

Using alternative data sources to produce consumer price indices

Liam and Lefteris

Overview of the Alternative Data Sources Project

Liam Greenhough

Consumer Prices Methods Transformation

How price statistics are measured

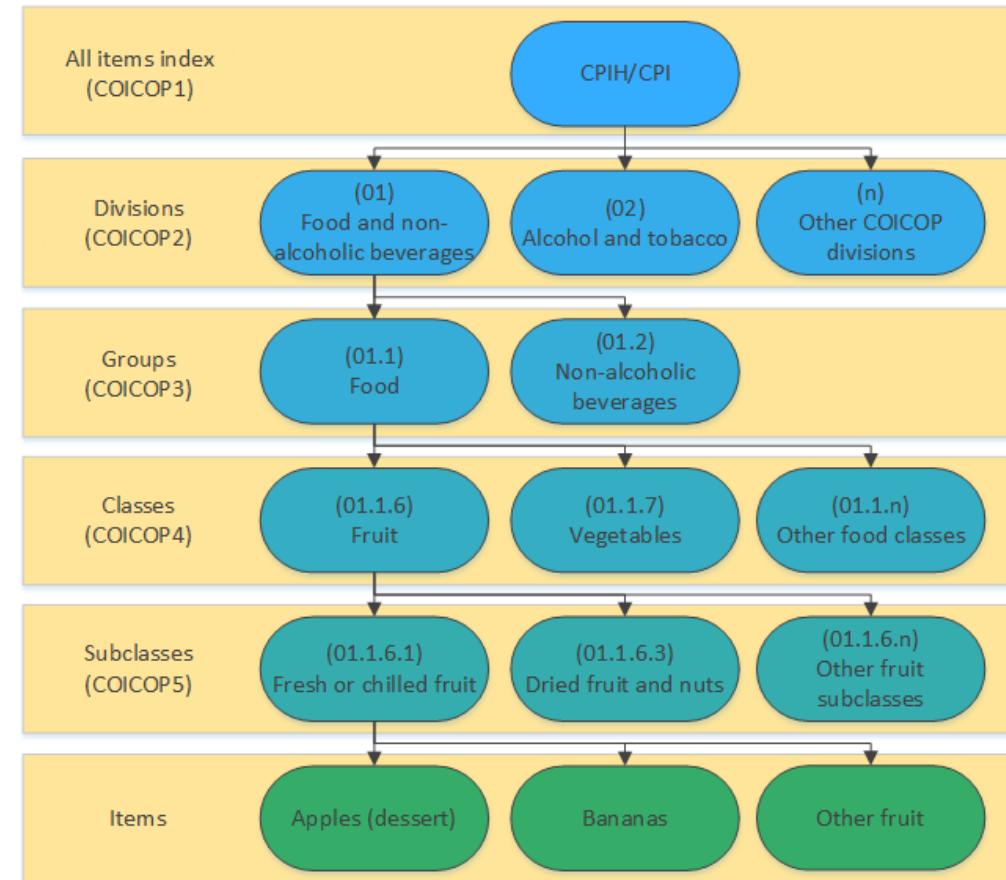
In January, select 700 “items” to track over year. Known as the fixed basket. Each year the basket is “refreshed” to account for changing consumer behaviours.



How price statistics are measured

For each item, select a group of products to track over the year.

Each item is an aggregate – but is also a “subset” of higher aggregates.



How price statistics are measured

Collect prices of products each month. These are collected:

