Instituto Nacional de Estatística
Statistics Portugal

eurostat

EUROPEAN CONFERENCE ON QUALITY IN OFFICIAL STATISTICS

2024 ESTORIL - PORTUGAL

# THE NEW ISTAT OPEN SOURCE AND STANDARDS BASED ARCHITECTURE FOR HIGH QUALITY WEB DISSEMINATION OF OFFICIAL STATISTICAL DATA

Mr Carlo Boselli, cboselli@istat.it

Mr Alessio Cardacino, alcardac@istat.it


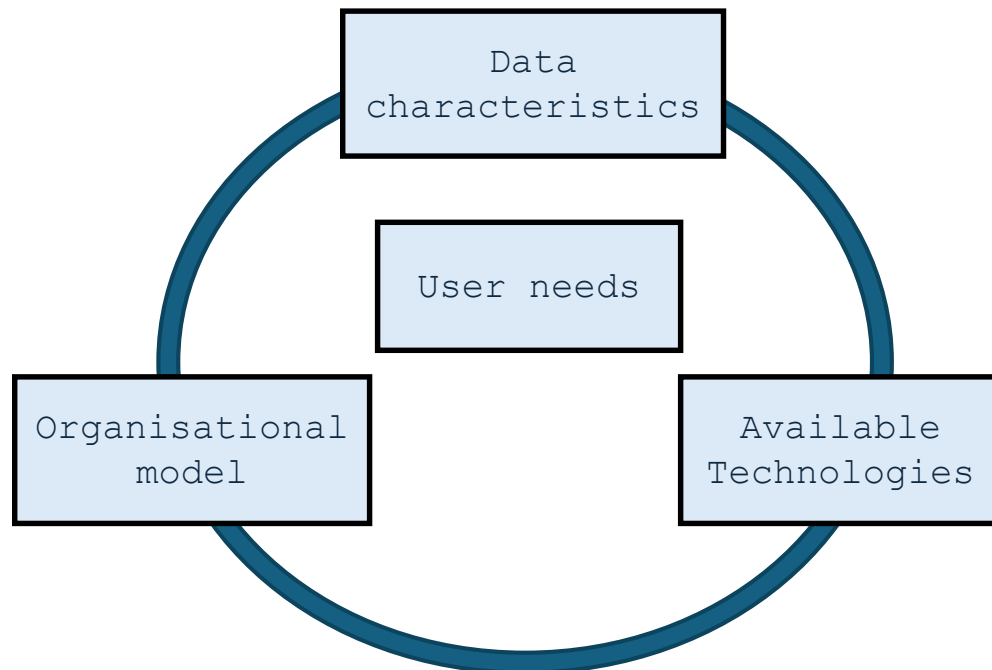ISTAT - Italian national institute of statistics, Italy

# Index

- The choice of a dissemination platform

- High level needs: Stability and flexibility

- Internal and external user needs

- Techonological solution: Istat SDMX Toolkit - Data Browser, Meta and Data Manager

- Data Browser features

- Meta Data Manager features

# The choice of a dissemination platform: Multidisciplinary approach

Different factors to consider before choosing a dissemination platform



Pre-existing complex system:

Different kinds of data: Admistrative data, survey data, micro and aggregate data

