



EUROPEAN CONFERENCE ON QUALITY IN OFFICIAL STATISTICS 2024 ESTORIL - PORTUGAL



Enhancing Data Quality:

A DataOps approach with R and GitLab



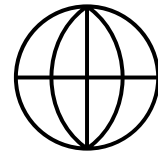
EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

Summary

1. Motivation
2. Data Pipelines in International Trade Data Collection
3. GitLab as a DataOps board
4. Improving data quality with R packages
5. Conclusions

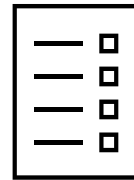


Motivation



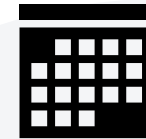
International Trade Data Collection Unit (ITDCU)

- INTRA: imports/exports with EU members
- EXTRA: imports/exports with non-EU members



Score Model

- Implementation of the Swedish Foreign Trade Statistics (SFTS) (*Norberg & Jader, 2005*)
- Considers not only the suspicious error but also the potential impact of the record



Daily Routine

- Data is updated twice a day
- May contain records from up to a year ago

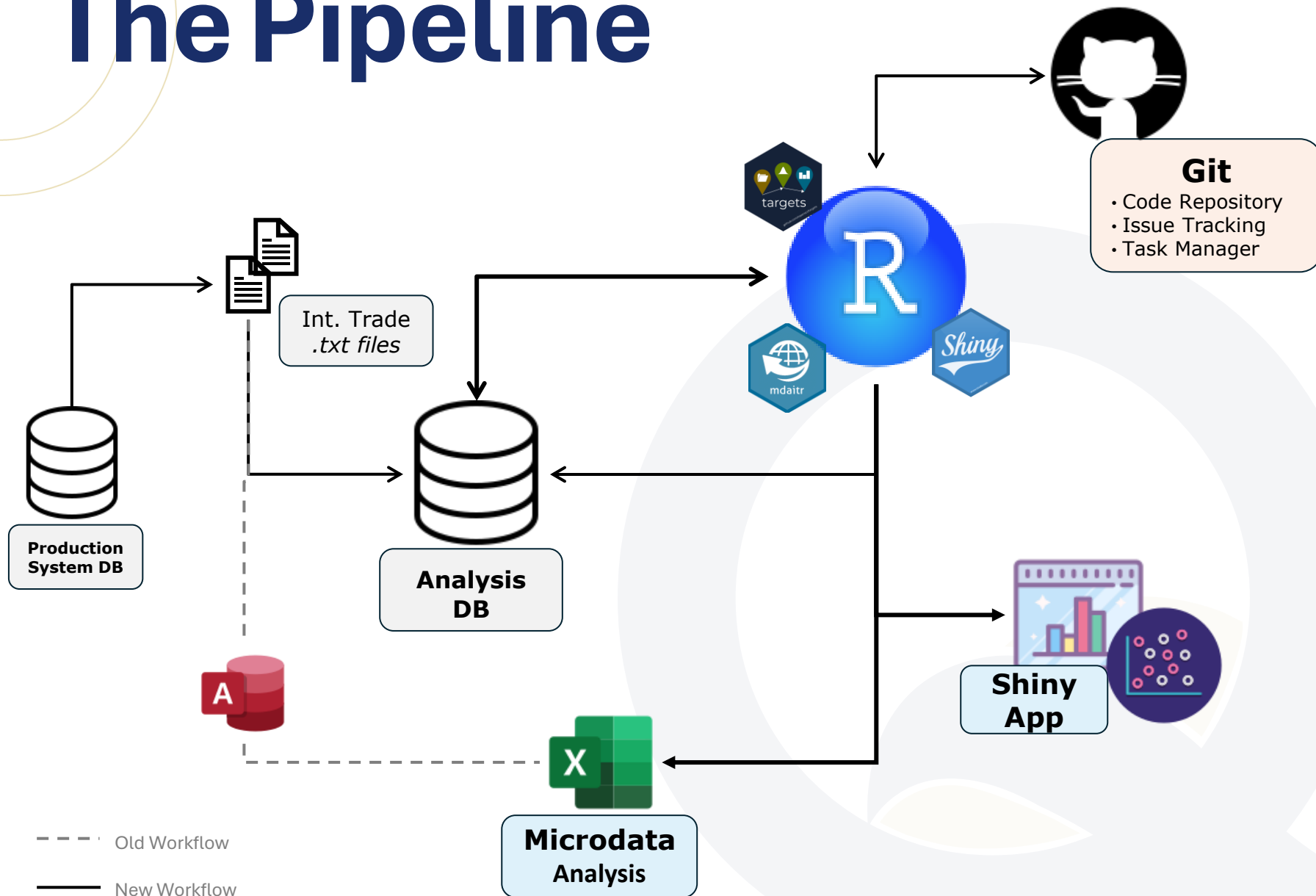


Microdata Analysis

- ITDCU perform daily analysis
- Suspicious get ranked and distributed to colleagues



