

Assessing Fitness for Integration - A Metadata- driven Approach

Thomas Gottron, Andrea Novello, Ilias Aarab, Bernadette Lauro
European Central Bank



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



EUROPEAN CENTRAL BANK

EUROSYSTEM

Assessing Fitness for Integration

A Metadata-driven
Approach



06/06/2024

Thomas Gottron, Andrea Novello, Ilias Aarab, Bernadette Lauro
European Central Bank

The challenge

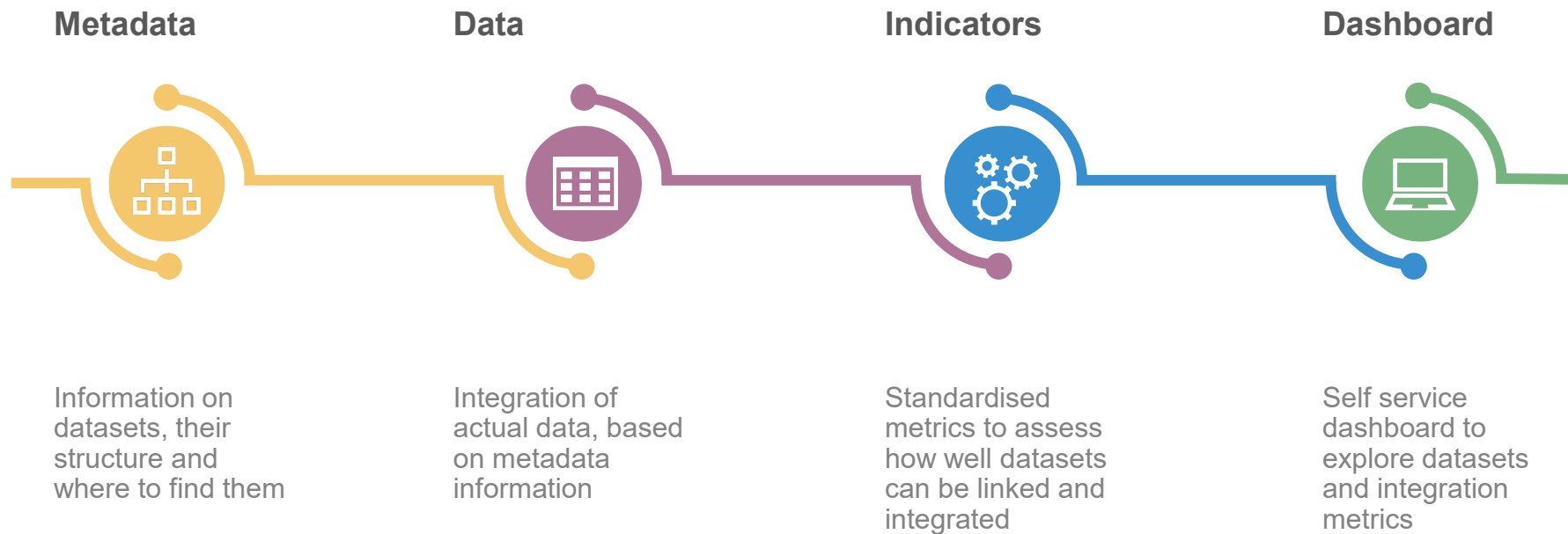


Data analytics is often based on **combining granular datasets**.

Users start by addressing **three questions**:

- (1) Which data sets can be integrated?
- (2) How to perform integration in a semantically meaningful way?
- (3) How good is the linkage rate?

Idea: Provide a dashboard to answer the questions



Fully automated process

How does it work – Step 1: Metadata

Metadata



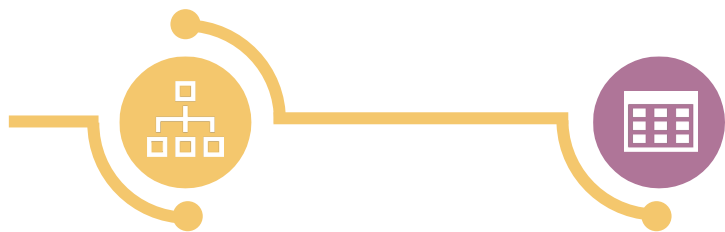
Identifier
Time
Unit
Stratification
...

- **Integrability** $I(A, B)$: indicates if two datasets can be integrated at all, based on common **identifiers**.
- **Linkage rate** $L_i(A, B)$: ratio of **identifier** values present in A that can also be found in B , considering **time** information, e.g. reference period, validity ranges.
- **Weighted linkage rate** $L_{u,f}(A, B)$: based on a measurable **unit**, e.g. reflecting linkage weighted by market values, outstanding nominal amount etc.
- **Breakdowns: stratification** attributes are used to analyse integration on subsets of data.