Different strategies and processes in the same organization

There is no better technology ever for data dissemination, but the one having the best balance between different aspects.
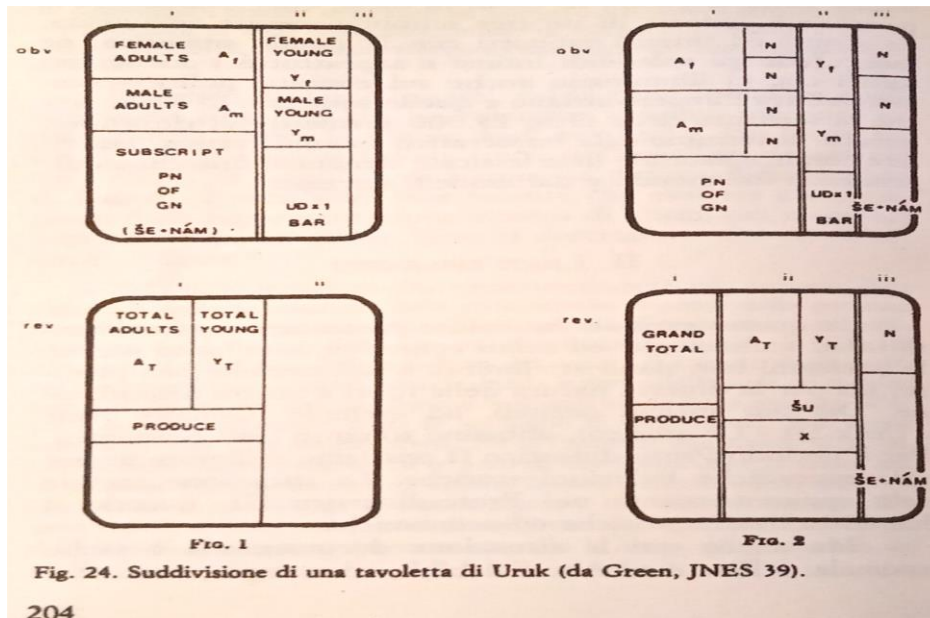
Technology must adapt to the objectives, the characteristics of the data and processes

**Multidisciplinary approach:**

Experts in business organization, statisticians, data modeling, IT

EUROPEAN CONFERENCE ON QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

INSTITUTO NACIONAL DE ESTATÍSTICA
Statistics Portugal

eurostat

The conference is partly financed by the European Union

# The choice of a dissemination platform: Multidisciplinary approach



Fig. 24. Suddivisione di una tavoletta di Uruk (da Green, JNES 39).

204

Structure of Uruk tablets for the dissemination of agricultural data. (M.W. Green, Journal of Near Eastern Studies 39 - Animal Husbandry at Uruk in the Archaic Period), in «I Sumeri» - Giovanni Pettinato - Bompiani

URUK III 4500 years ago, tablets for the dissemination of agricultural data:

Subdivision of animals by sex, by age, the animals born from the precedents, the related production of butter and cheese, those responsible for registration and administrative information

A long form and a short form of the tablets structure

Metadatation, data modelling, data comparison and linkage

Standardization and stability of the dissemination process

Some topics of data dissemination are a kind of archetype, involving more the users needs, data characteristics, the organizational system than the available technology:

**Consider context, user needs and technological solutions**

# High level needs

**Stability in reporting:**

National Statistical Institutes are required to respond to other organization in a reporting format. From this point of view, a protocol is recommended that stabilize data provisioning and reduce the entropy.

**Need:** Architecture based on a standard to guarantee data and metadata transfer via machine-to-machine.

**Flexibility in dissemination:**

Each National Statistical Institute has the mandate to collect data from institutional units and provide back to society aggregated data with a statistical value added, following a dissemination flexible approach. This process constantly redefine data structure, depending on new possibility to join data, new available technology, new needs from users.

**Need:** Tools useful to regain in terms of flexibility and guarantee the possibility, for example, to add new modalities in the codelists, update data structure definition (DSD), and in the future to modify the star schema and re-align all the depending artefacts reducing impact on external users.

EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

INSTITUTO NACIONAL DE ESTATÍSTICA
Statistics Portugal

eurostat

The conference is
partly financed by the
European Union

# SDMX – Statistical Data and Metadata eXchange

Migrate to a standard (SDMX) in order to:

- take advantage of technological modularity (changes or improvement in each component of the platform) without changes in the data modeling and preserving data dissemination processes in the future

- allow machine-to-machine access using API, and avoid misalignments between what is stored in DB and what is exposed via webservice

# Internal user needs

Specific context in Istat for a large-scale dissemination of aggregate data:

**Pre-existing processes:**

Dissemination of aggregate data from a large number of surveys under Eurostat Regulation. A large number of processes already existing of data aggregation and dissemination.

