



EUROPEAN CONFERENCE ON QUALITY IN OFFICIAL STATISTICS 2024 ESTORIL - PORTUGAL



EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL



A step-by-step process to deal with the protection of a set of tabular data

Julien Jamme, Insee
Clara Baudry, Insee





EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL



Context

- Very large tabular data dissemination
- SDC with suppressive methods
- High sensitivity of the protection step to the actual dissemination content
- Protection methodology and tools require high level of expertise



Issues

- Problems not generalized
 - Processes not automated
 - Complexity of the task
- => Protection process supported by a team in the department of statistical methods



Objective

- The long-term goal : Let the producers protect data on their own
- To reach it :
 - Develop methodology and tools to reduce the expertise level required to handle protection of large tabular data dissemination



Methodological approach

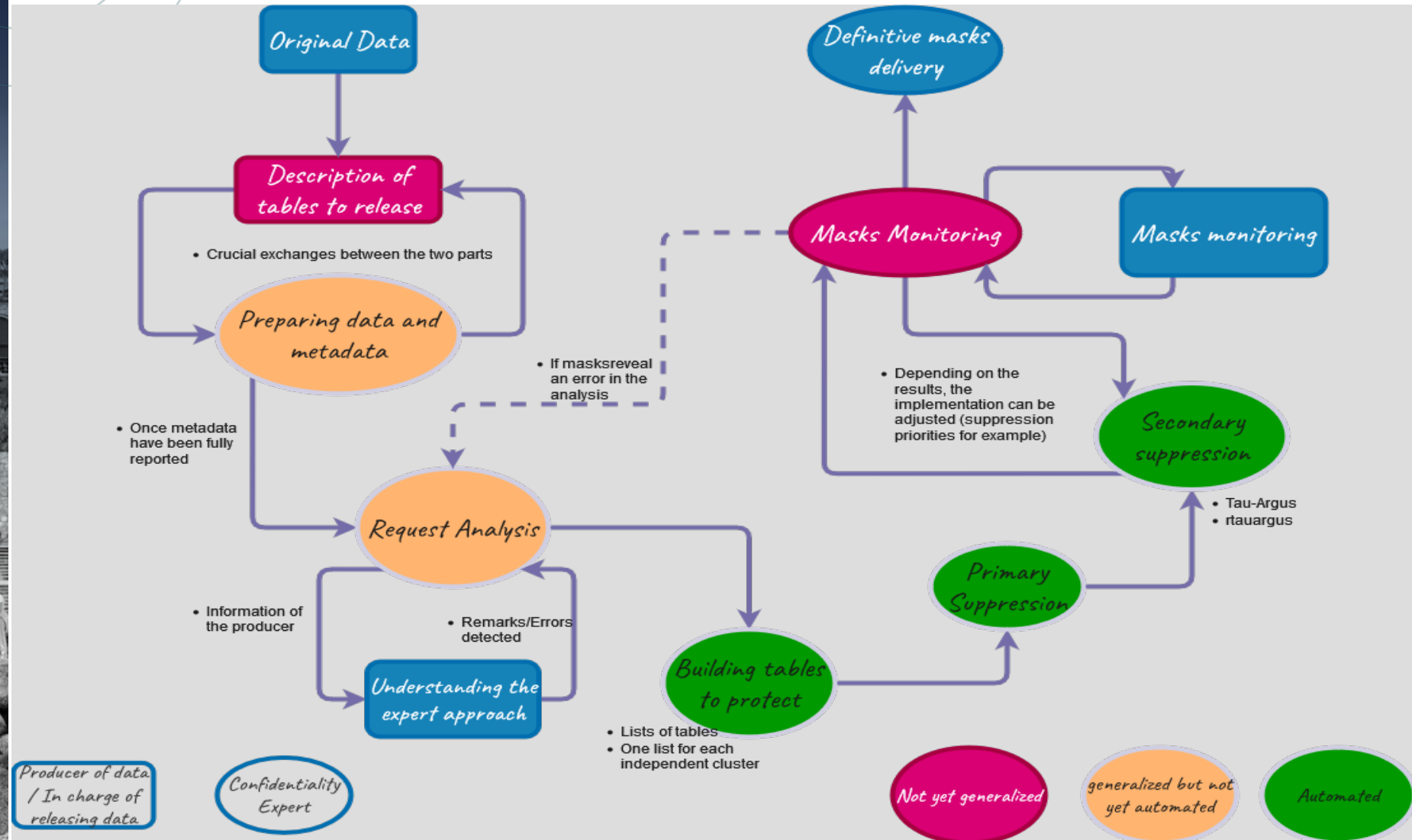
- **Describe** the process step by step to identify issues and sticking points
- **Generalize** each step to handle every case of dissemination
- **Automate** where possible to reduce the implementation burden and make the code reproducible as much as possible
- **Disseminate** methods, tools and practices by training people

