Webinar on Building a Central Statistical Data Warehouse
Doha, 17 May 2022

حلقة عمل عن بعد بشأن بناء قاعدة بيانات إحصائية مركزية
الدوحة، 17 مايو 2022

Developing a statistical Data warehouse
17 May 2022

# The Statistical Dissemination Data Warehouse in Istat

Domenico Fedele, Francesco Rizzo

Istat

# Summary

- Standardisation activities in Istat
- Why supporting the dissemination business process with a Statistical DWH
- Dissemination DWH implementation steps and timeline
- Dissemination DWH architecture
- Distributed Dissemination DWH within the National Statistical System
- Dissemination DWH: the free and open source toolkit
- Design and implementation principles of the Toolkit
- Processing – Cubus Toolkit
- Transformation – Excel2CSV tool
- Meta and Data Manager - main features, architectures
- Data Browser - main features
- Lesson learnt

Istat

# Standardisation activities in ISTAT

ISTAT has been running a modernisation program based on the guidelines of the UNECE "High-level Group for the Modernisation of Statistical Production and Services" in order to:

❑ Satisfy new demands for statistical information:

    ❑ Produce more statistical indicators with more sectorial and territorial details

    ❑ Improve quality (*Coherence, Timeliness, Comparability, Accessibility*)

    ❑ Provide Governments needs to help the formulation of good policy, not only at national level

❑ Streamline the expenses due to a reduction of the financial allocation

    ❑ Leveraging the ICT development that has allowed a cost reduction in producing statistics and its easily accessibility and dissemination

❑ Seizing the opportunities that Internet is offering in terms of new information sources and new ways of combining and using information
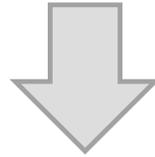
<span style="color:red">Modernisation = Standardisation + Industrialisation</span>

Istat

# Why supporting the dissemination business process with a Statistical DWH

❑ Facilitating data integration and process integration

❑ Providing a streamlined infrastructure for the standardisation and industrialisation of the data life cycle

❑ Driving the (data-centric) workflow through standardised information objects between different processes (metadata-drive approach)

❑ Increasing the quality of the statistics by reducing manual transformation steps with automated and stable processes

❑ Reducing the files dispersion on the local disks of statisticians, offering a dedicated collaborative infrastructure using databases

❑ Overcoming the issues related to dated technologies of the existing legacy systems by implementing a new services based IT platform in a loosely coupled architecture

❑ Offering the same platform to the Organisations part of the National Statistical System (implementing a distributed DWH) in order to improve quality in dissemination