- Locally, and
- Centrally

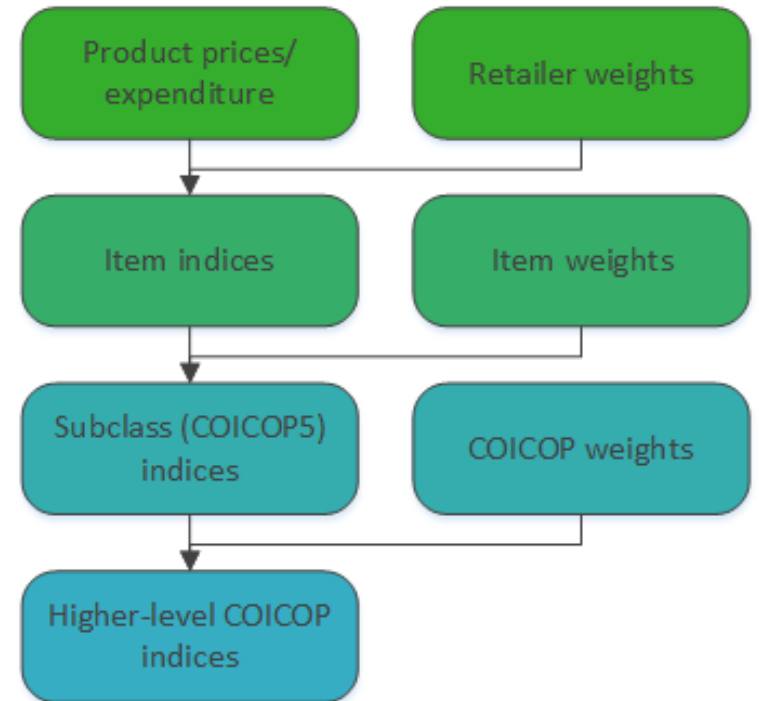


Approximately 180,000 price quotes are collected per month.