**Need:** The new platform must be linkable with pre-existing processes and allow data migration form legacy systems without altering production flows or increase burden on the production sectors and reuse of existing procedures for populating and updating data and metadata

**Administrative data:**

Dissemination of aggregate data from a combination of administrative and survey data that is possible to integrate, for specific domains of analysis, to obtain a more granular data dissemination.

**Need:** High data storage capacity per single data cube. In particular, for the population and foreign trade aggregate data, there are datasets with more than 60 million of data points, and it is necessary to store, query and browse them on the web as well

# External user needs

- Homogeneous environment to reduce access costs
- Possibility to publish different nodes in the same instance
- Synthetic data representation, as dashboards
- Possibility to visualize tables with an high number of data points (at least one million)
- High performance in data retrieving and visualization using an innovative cache system
- Text search system that allows the use of artificial intelligence in the future
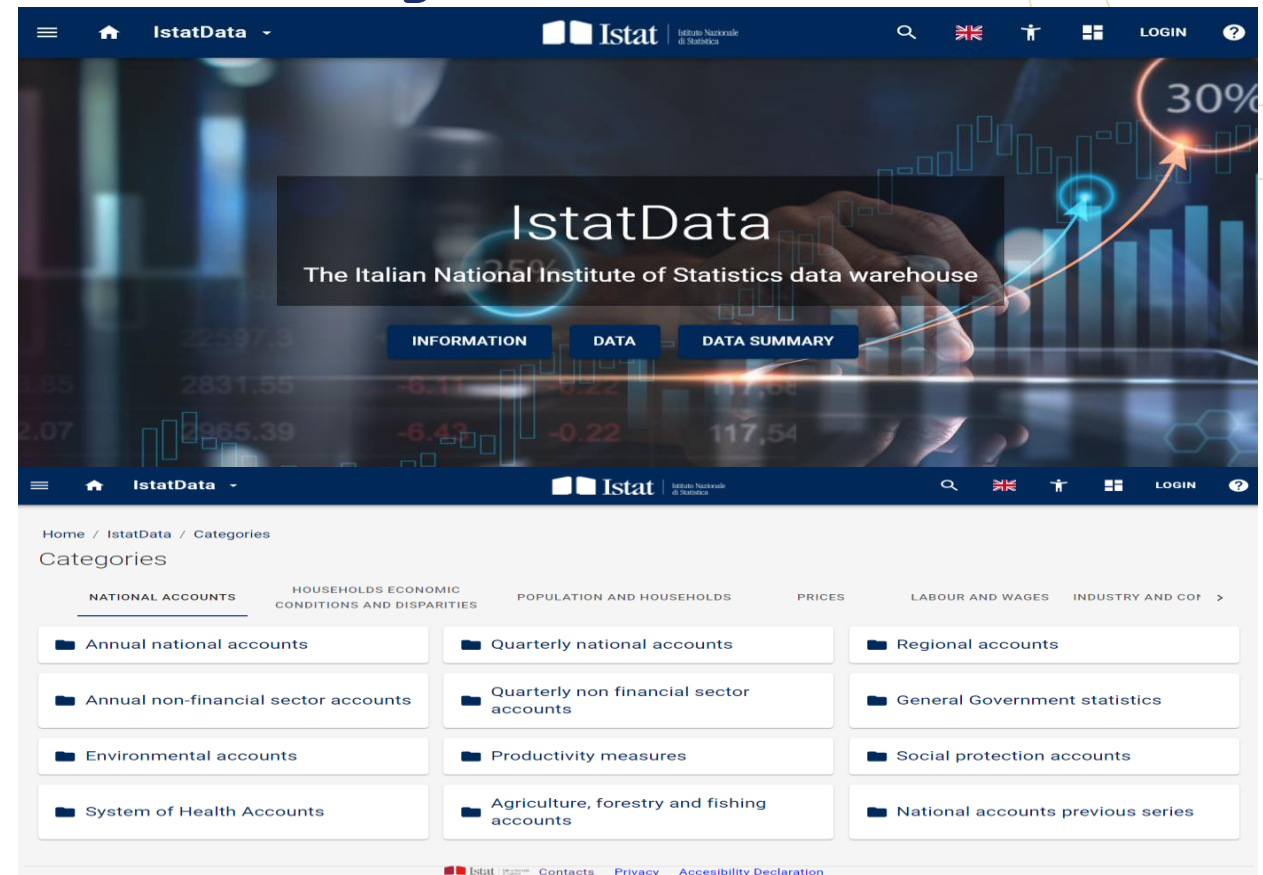- An ease installation and configuration of the tools

EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

INSTITUTO NACIONAL DE ESTATÍSTICA
Statistics Portugal

eurostat

The conference is
partly financed by the
European Union

# The techonological solution: SDMX Istat Toolkit – Data Browser, Meta and Data Manager

Istat SDMX Toolkit (https://sdmxistattoolkit.github.io) is an open source platform developed by Istat based on two modules:

- Data Browser (Front end): to allow external users to browse data

- Meta & Data Manager (Back end): to create SDMX data structure, data cubes, data mapping, data loading and data flows (for Istat internal users)

Front and back-end technologies are based on the SDMX standard.

Starting from 2021, IstatData (https:\\esploradati.istat.it) is the new Istat corporate data warehouse for aggregate data dissemination

EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

INSTITUTO NACIONAL DE ESTATÍSTICA
Statistics Portugal

eurostat

The conference is
partly financed by the
European Union

# Different dissemination systems towards a new single platform

**Istat Legacy systems:**

- I.Stat (800.000 users per year)

- Permanent census

- COEWEB (data of foreign trade)
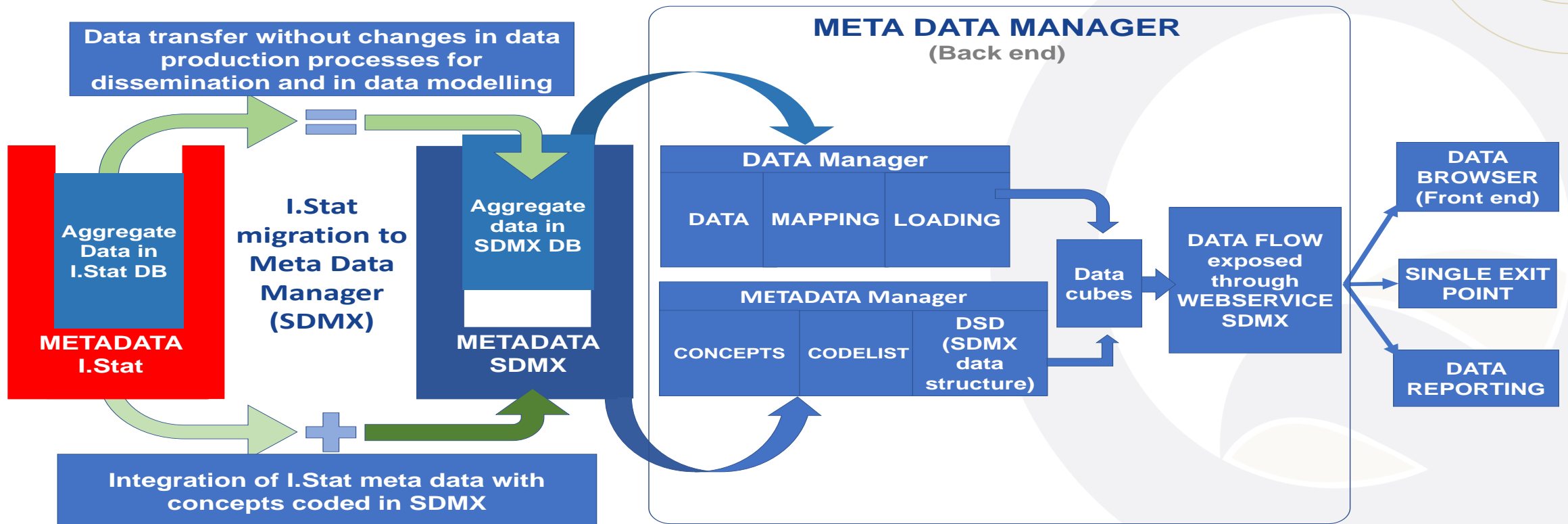
- Other data bases storing aggregate data