Istat

# Dissemination DWH implementation steps and timeline

❑ Analysis of different options: commercial solutions, tools available within the statistical community, evolution of a prototype developed in-house
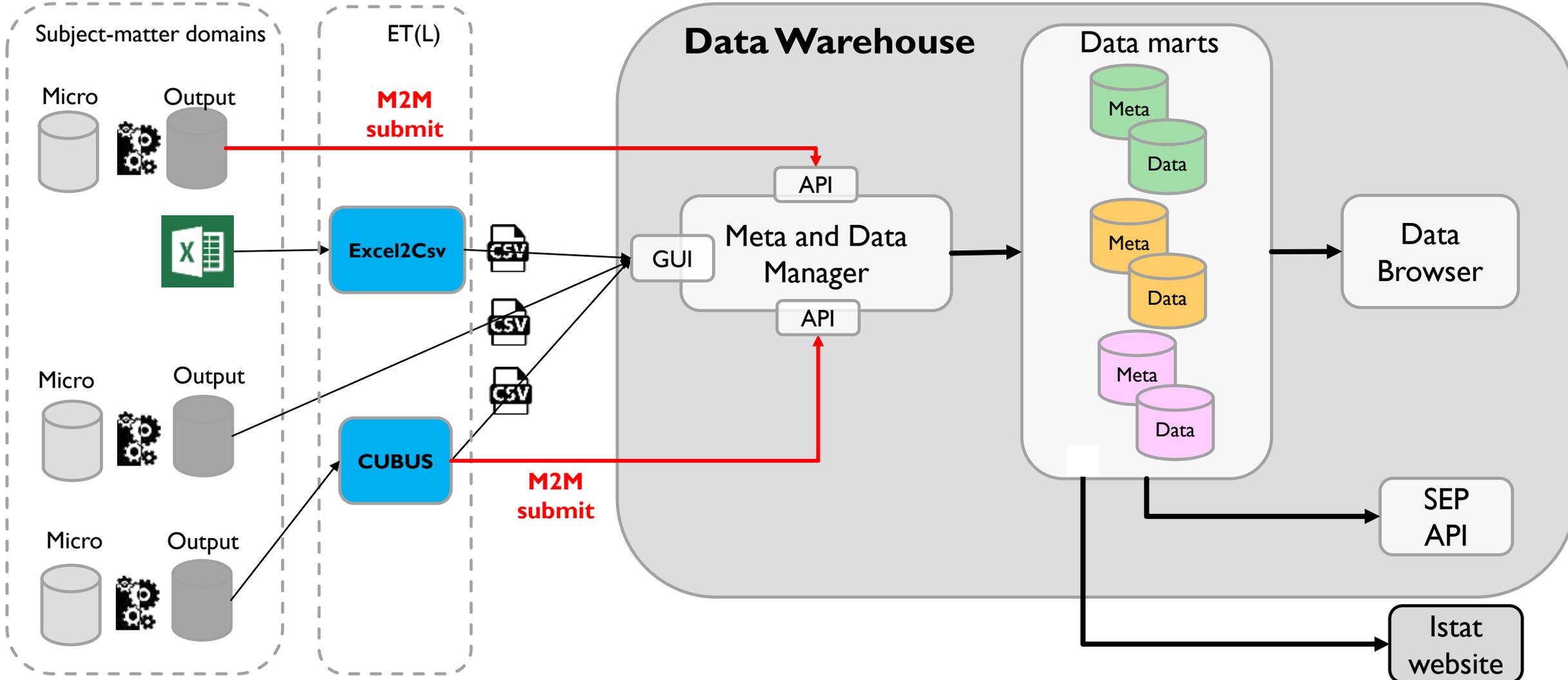
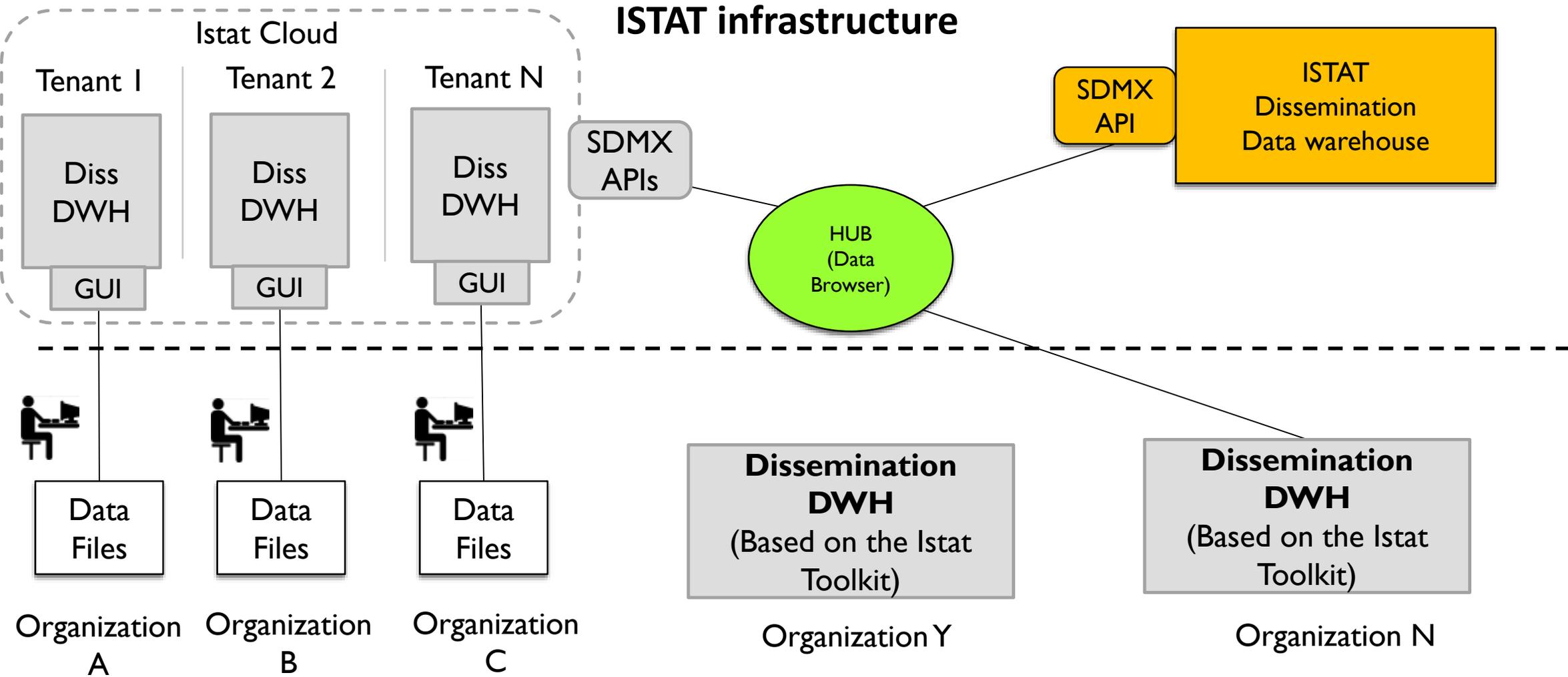**evolve the prototype and complement it with tools developed in other context**