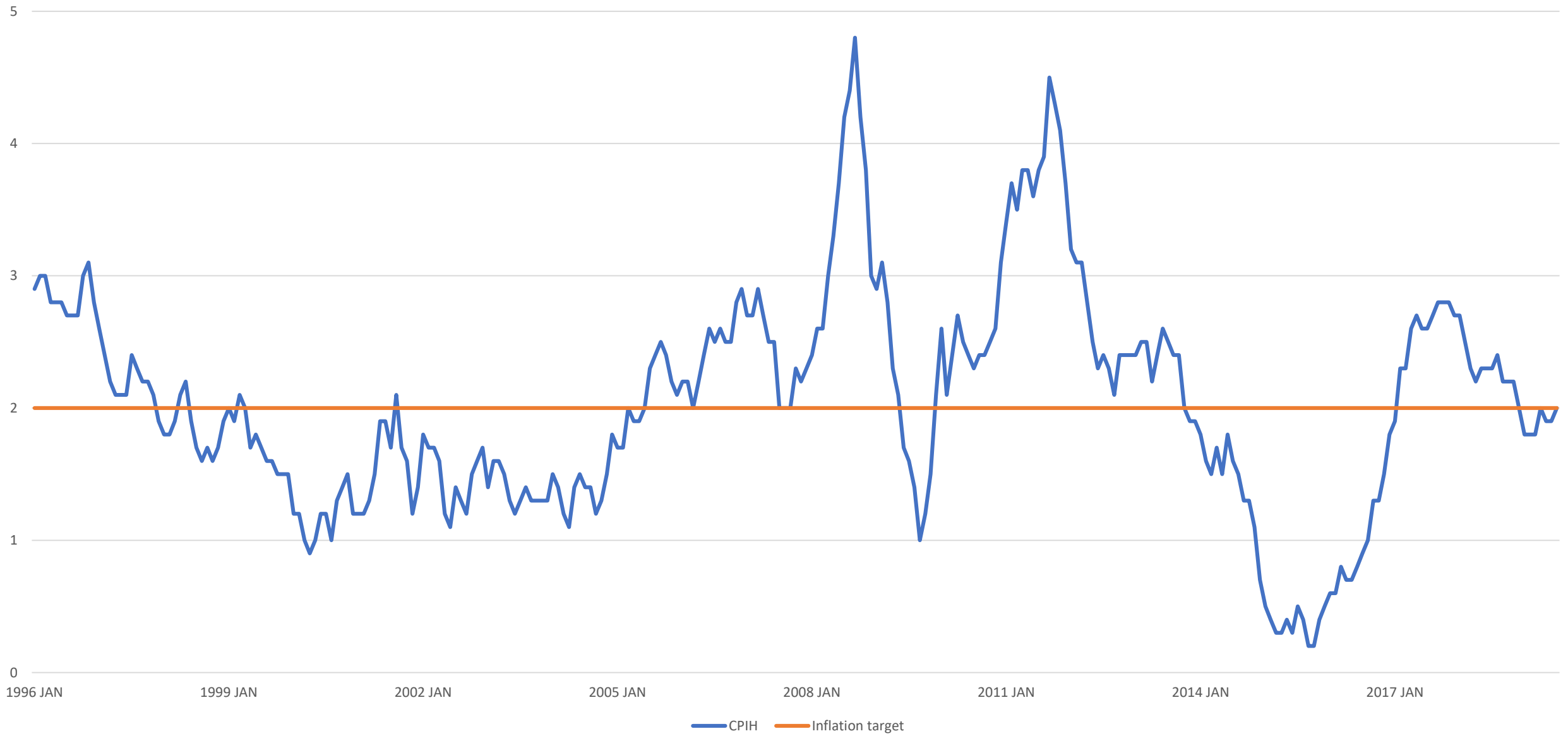
How price statistics are measured

Use index formulae to compare prices of products across months. Most common index is the Jevons.

Use weights to aggregate upwards to higher-level indices.



CPIH compared to the current Bank of England inflation target



Consumer Price Statistics: Alternative Data Sources

Alternative Data

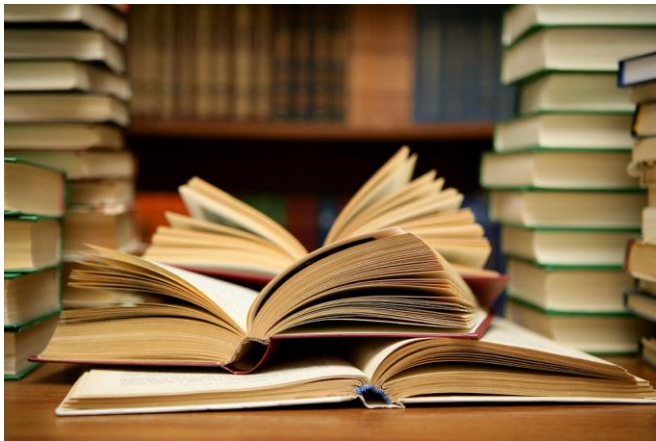
Looking to implement two new data sources:

- Scanner data – transactional data from large retailers
- Web scraped data – data scraped from online retailers

Aim to use in conjunction with traditional!



Alternative Data – targeted items



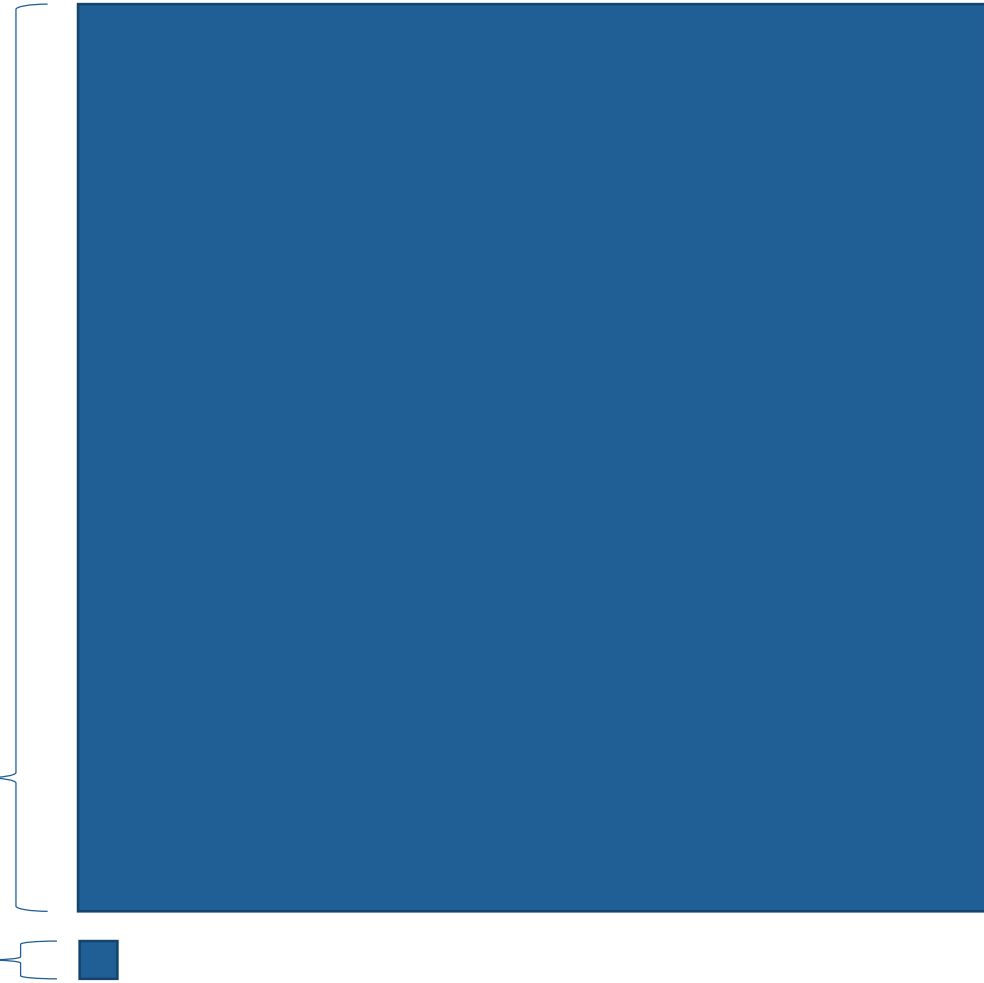
<i>Data dimension</i>	Traditional	Scanner data	Web scraping
<i>Data acquisition</i>	Manual	Automated	Automated
<i>Completeness/scope</i>	Sample from all retailers	All transactions (bulk) from medium to large retailers	Bulk or sample from online retailers
<i>Metadata</i>	Item description	Item description + limited attributes	Item description + attributes
<i>Quantity data</i>	N/A	Quantities sold	N/A
<i>Timing</i>	Single collection day	Daily	Daily