How does it work – Step 2: Determine integrability

Metadata

Data



Identifier
Time
Unit
Stratification
...

Dataset A

ID	U	S	T _A

Dataset B

ID	T _B



ID	U	S	T _A	ID	T _B

How does it work – Step 3: Compute metrics

Metadata

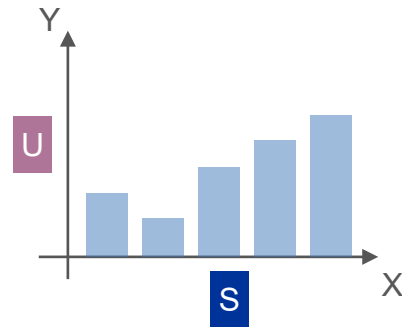
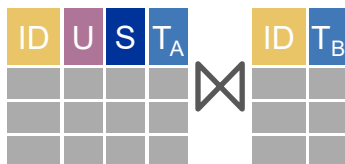
Data

Indicators



Source	Target	Metric	U	S	T
Table A	Table B	Linkage rate	Y ₁	X ₁	01/23
Table A	Table B	Linkage rate	Y ₂	X ₂	01/23
Table A	Table B	Linkage rate	Y ₃	X ₃	01/23

Identifier
Time
Unit
Stratification
...



How does it work – Step 4: Visualise

Metadata

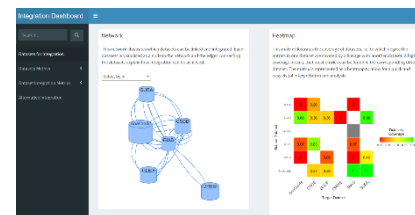
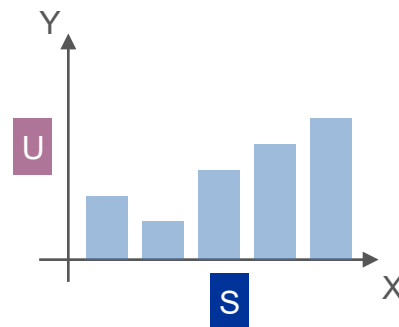
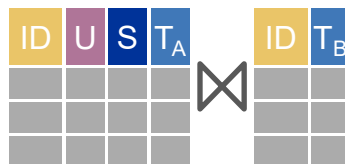
Data