❑ Design a suitable architecture

❑ organize the needed building blocks as a reusable Toolkit (free and open source)

❑ Plan migration from legacy systems to the new DWH platform

# Dissemination Data Warehouse architecture in Istat

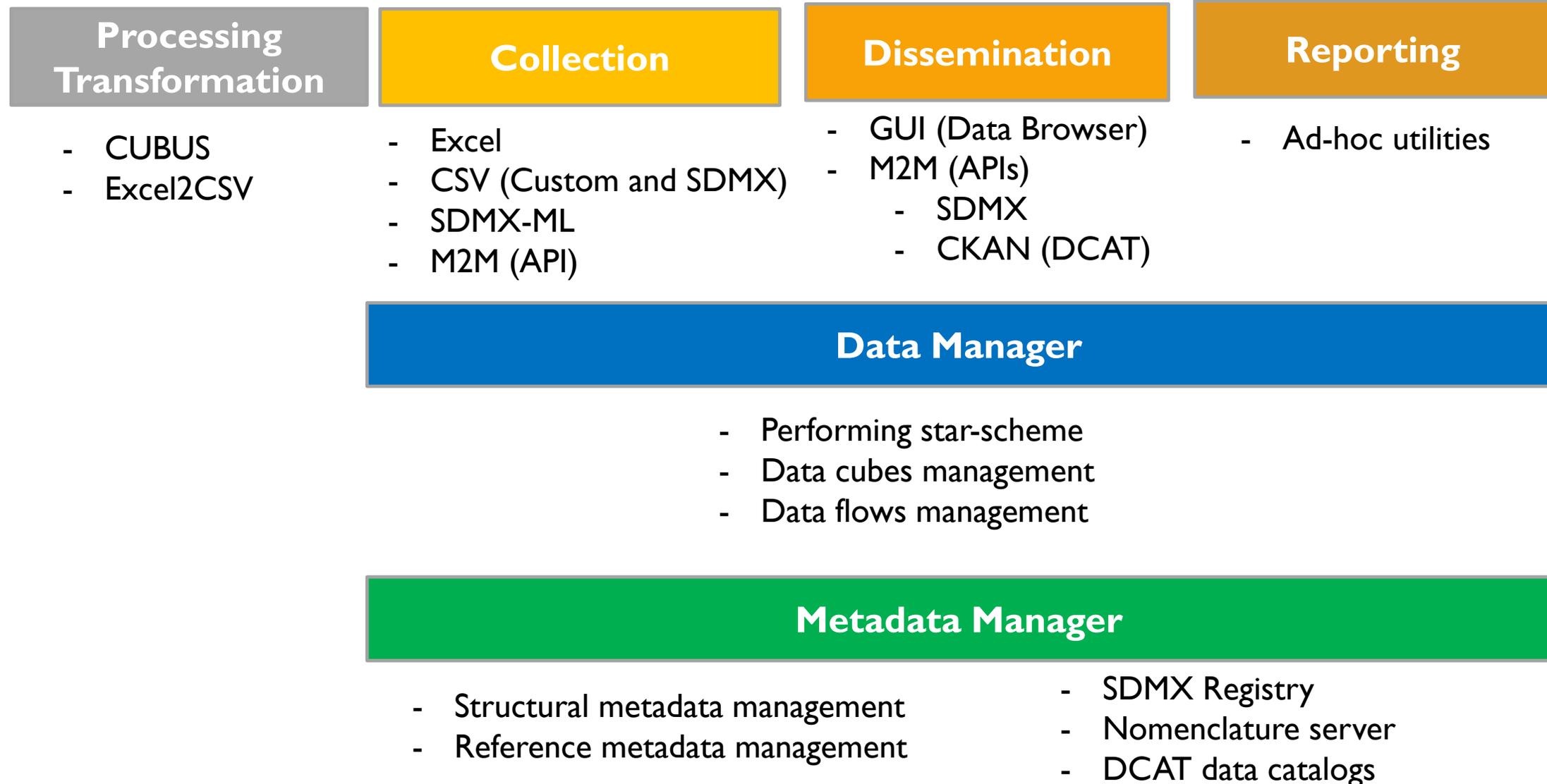# Distributed Dissemination DWH within the National Statistical System