Methodology improvements while producing confidentiality masks

Ongoing work



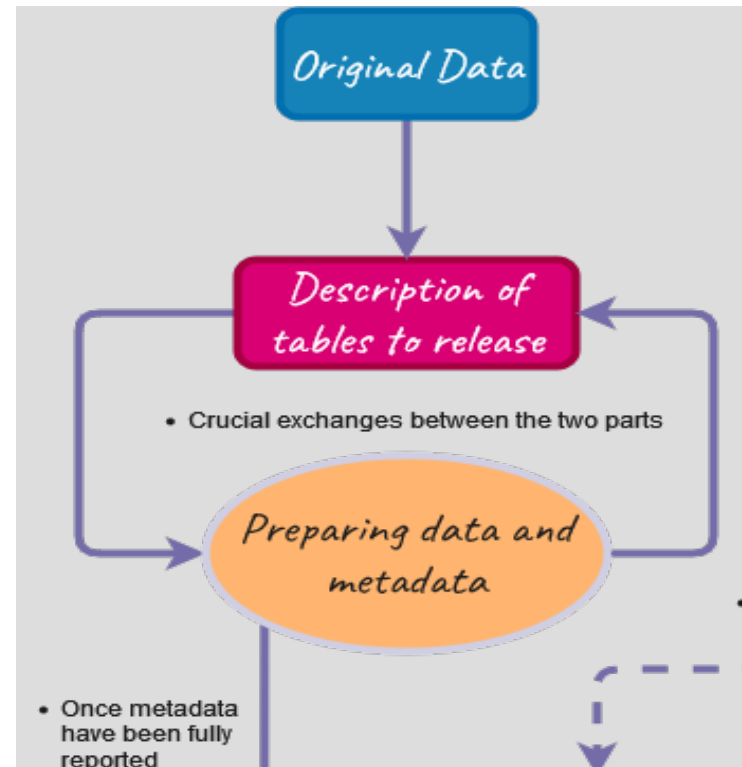
Describe





Issues and Sticking points

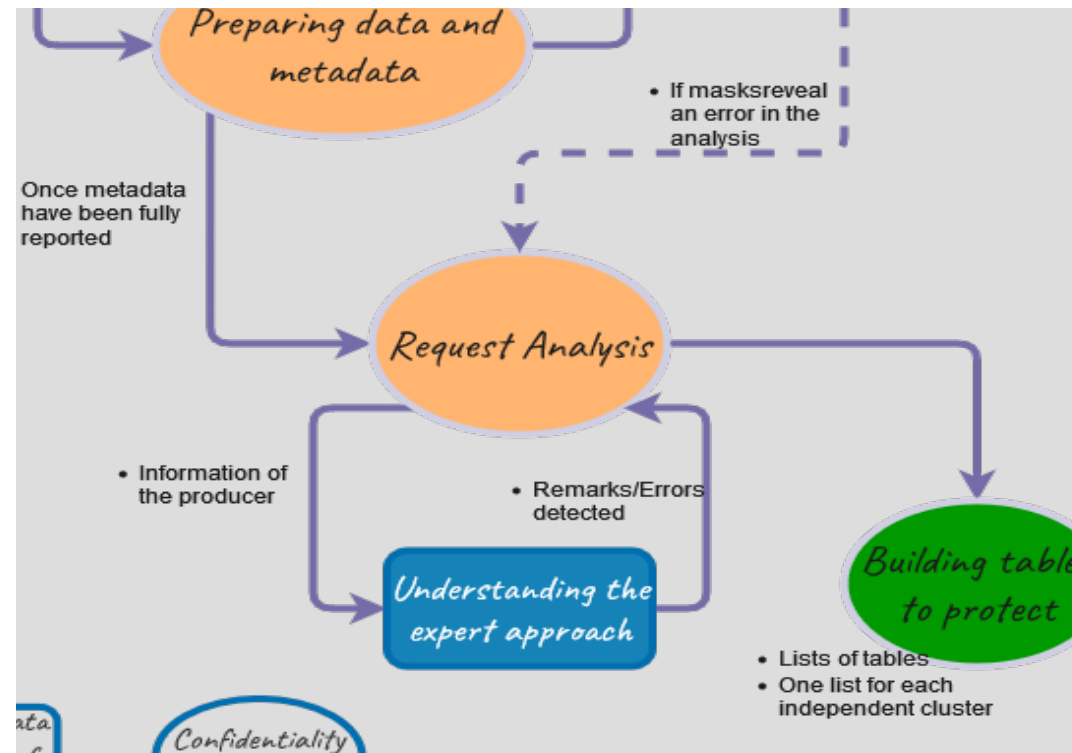
1) The initial contact between producer and confidentiality expert