The Pipeline





DataOps

CI/CD approach

- High quality, automated, reproducible processes
- Culture of continuous improvement in data processing
- Speed and collaboration



GitLab

DataOps Board

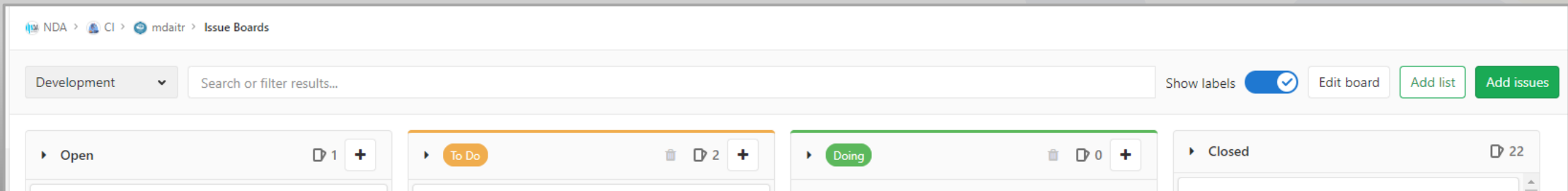
- Define, assign and track tasks or objectives
- At-a-glance view of project progress
- Issue Tracking and notifications
- Version Control



GitLab

DataOps Board

- Define, assign, and track tasks or objectives
- At-a-glance view of project progress
- Issue Tracking and notifications
- Version Control





R

ETL Workflow

Load data from RDBMS
databases

Calculate the *score value*

Export to .xlsx

Update Shiny App



R

ETL Workflow

Load data from RDBMS
databases

Calculate the *score value*

Export to .xlsx

Update Shiny App

Calculate quantiles for each
Combined Nomenclature

Calculate Suspicious Error

Calculate Potential Impact

Calculate `score()`



- Adapted from an implementation of the Swedish Foreign Trade Statistics (SFTS) (*Norberg & Jader, 2005*)

R

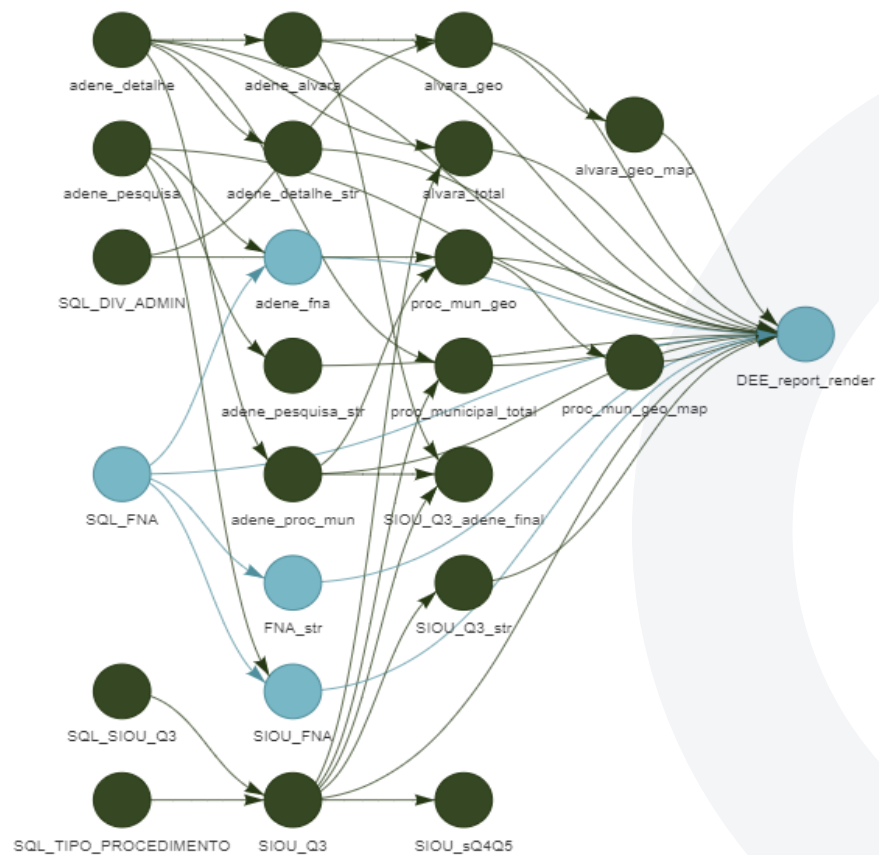
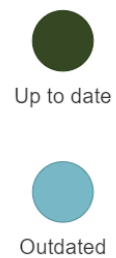
`{mdaitR}`