**Istat new aggregate data dissemination systems:**

- IstatData (https:\\esploradati.istat.it) (2 MLD of data points, 450 datacubes, 3000 dataflows)

- Permanent census of population and housing (https://esploradati.censimentopopolazione.istat.it) (800 MLN records of data points)

- COEWEB (data of foreign trade) (https:\\esploradati.istat.it\coeweb) (10 MLD of data points, 40 datacubes)

# Meta And Data Manager Tools

Set of tools useful in data modelling and data migration from legacy systems
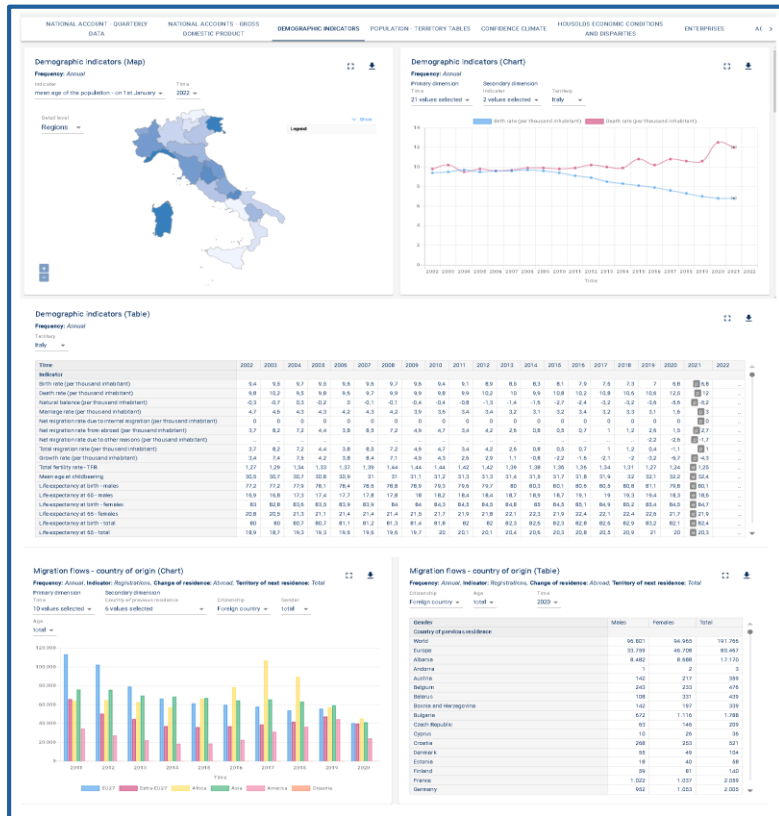
# Meta And Data Manager main features

- Modelling of statistical tables through the definition of the relevant metadata, which specifies the role that each variable plays in the structure of the tables, as well as their representation through appropriate lists of codes or classifications. Creation and publication of descriptive metadata (conceptual, methodological, quality)

- Possibility to download and upload each artefact (as Data Structure Definition DSD, codelists, concept scheme) in different format, SDMX 21., SDMX 2.0 XML

- Creation and management of statistical databases organized as multidimensional cubes in which to insert and update data starting from CSV, Excel and XML files

- Display of multidimensional statistical tables through a specific API according to the SDMX standard

- Creation and publication of catalogs of the datasets that you intend to disseminate. These catalogs are published using "DCAT vocabulary" according to the RDF standard for Linked Open Data. In particular, as default, everything necessary to implement the DCAT-AP_IT application profile created by AgID (Italian agency for digitalization) is provided;

- Creation, management and publication of thematic glossaries

- Ability to manage data structure definition (DSD), codelists and concepts through an intuitive interface that allows the internal user to easily edit, modify, upload or clone for easy reuse of the artefacts that allows you to create new versions without repeating all the steps