Issues and Sticking points

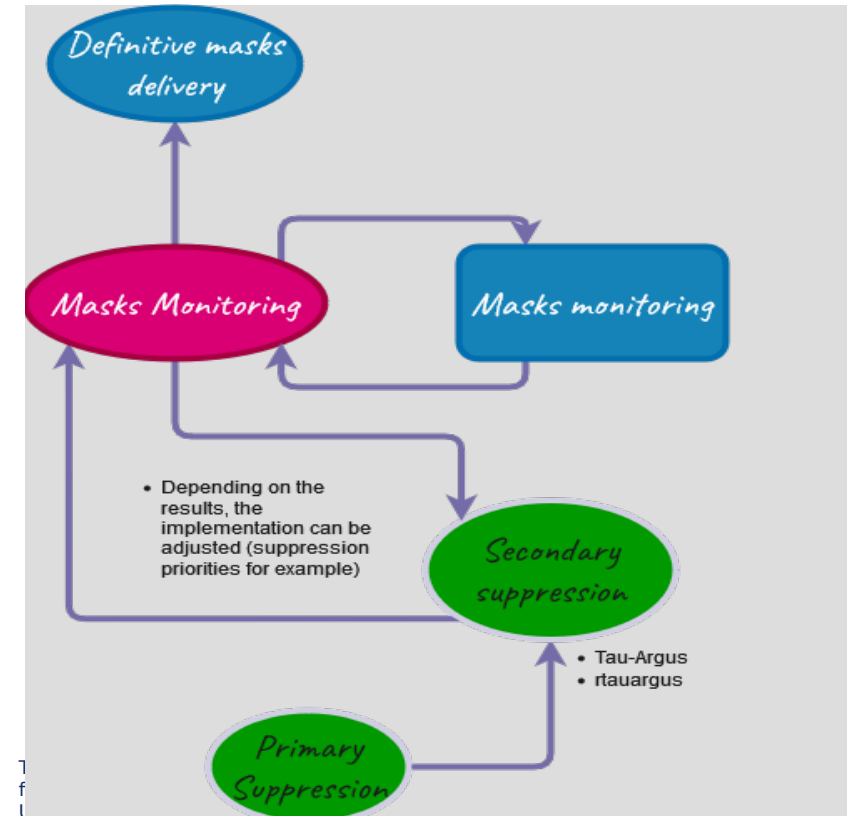
2) The analysis of the dissemination





Issues and Sticking points

3) The implementation of the suppression process on very large dissemination





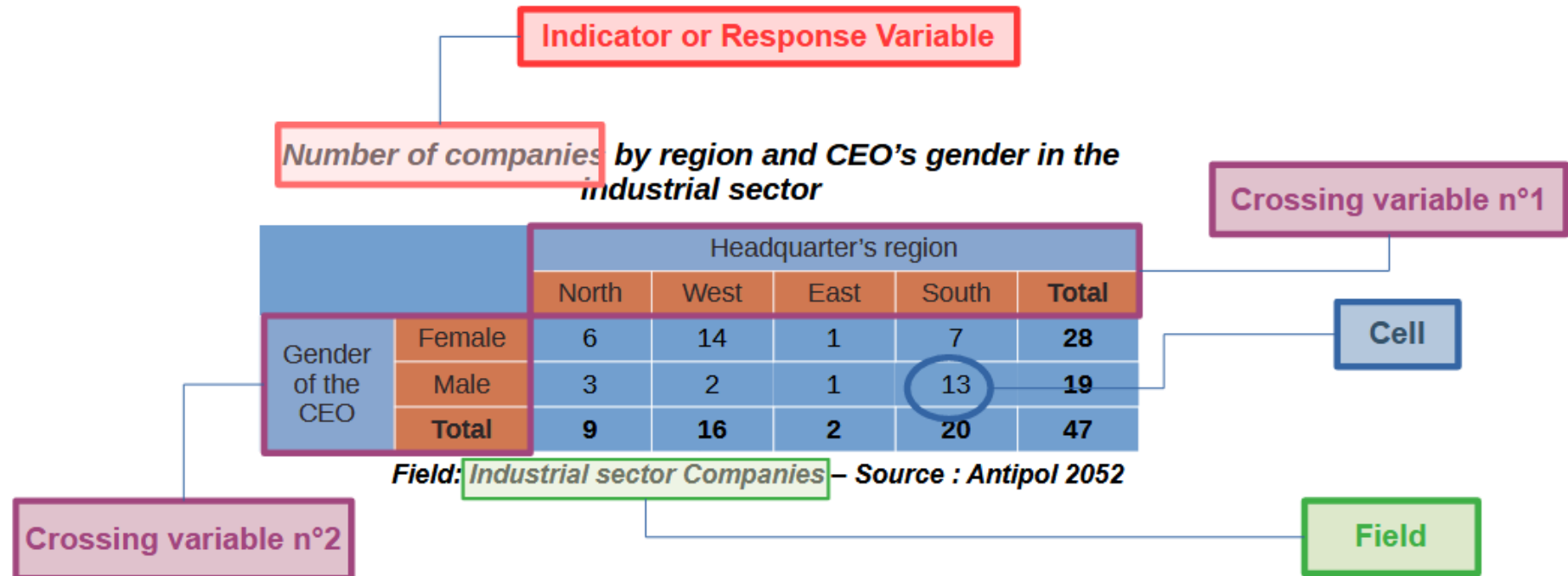
Step 1 : Building relevant metadata

- First issue : Make the initial contact productive
- Objective :
 - Producer : to give a complete description of the dissemination
 - Expert : to get all the information necessary for the protection process (and no more)
- data expert vs confidentiality expert



Step 1 : Building relevant metadata

- Find a way to universally describe a table (from a protection perspective)
- What is required ?





Step 1 : Building relevant metadata

- Find a way to universally describe a table (from a protection perspective)
- Formalize the description

$$RV^{hrcRV} \otimes_{Field^{hrcF}} \left\{ CV1_{tot1}^{hrc1} \times CV2_{tot2}^{hrc2} \right\}$$



Step 1 : Building relevant metadata

- Find a way to universally describe a table (for the protection perspective)
- Formalize the description

$$\text{Freq}_{Ind.Sect} \otimes \{ \text{Region}_{total} \times \text{Gender}_{total} \}$$



Step 1 : Building relevant metadata

- Relevant metadata contain a complete description of all tables :
 - spanning and response variables
 - field definition
 - hierarchies (nested or non-nested),
 - etc.
- They must allow to analyze the links between the tables



Step 1 : Building relevant metadata

- Producer has to be able to fill in the file
=> Choice for a spreadsheet format

<i>Table</i>	<i>Field</i>	<i>Response Variable (RV)</i>	<i>Crossing Variable (CV)</i>	<i>CV Total Code</i>
T1	Companies of Industrial Sector	Frequencies	HQ's Region	Total
T1	Companies of Industrial Sector	Frequencies	CEO's Gender	Total



Step 2 : Analyzing the dissemination

- Objective of the analysis :
 - Determine the links between tables
 - Linked tables => simultaneous protection
 - Unlinked tables => independent protection
 - Ultimate goal : describe the tables on which suppression is really applied



EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

Step 2 : Analyzing the dissemination

- Issues

- High expertise and experience required
- Very sensitive to the actual content
- Lack of a methodology to handle this step
- Fit the method to every situation (non nested hierarchies, holding variable, links on response or spanning variable, complementary fields, etc.)



EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

Step 2 : Analyzing the dissemination

- Our proposition
Build a complete methodology to analyze dissemination :
 - By describing steps and objectives
 - By generalizing the approach to handle every situation as much as possible
- Automation of the analysis is under way



Step 2 : Analyzing the dissemination

$\left\{ \begin{array}{l} \text{to_margarita} \otimes \{\text{NUTS2} \times \text{SIZE}\} \\ \text{to_margarita} \otimes \{\text{NUTS3} \times \text{SIZE}\} \\ \text{to_calzone} \otimes \{\text{NUTS2} \times \text{SIZE}\} \\ \text{to_calzone} \otimes \{\text{NUTS3} \times \text{SIZE}\} \\ \text{to_pizzas} \otimes \{\text{NUTS2} \times \text{SIZE}\} \\ \text{to_pizzas} \otimes \{\text{NUTS3} \times \text{SIZE}\} \end{array} \right.$



Step 2 : Analyzing the dissemination

$$\left\{ \begin{array}{l} \text{to_margarita}^{pizzas} \otimes \{ \text{NUTS2}_{all}^{nuts} \times \text{SIZE}_{all} \} \\ \text{to_margarita}^{pizzas} \otimes \{ \text{NUTS3}_{all}^{nuts} \times \text{SIZE}_{all} \} \\ \text{to_calzone}^{pizzas} \otimes \{ \text{NUTS2}_{all}^{nuts} \times \text{SIZE}_{all} \} \\ \text{to_calzone}^{pizzas} \otimes \{ \text{NUTS3}_{all}^{nuts} \times \text{SIZE}_{all} \} \\ \text{to_pizzas}^{pizzas} \otimes \{ \text{NUTS2}_{all}^{nuts} \times \text{SIZE}_{all} \} \\ \text{to_pizzas}^{pizzas} \otimes \{ \text{NUTS3}_{all}^{nuts} \times \text{SIZE}_{all} \} \end{array} \right.$$



EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

Step 2 : Analyzing the dissemination

$$\left\{ \begin{array}{l} \text{to_margarita}^{pizzas} \otimes \{ \text{NUTS}_{all}^{nuts} \times \text{SIZE}_{all} \} \\ \text{to_calzone}^{pizzas} \otimes \{ \text{NUTS}_{all}^{nuts} \times \text{SIZE}_{all} \} \\ \text{to_pizzas}^{pizzas} \otimes \{ \text{NUTS}_{all}^{nuts} \times \text{SIZE}_{all} \} \end{array} \right.$$



EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

Step 2 : Analyzing the dissemination

$$\text{to } \otimes \left\{ \text{NUTS}_{all}^{nuts} \times \text{SIZE}_{all} \times \text{PIZZAS}_{pizzas}^{(h)} \right\}$$



Step 3 : Suppression

- Initially
 - the most intensive step in terms of code
 - One dissemination => One algorithm to implement to handle all the links between tables
- Context
 - Tau-Argus as reference
 - But some limitations with numerous linked tables dissemination

