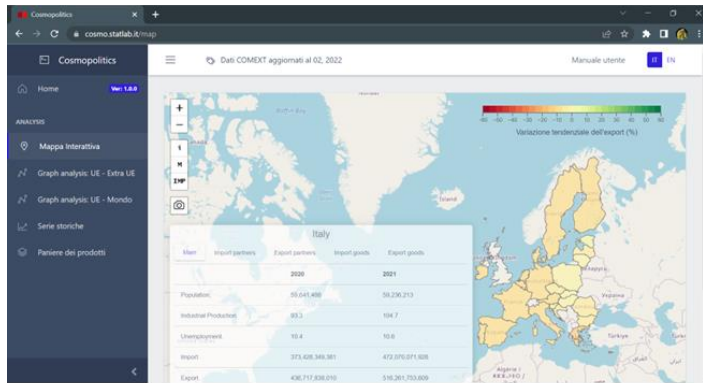- o Further suggested analysis: study of the **impact on COVID** products, etc.

# Improving data dissemination pipeline: TERRA

## TERRA in action:

# Improving data dissemination pipeline: TERRA

## TERRA is online (almost ready for production):



https://www.terra.statlab.it/

# Improving data dissemination pipeline: TERRA

## Open challenges:

o In its first version, the time series analysis section provided a forecasting up to 6 months a-head of the future international trade flows. This functionality used **Google COVID-19 Mobility Reports** to build a synthetic indicator capable of explaining the level of restriction imposed in each country using the principal component analysis methodology

o **Integration of new data sources (e.g., Stringency Index):**

o TERRA could provide new tools for performing scenario analysis. Indeed, the time series analysis section could be enriched by the inclusion of a subsection dedicated to one or more open indicators such as the Oxford University Stringency Index

Istat

# Concluding remarks...

# Methodological issues (recap)

**Open challenges:**

o [**Web scraping**] One of the most challenging aspects of the project concerns the issue of integrating a Big Data source with survey data. In the project, survey data are used as a training set of a Machine Learning classifier executed on Web extracted data

o [**Social Mood**] In relation to SMEI, there are two major research directions that we would like to explore, namely: 1) evaluation of the quality of Twitter's filters; 2) improvement of the index interpretability

o [**Land Cover**] One of the major problems in automated Land Cover (LC) estimation project is the lack of a benchmark to validate the algorithm, according to the chosen resolution and type of classification

# Conclusions

○ In the last 10 years huge investments on internal capacity building

○ Big Data and ML projects are **now a strategic asset in Istat**

  ○ **Experimental statistics are already there**

  ○ Production processes based on ML inference are planned in our 2021-2024 **Trusted Smart Statistics roadmap**

# Thanks!

MAURO BRUNO| mbruno@istat.it

Istat | Istituto Nazionale di Statistica