# Towards Better Identification of Fact Tables in Statistical Spreadsheets
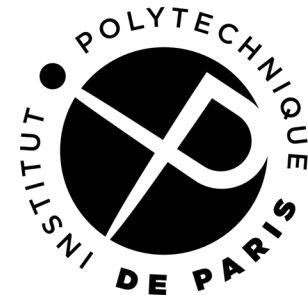
Paul Kronlund-Drouault

Supervisor: Ioana Manolescu
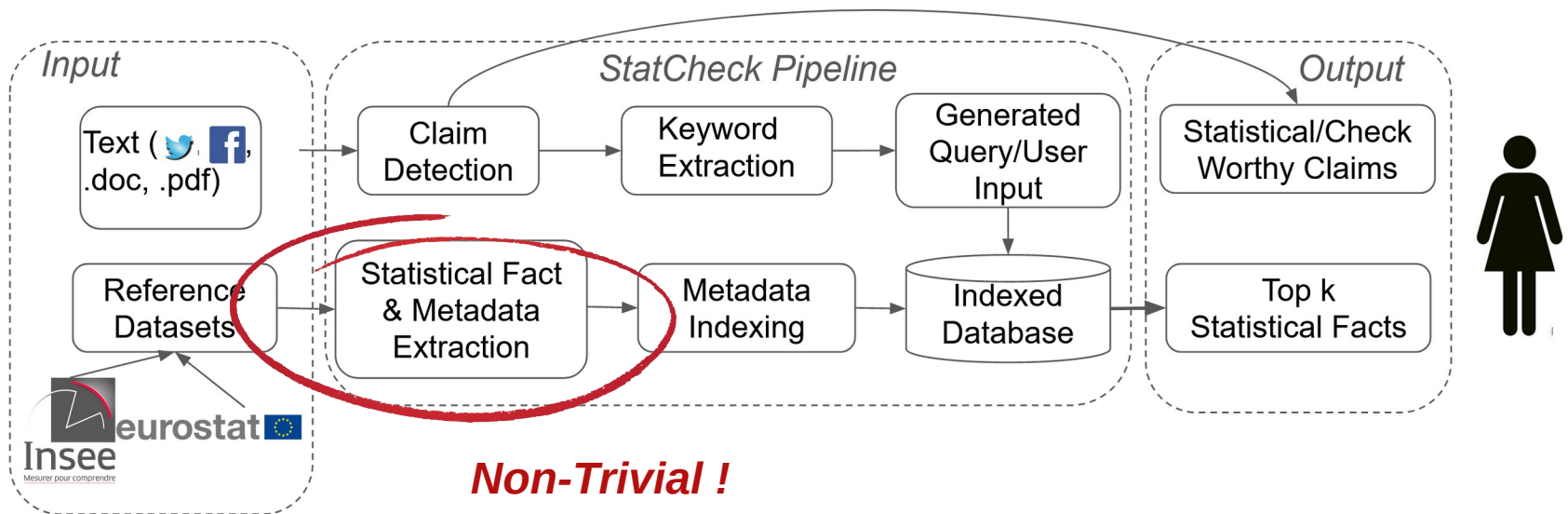
# Abstract & Motivation