- Developed under the Functional Programming paradigm
- Creates tested and documented functions
- Can be applied to other data domains



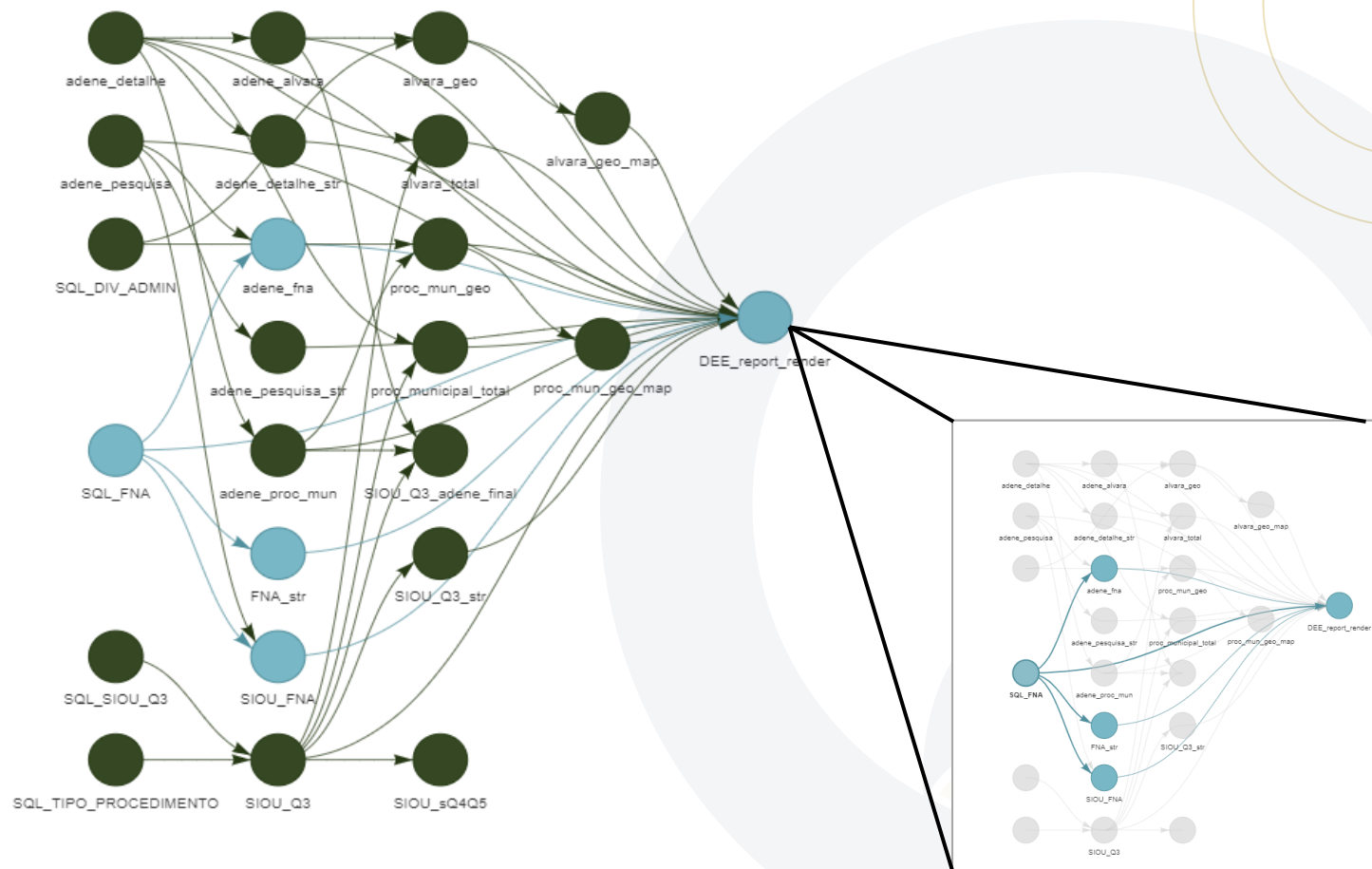
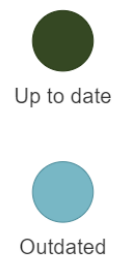
R

{targets}





R {targets}





Conclusions

- RAPs are based on best practices with the aim of ensuring data pipelines that are reproducible, auditable, efficient and of high quality
- Anything that can be automated should be automated!



EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

Thank You!



INSTITUTO NACIONAL DE ESTATÍSTICA
STATISTICS PORTUGAL

eurostat 

The conference is partly
financed by the European Union



EUROPEAN CONFERENCE ON QUALITY IN OFFICIAL STATISTICS 2024 ESTORIL - PORTUGAL