**ISTAT infrastructure**

Istat Cloud

Tenant 1 | Tenant 2 | Tenant N

Diss DWH | Diss DWH | Diss DWH

GUI | GUI | GUI

SDMX APIs

SDMX API

**ISTAT Dissemination Data warehouse**

HUB (Data Browser)

Data Files | Data Files | Data Files

Organization A | Organization B | Organization C

**Dissemination DWH** (Based on the Istat Toolkit)

Organization Y

**Dissemination DWH** (Based on the Istat Toolkit)

Organization N

**National Statistical System Organizations premises**

Istat

# Dissemination DWH

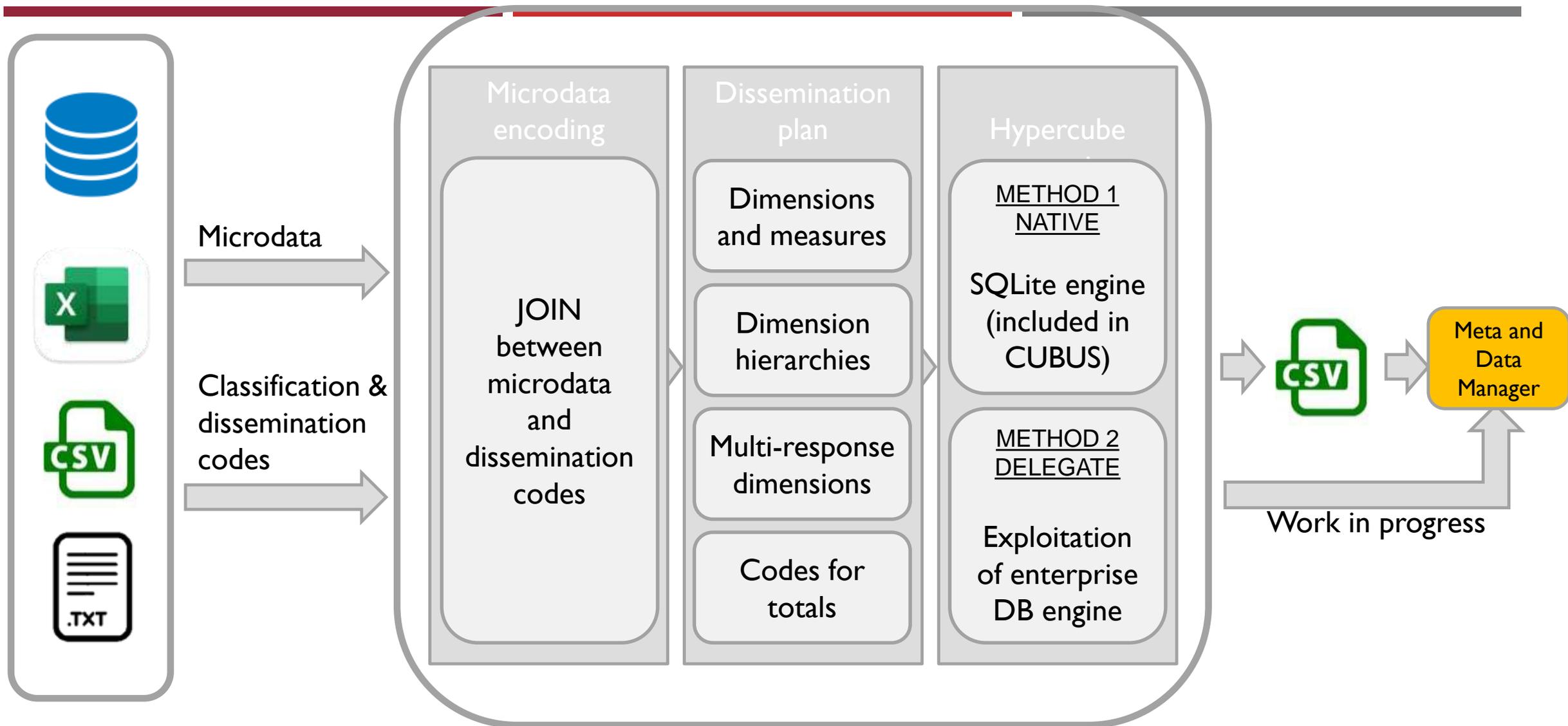# The free and open source toolkit

Istat

# Design and implementation principles of the Toolkit

- ❑ **Reengineering by experience**
  - ▪ Old version used as prototype
  - ▪ Feedbacks collected from +500 users (trainings and national and international projects)
- ❑ **Standards-based**
  - ▪ SDMX, DCAT
- ❑ **Easy to install (Plug and Play)**
  - ▪ Not more than 20 min
  - ▪ Technical skills easy to find in statistical organizations
- ❑ **Easy to configure**
  - ▪ Default configuration (Play)
  - ▪ Extended configuration through GUI
- ❑ **Easy to change the brand and to localize in different languages**
  - ▪ CSS files
  - ▪ resource files
- ❑ **Adequate performance**
  - ▪ Data rendering of a multidimensional table with 100.000 cells in less than 10 seconds
  - ▪ Benchmark with the legacy Istat dissemination I.Stat, and other tools available in the statistical community