Indicators

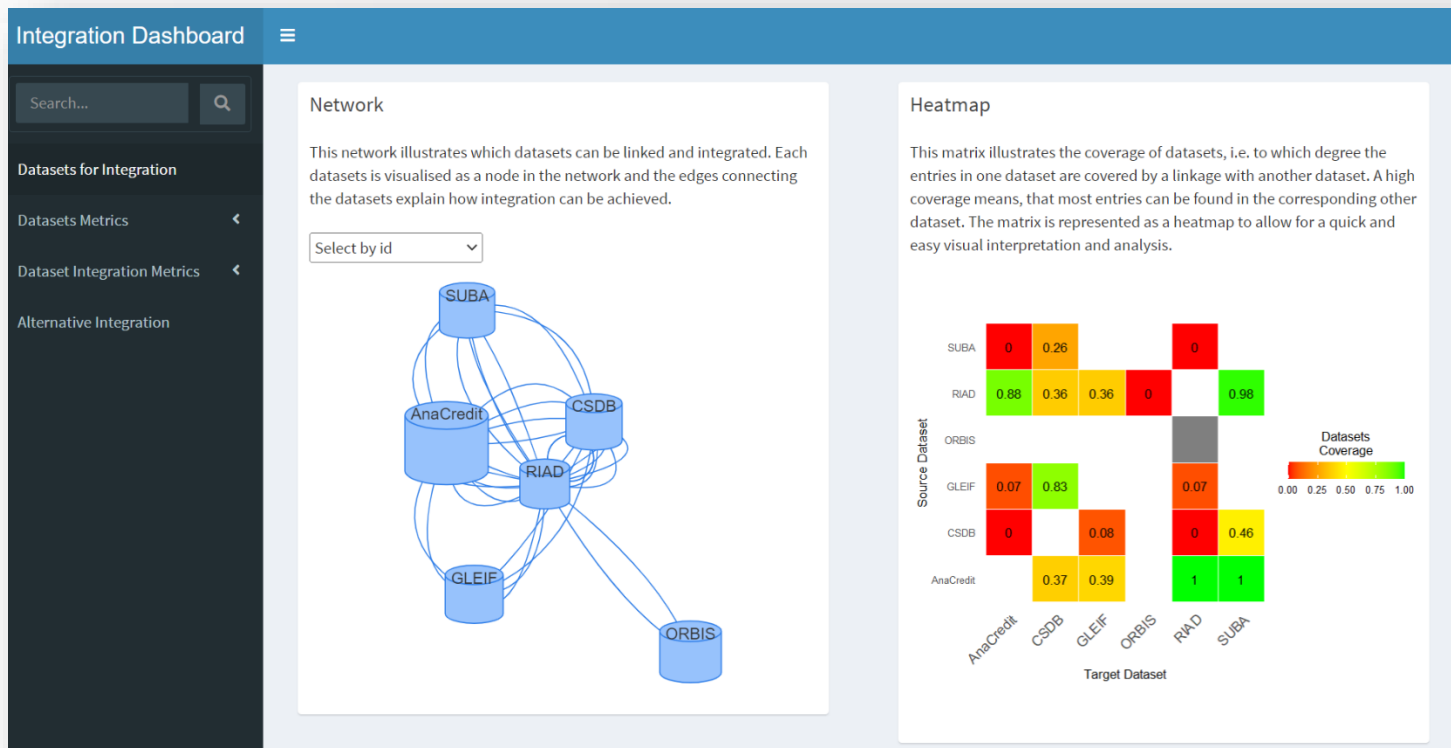
Dashboard



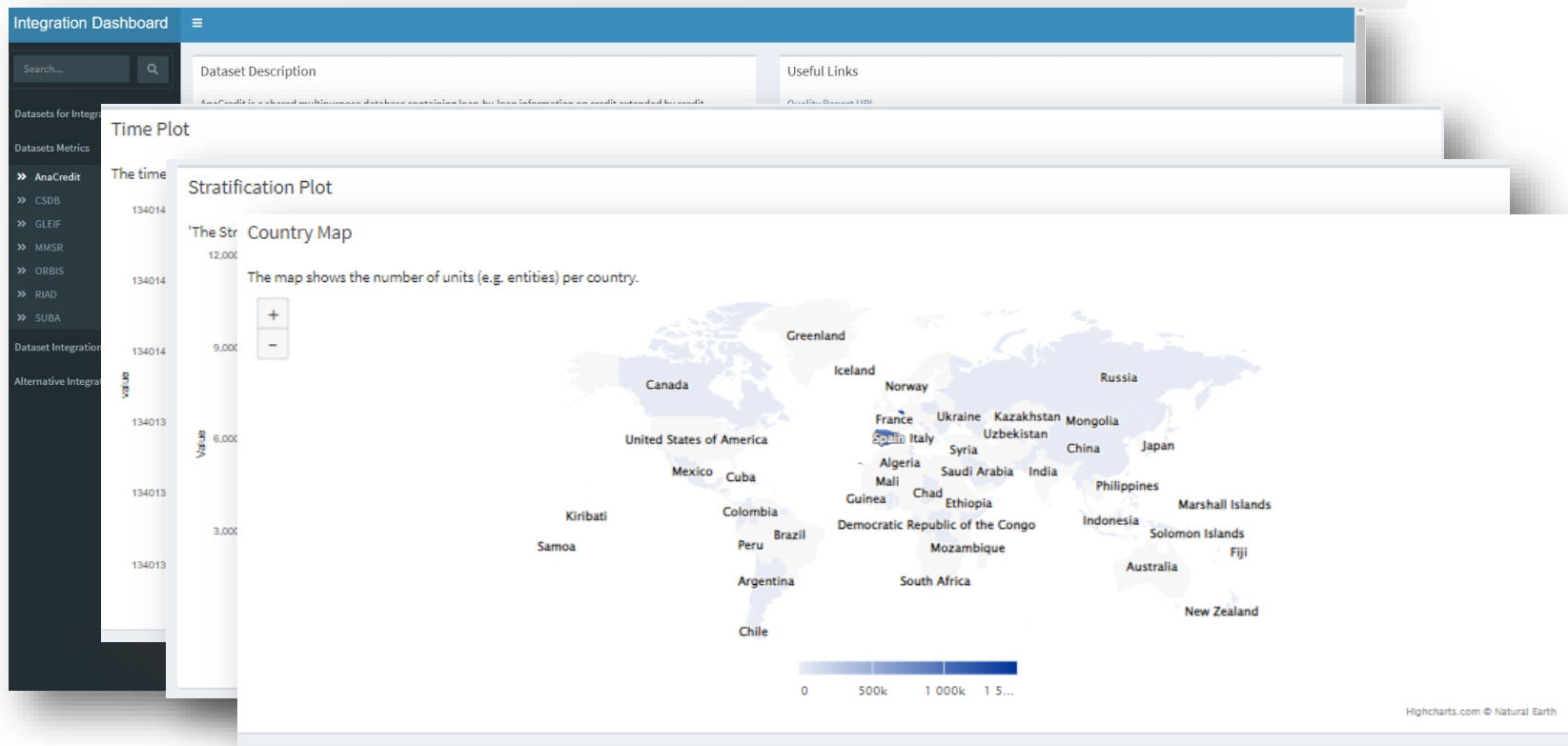
Identifier
Time
Unit
Stratification
...