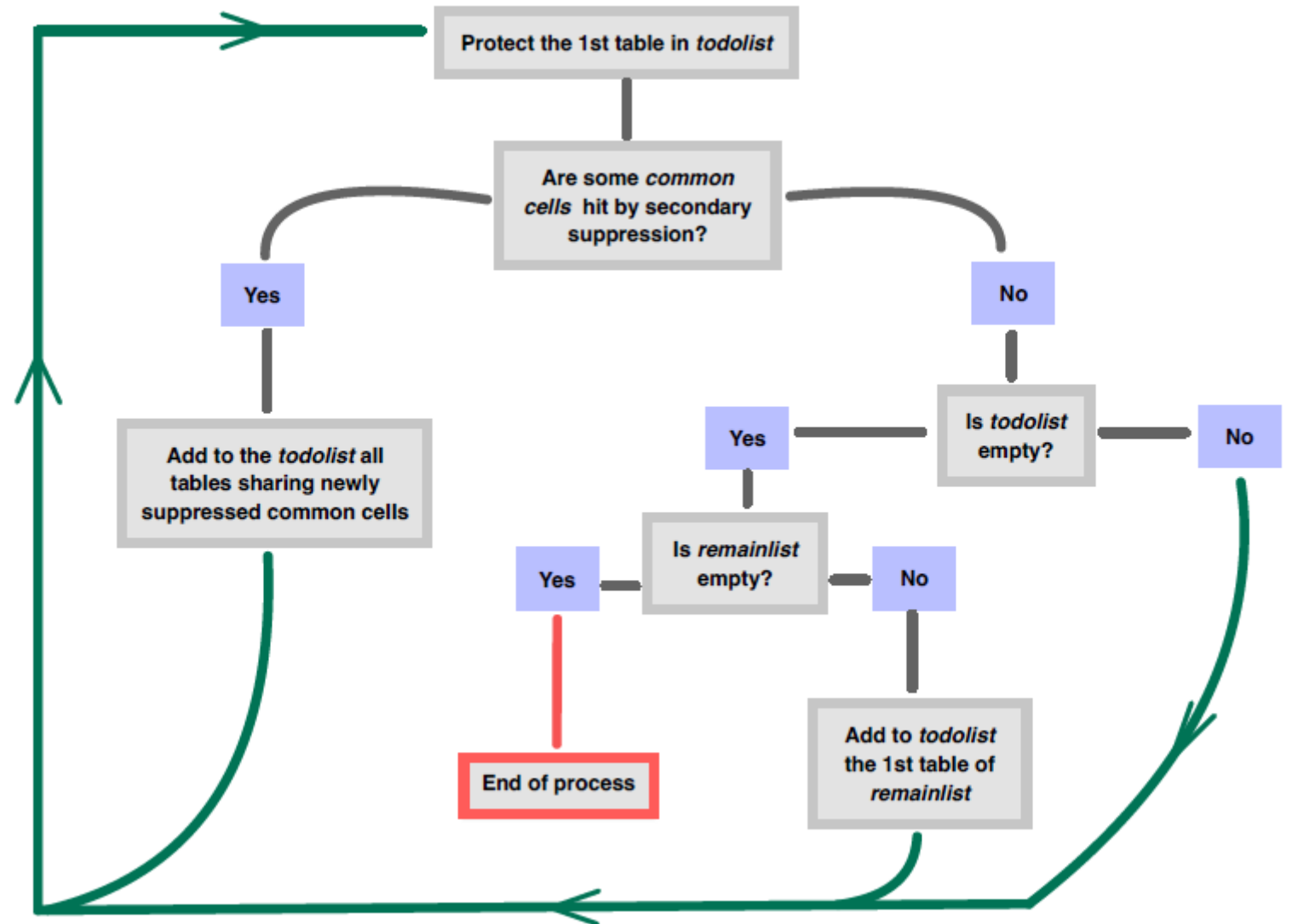


Step 3 : Suppression

- Solution
 - **rtauargus** : R package developed to interface R with Tau-Argus
 - + algorithm implemented to manage the linked tables
 - Versionning with **git**
- Productivity and efficiency gain
- Reduction of code length



Step 3 : Suppression





Review

- Description of the whole process => to identify issues to be addressed
- Generalization :
 - Formalization of table description
 - Building of metadata file
 - Describing steps to analyze a dissemination
- Automation :
 - Analysis step (under way)
 - Suppression step (completed)
- Dissemination of expertise :
 - Internal training
 - Guidelines



EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL



Outlook

How far are we from the ultimate goal of
our quest ?



EUROPEAN CONFERENCE ON QUALITY IN OFFICIAL STATISTICS 2024 ESTORIL - PORTUGAL