Big data

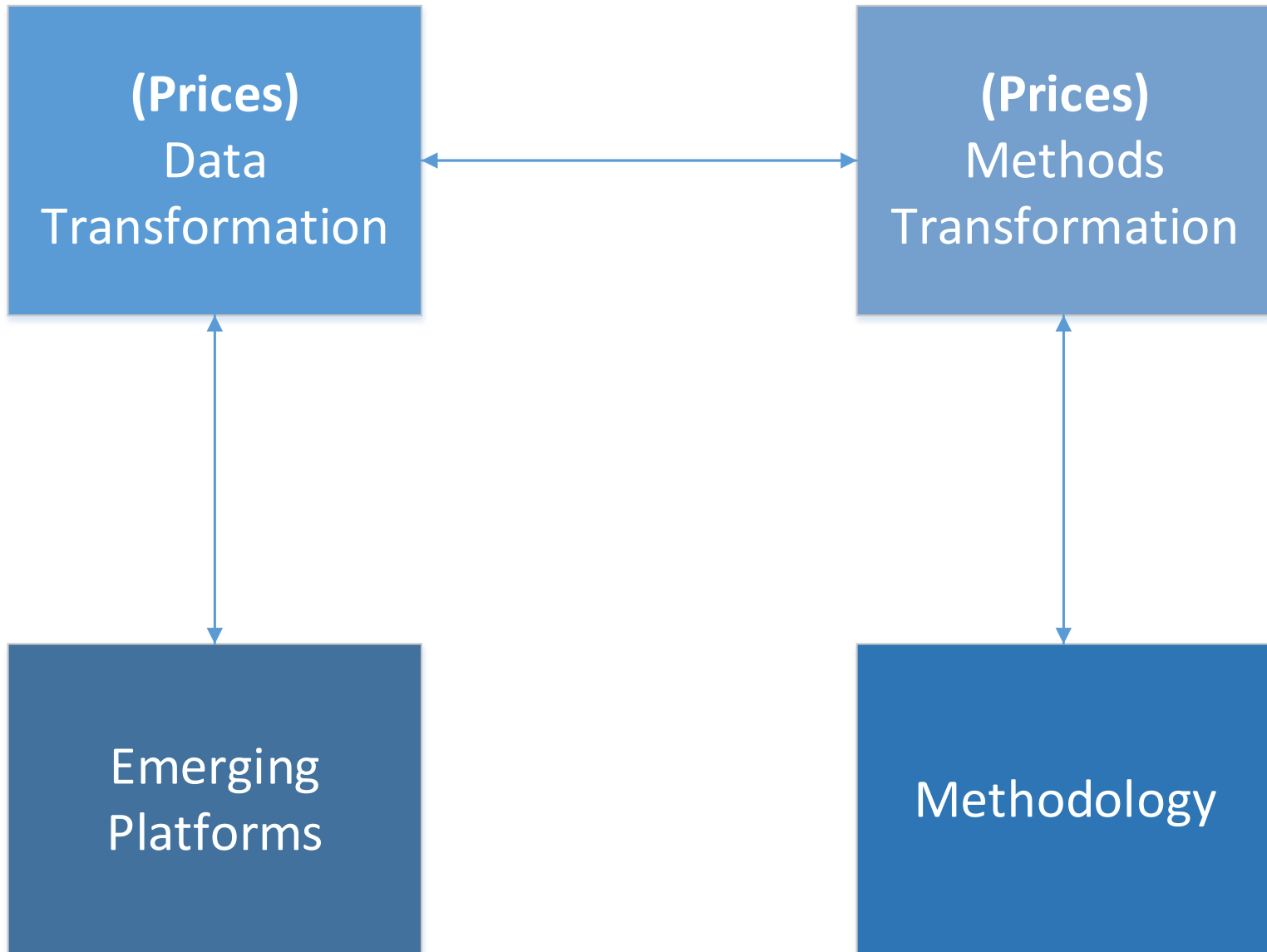
System needs to process big data.