Prototype dashboard – Overview of datasets

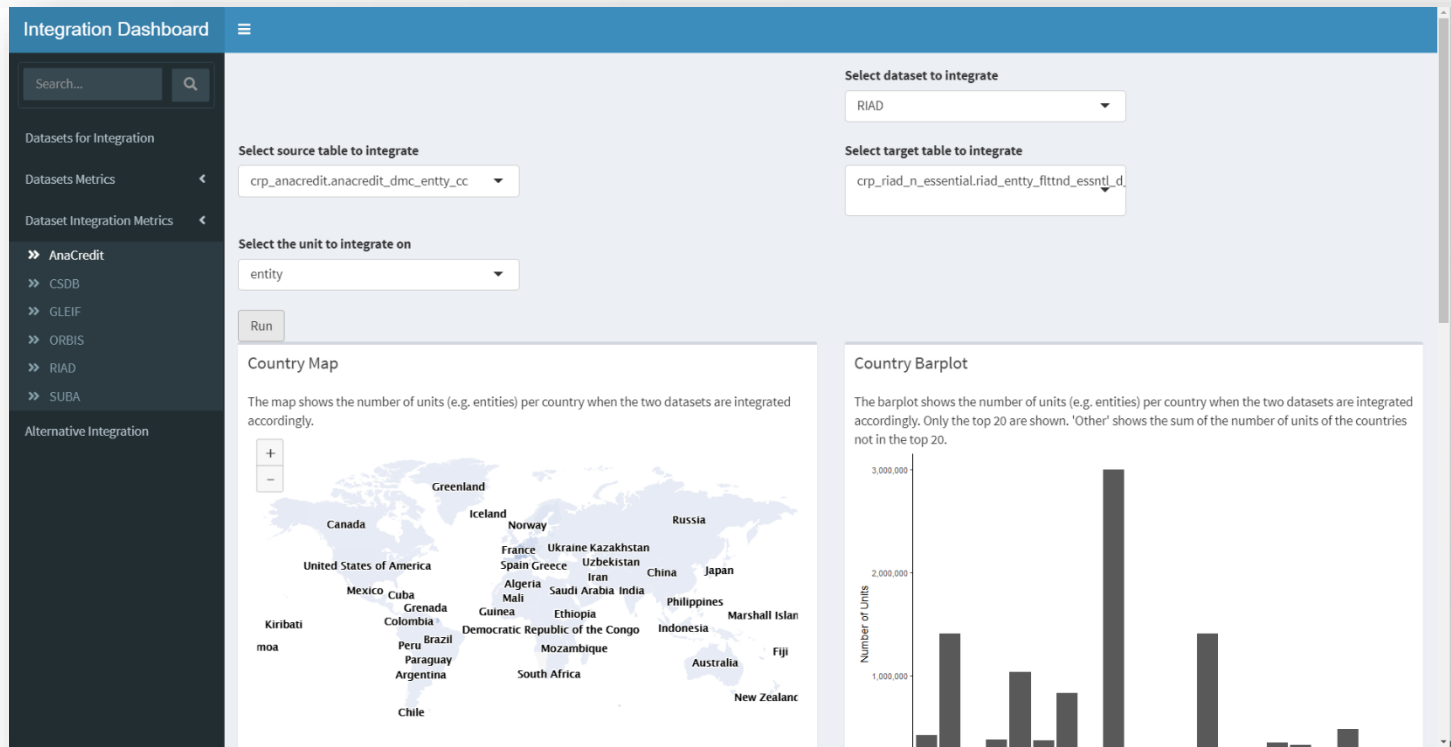


Prototype dashboard – Details for a single dataset



Note: due to confidentiality of some datasets, the visualisations are based on artificial data

Prototype dashboard – Integration of two datasets



Fitness for Integration



- The Fitness for Integration dashboard **supports users** in common data integration and exploration tasks
- The report is **metadata driven and automated** to reduce the need for manual investigations
- A **prototype has been implemented** and tested with users

Thank you and happy to take questions 😊



EUROPEAN CONFERENCE ON QUALITY IN OFFICIAL STATISTICS 2024 ESTORIL - PORTUGAL