- ❑ **Easy to use (GUI-design driven by users)**

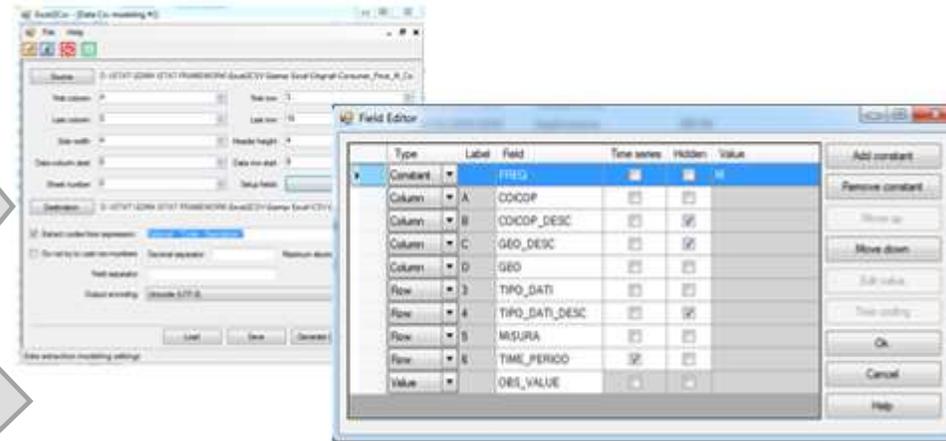Istat

# The Toolkit and the statistical processe

## Processing Transformation
- CUBUS
- Excel2CSV

## Collection
- Excel
- CSV (Custom and SDMX)
- SDMX-ML
- M2M (API)

## Dissemination
- GUI (Data Browser)
- M2M (APIs)
  - SDMX
  - CKAN (DCAT)

## Reporting
- Ad-hoc utilities

## Data Manager
- Performing star-scheme
- Data cubes management
- Data flows management

## Metadata Manager
- Structural metadata management
- Reference metadata management
- SDMX Registry
- Nomenclature server
- DCAT data catalogs

Istat

# Processing – Cubes Tool

# Transformation – Excel2Csv Tool



SDMX
Data
Structure Definition

### Code List (CSV)

```
CODE;NAME
IT;Italy
ITC1;Piemonte
ITC11;Torino
ITC12;Vercelli
ITC13;Biella
ITC15;Novara
ITC16;Cuneo
ITC18;Alessandria
```

### Data File (CSV)

```
1  FREQ;COICOP;GEO;TIPO_DATI;MISURA
2  M;01;ITC;9;4;Ott-2011;103,6
3  M;01;ITC;9;4;Nov-2011;104,3
4  M;01;ITC;9;4;Dic-2011;104,3
5  M;01;ITC;9;4;Gen-2012;104,6
6  M;01;ITC;9;4;Feb-2012;105,2
7  M;01;ITC;9;6;Ott-2011;0,6
8  M;01;ITC;9;6;Nov-2011;0,7
9  M;01;ITC;9;6;Dic-2011;0
10 M;01;ITC;9;6;Gen-2012;0,3
11 M;01;ITC;9;6;Feb-2012;0,6
12 M;01;ITC;9;7;Ott-2011;3,3
13 M;01;ITC;9;7;Nov-2011;3,7
14 M;01;ITC;9;7;Dic-2011;3,5
15 M;01;ITC;9;7;Gen-2012;3
16 M;01;ITC;9;7;Feb-2012;3,1
```