- Ability to add modalities in codelists in a specific DSD and datacube, without modifying the versioning, guaranteeing external users a reuse of the pre-existing dataflows

- Ability to add new attributes, with new footnotes or multiple encodings associated for a specific dimension of analysis, directly updating a pre-existing DSD, without destroying data cubes, dataflows and other depending artefacts

- Using Meta-Data Handler tools, ability to make available data, already published according to a specific data structure, on the base of a different DSD defined by another international organization for reporting purposes

- Ability to upload pre-existing csv files, with different formats, using file-mapping tools. This functionality allows internal user to upload in the Meta and Data Manager, without changes in the format, files already defined from production units to populate legacy dissemination systems, reducing the burden during data migration

# DATA BROWSER: Data comparison and synthetic prospects (Dashboard)



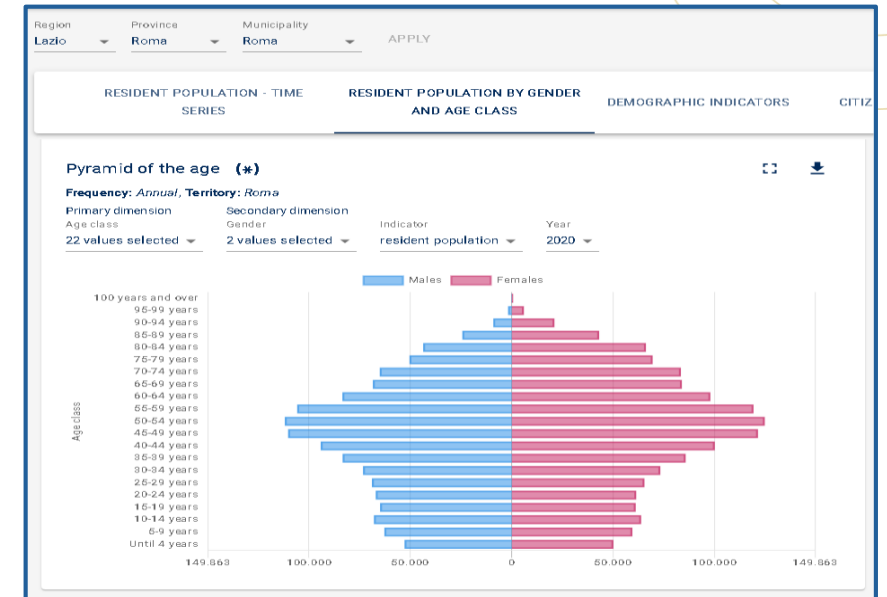Combination of maps, charts tables and explanatory text

Objects created by combining different dataflow from different datasets and different nodes

Interactive selection

Full screen enlargement for each object

Data export in different format

Image export for charts and maps



Territorial data selection system at regional, provincial and municipality level: all the dashboard data automatically updated for a specific territory

# Data Browser main features

- This application can be used within a single organization in order to disseminate datasets stored in one or more databases (e.g. https://esploradati.istat.it), or within a distributed architecture on web servers reachable via the http protocol on the Internet (e.g. https://idp.sister.it).Possibility to download and upload each artefact (as Data Structure Definition DSD, codelists, concept scheme) in different format, SDMX 21., SDMX 2.0 XML

- It is possible to browse datasets in the format of multidimensional tables, associated with various graph modes, thematic maps, and customize the layout of all forms of visualization. It is also possible to download data in various formats, including open data such as XML (SDMX-ML), SDMX-CSV, SDMX-JSON, CSV (custom), Excel, image and pdf (for graphs and maps).

- For the aspects related to performances, the module allows browsing a dataset with more than 300 million of data points, visualizing on the web a table with more than one million of data points per time and providing search criteria to change the selection to browse all stored data.

- The system is also equipped with a sophisticated cache, in order to make queries and data visualizations highly efficient even in the presence of particularly large datasets. This solution allows storing and reusing, for subsequent queries, the results of all the queries carried out by users. This system is consistent with updating criteria of the datasets themselves, which invalidate the caches previously created.