1) Fact-checking is more relevant than ever

2) Computers can help us to fact-check efficiently with systems like *StatCheck*

3) To fact-check we need facts
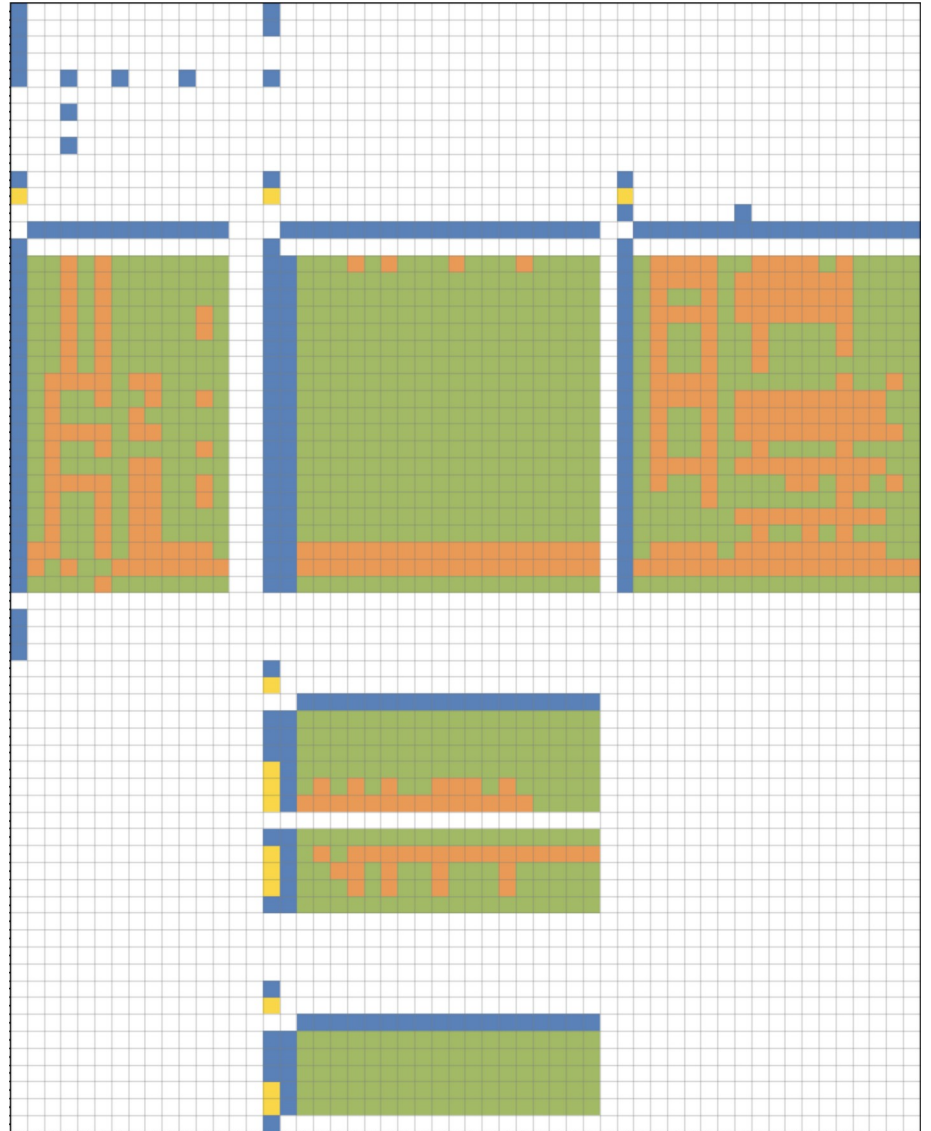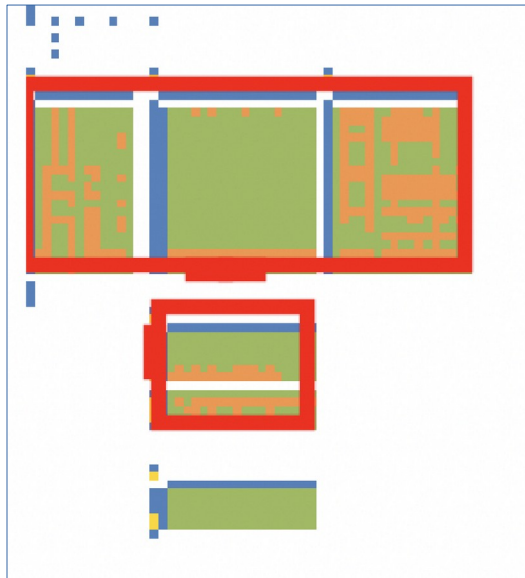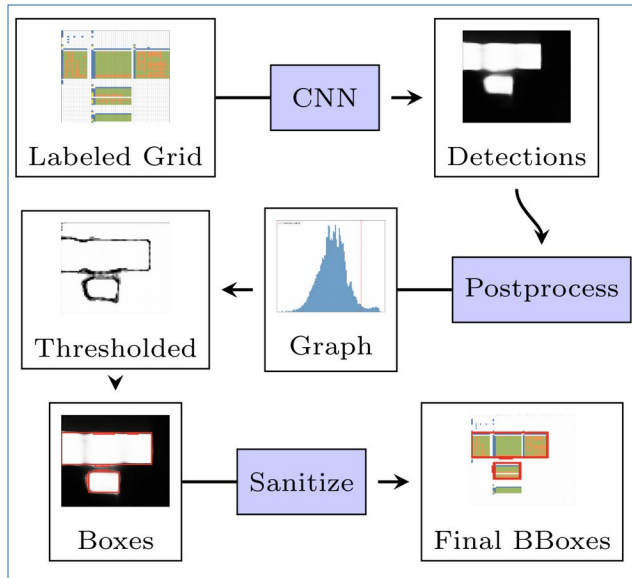


*Non-Trivial !*

# Abstract & Motivation

1. Fact-checking is more relevant than ever

2. Computers can help us to fact-check efficiently with systems like *StatCheck*

3. To fact-check we need facts

4. Trusted institutions like *INSEE* or *Eurostat* produce usable facts

5. For *StatCheck* we need to extract those facts into a structured position-aware data format

6. The Sheets were the facts are stored are made for humans and thus can be very messy

# INSEE statistical spreadsheets have very heterogeneous layouts, making it hard for heuristic-based methods to parse

# Grid-based table detection



Labeled Grid → CNN → Detections → Postprocess → Graph → Thresholded → Boxes → Sanitize → Final BBoxes

# Structure Classification