### SDMX Data File

Istat

# Excel2CSV tool – example of Excel multidimensional table

| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | | | | | | | | | | | |
| 3 | | | | | | | | **9** | | | |
| 4 | | **Index Type** | | | Consumer price index for the whole nation (base 2015=100) - monthly data | | | | | | |
| 5 | | **Measure** | | | [4]Index Number | | | | | [6] Percentage c | |
| 6 | | **Time period** | | | Oct-2011 | Nov-2011 | Dec-2011 | Jan-2012 | Feb-2012 | Oct-2011 | Nov-2011 |
| 7 | | **Istat** | **Territory** | | | | | | | | |
| 8 | | Food and non-alcoholic beverages | Nord-ovest | ITC | 103,6 | 104,3 | 104,3 | 104,6 | 105,2 | 0,6 | 0,7 |
| 9 | 01 | | Nord-est | ITD | 103,1 | 103,9 | 104 | 104,3 | 105,3 | 0,4 | 0,8 |
| 10 | | | Centro | ITE | 102,5 | 103,4 | 103,3 | 103,5 | 104,6 | 0,3 | 0,9 |
| 11 | | Alcoholic beverages and tobacco | Nord-ovest | ITC | 107,2 | 107,3 | 107,4 | 107,5 | 107,7 | 3,6 | 0,1 |
| 12 | 02 | | Nord-est | ITD | 107 | 107,1 | 107,1 | 107,2 | 107,3 | 3,5 | 0,1 |
| 13 | | | Centro | ITE | 107,4 | 107,4 | 107,4 | 107,6 | 107,6 | 3,7 | 0 |
| 14 | | Clothing and footwear | Nord-ovest | ITC | 102,6 | 102,7 | 102,8 | 102,9 | 102,7 | 0,6 | 0,1 |
| 15 | 03 | | Nord-est | ITD | 102,4 | 102,5 | 102,6 | 102,8 | 102,7 | 0,8 | 0,1 |
| 16 | | | Centro | ITE | 103,1 | 103,2 | 103,2 | 103,3 | 103,4 | 1,4 | 0,1 |

# Excel2CSV tool – example of Excel multidimensional table

| | Gross domestic product | | | Gross value added Total A10 | | | Agriculture, forestry and fishing | | | Mining and quarrying; manufacturing; electricity, gas, steam and air conditioning supply; water supply; sewerage, waste management and remediation activities | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | Total | | | of which: Manufacturing | | |
| STO ▶ | B1GQ | S T A T U S | S T A T U S | B1G | S T A T U S | S T A T U S | B1G | S T A T U S | S T A T U S | B1G | S T A T U S | S T A T U S | B1G | S T A T U S | S T A T U S |
| ACTIVITY ▶ | _Z | | | _T | | | A | | | BTE | | | C | | |
| ACCOUNTING ENTRY ▶ | B | | | B | | | B | | | B | | | B | | |
| TIME ▼ | 1=2+16+17 | | | 2=3+4+6+..+13 | | | 3 | | | 4 | | | 5 | | |
| 1995 | 1244538 | A | F | 1122615 | A | F | 26413 | A | F | 252752 | A | F | 223469 | A | F |
| 1996 | 1258660 | A | F | 1135998 | A | F | 26879 | A | F | 251989 | A | F | 221712 | A | F |
| 1997 | 1282146 | A | F | 1154423 | A | F | 27604 | A | F | 254242 | A | F | 224090 | A | F |
| 1998 | 1300714 | A | F | 1168621 | A | F | 28311 | A | F | 256061 | A | F | 226052 | A | F |
| 1999 | 1319588 | A | F | 1182242 | A | F | 30074 | A | F | 256249 | A | F | 225407 | A | F |
| 2000 | 1367801 | A | F | 1229008 | A | F | 29368 | A | F | 264569 | A | F | 233876 | A | F |
| 2001 | 1393278 | A | F | 1252220 | A | F | 28607 | A | F | 262305 | A | F | 232056 | A | F |
| 2002 | 1399568 | A | F | 1257988 | A | F | 27786 | A | F | 261154 | A | F | 230313 | A | F |
| 2003 | 1398916 | A | F | 1255411 | A | F | 26493 | A | F | 255273 | A | F | 224565 | A | F |
| 2004 | 1423126 | A | F | 1278452 | A | F | 29908 | A | F | 259508 | A | F | 227912 | A | F |
| 2005 | 1436379 | A | F | 1291692 | A | F | 28600 | A | F | 261909 | A | F | 229848 | A | F |
| 2006 | 1467964 | A | F | 1320418 | A | F | 28276 | A | F | 272010 | A | F | 239639 | A | F |
| 2007 | 1492671 | A | F | 1344313 | A | F | 28332 | A | F | 279679 | A | F | 247336 | A | F |
| 2008 | 1475412 | A | F | 1329002 | A | F | 28729 | A | F | 271375 | A | F | 238470 | A | F |
| 2009 | 1394347 | A | F | 1254718 | A | F | 28007 | A | F | 230422 | A | F | 198986 | A | F |
| 2010 | 1418376 | A | F | 1276477 | A | F | 27952 | A | F | 244266 | A | F | 214249 | A | F |
| 2011 | 1424752 | A | F | 1284355 | A | F | 28105 | A | F | 247946 | A | F | 217861 | A | F |
| 2012 | 1391018 | A | F | 1256553 | A | F | 26908 | A | F | 240500 | A | F | 210230 | A | F |
| 2013 | 1365227 | A | F | 1236836 | A | F | 26980 | A | F | 232792 | A | F | 203609 | A | F |