- Other features are the possibility to browse and search datasets within each single data source by typing keywords or directly selecting the dataset of interest with the help of specific hierarchical "theme trees", and the possibility for registered users to save their data views within a portfolio that can be consulted once logged into the system.

EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

INSTITUTO NACIONAL DE ESTATÍSTICA
Statistics Portugal

eurostat

The conference is
partly financed by the
European Union

# Webservice and single exit point

Architecture natively based on an internationally recognized data transfer standard (SDMX)
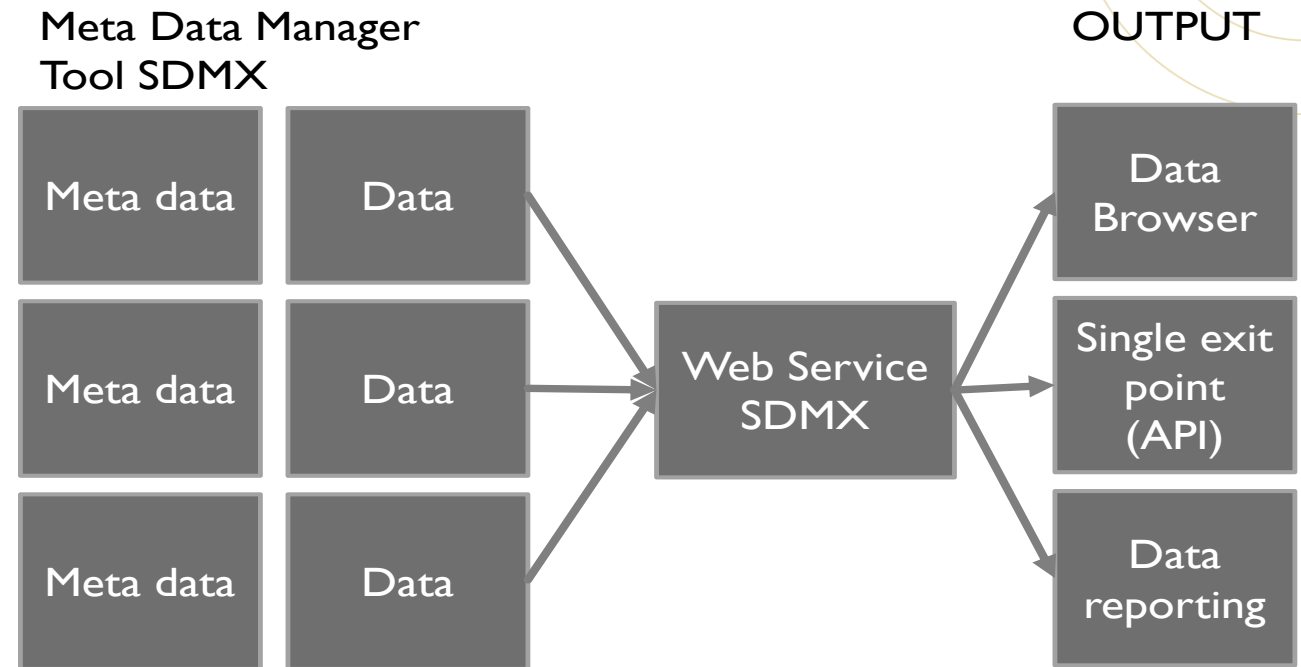
- High data store capacity per single datacube

- Modular infrastructure with access via webservice

- New technological framework is developed by ISTAT and completely Open Source

IT reference colleagues for technological aspects:

Francesco Rizzo (rizzo@istat.it)

Alessio Cardacino (alcardac@istat.it)

Simone Coccia (sicoccia@istat.it)

**Meta Data Manager Tool SDMX**

**OUTPUT**

| Meta data | Data |
| Meta data | Data |
| Meta data | Data |

→ Web Service SDMX →

- Data Browser
- Single exit point (API)
- Data reporting

https://github.com/SDMXISTATTOOLKIT