Istat

# Meta Manager – main features

Data modelling

Structural metadata utilities

Reference Metadata Editor

Data catalog (DCAT) Editor

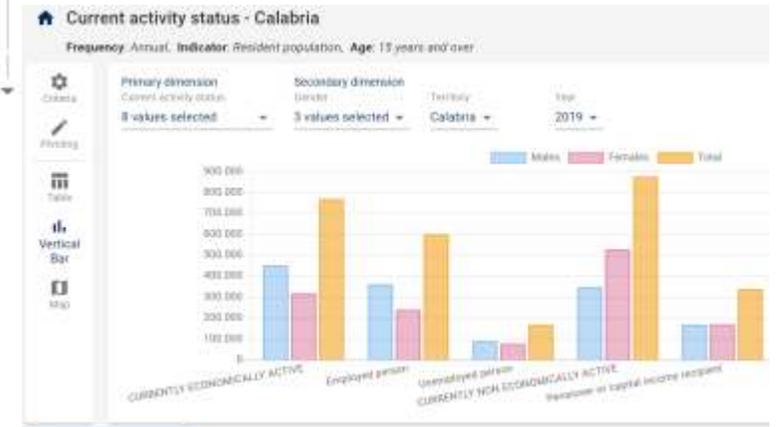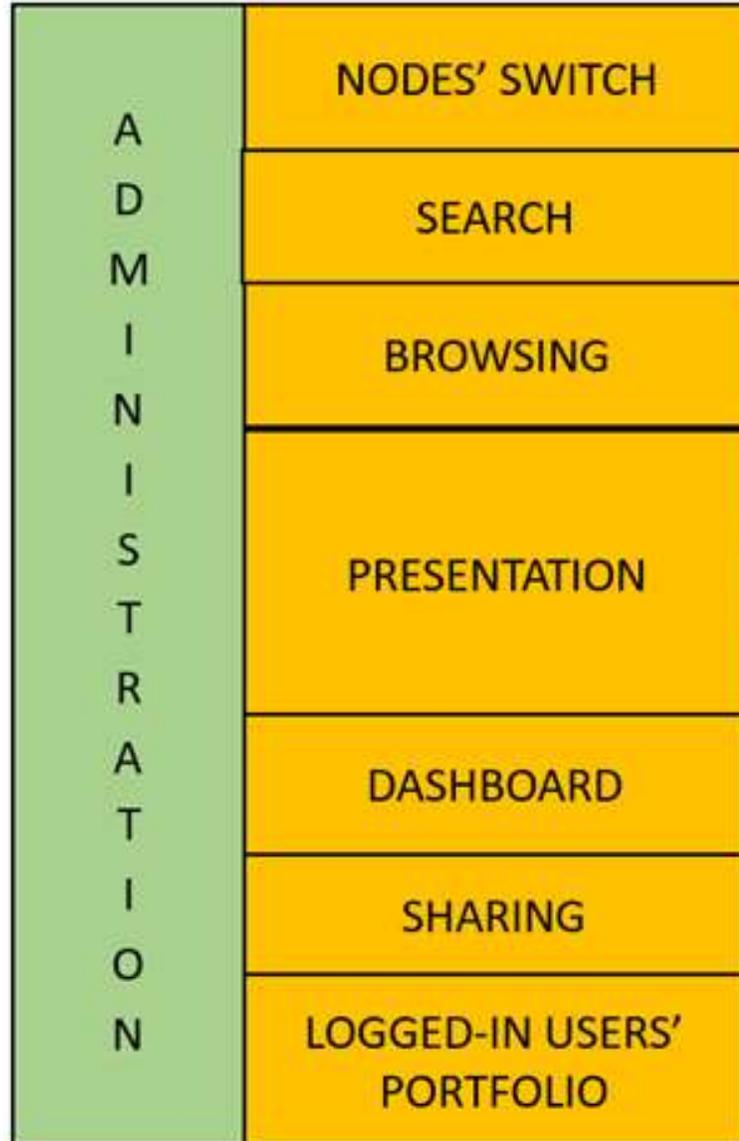Nomenclature server

"layout" Annotations editor

- ❑ **SDMX artefacts CRUD**: Concept Schemes, Code Lists, Concept Schemes, Category Schemes, Agency Schemes, Data Structure Definitions, Data Flows, Categorisations
  - ▪ Hierarchical CL, Metadata Structure Definitions, Metadata Flows can be used but not created
- ❑ **Structural metadata utilities**:
  - ▪ Derive item schemes (Code Lists, Concept Schemes, Category Schemes)
    - ✓ Add items preserving hierarchy
    - ✓ Include "Parents", "Children", "Descendants"
  - ▪ Merge item schemes (Code Lists, Concept Schemes, Category Schemes)
  - ▪ Compare item schemes (Code Lists, Concept Schemes, Category Schemes)
  - ▪ Compare DSDs
  - ▪ Upgrade DSD
- ❑ **Overcome SDMX 2.1 deficiencies**: add items to final item schemes, modify/add names, descriptions and annotations to final artefacts
- ❑ **Manage annotations** for presentation purpose (e.g. Order, Layout, etc.)
- ❑ **Reference metadata** editor
- ❑ **DCAT Data catalog** editor

Istat

# Data Manager – main features



❑ **Building:** Multidimensional data cubes based on SDMX DSDs (advance database star-schema synchronized with the structural metadata repository)

❑ **Mapping:**

- "Custom" CSV files
- Excel files (multidimensional statistical tables)

❑ **Loading:** Excel, CSV, SDMX-ML data files

❑ **Publishing**: Data flows

- As sub cube of a data cube
- As sub cube with less dimensions (a new DSD is generated automatically)

❑ **Download:** SDMX-ML (2.0 & 2.1), SDMX-CSV, SDMX-JSON, Custom CSV, RDF/XML, DCAT-JSON

# Data Browser – main features

# Lesson learnt

- ❏ Top management has to foster the initiative and must be involved regularly
  - ❏ The state of progress and bottlenecks must be reported in order to smooth the implementation (monitoring can be achieved by a cross-cutting working group)
- ❏ A suitable time must be devoted in analysing the experiences performed by other organisations
- ❏ Statistical standards must be the main ingredients of the innovation
  - ❏ No standards no industrialization
- ❏ "Wheel must not reinvented"
  - ❏ Reusing software available in the statistical community means save money and time
- ❏ A winning design is a good balance between innovation and integration
  - ❏ Step-by-step approach
- ❏ Capacity building actions facilitate the knowledge of the new methodologies and technologies

# Contacts and links

❑ Meta and Data Manager Tool
❑ Data Browser Tool
❑ Excel2Csv Tool

➡ Alessio Cardacino – alcardac@istat.it
Francesco Rizzo – rizzo@istat.it

❑ Cubus Tool

➡ Luca Ramadori - luca.ramadori@istat.it
Guido Drovandi – drovandi@istat.it

Useful links:

SDMX Istat toolkit download: https://sdmxistattoolkit.github.io/index.html

"Permanent census of population and housing": https://esploradati.censimentopopolazione.istat.it/databrowser/#/en

"Hub of the Public Statistics": https://sistanhub.istat.it/databrowser/#/en

Istat

# Thanks

Massimo FEDELI | fedeli@istat.it

Francesco RIZZO | rizzo@istat.it