Scanner data: ~100,000,000 price quotes per month

Traditional data sources: ~180,000 price quotes per month



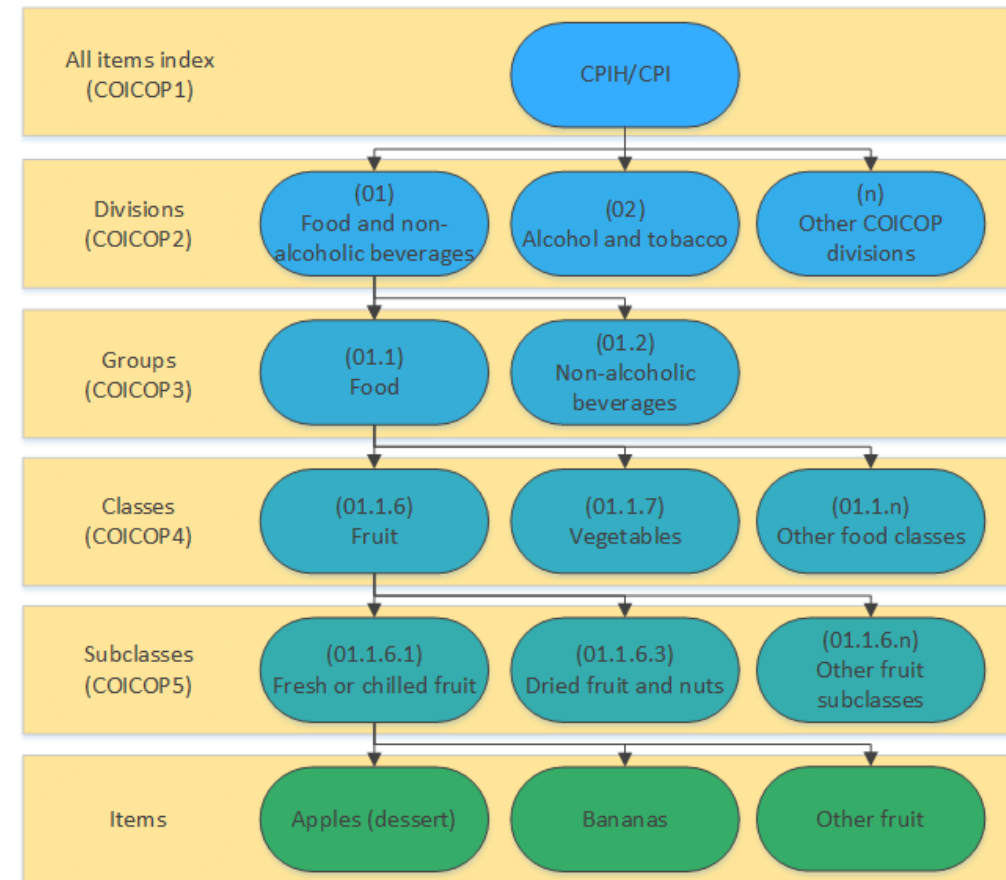
The Team



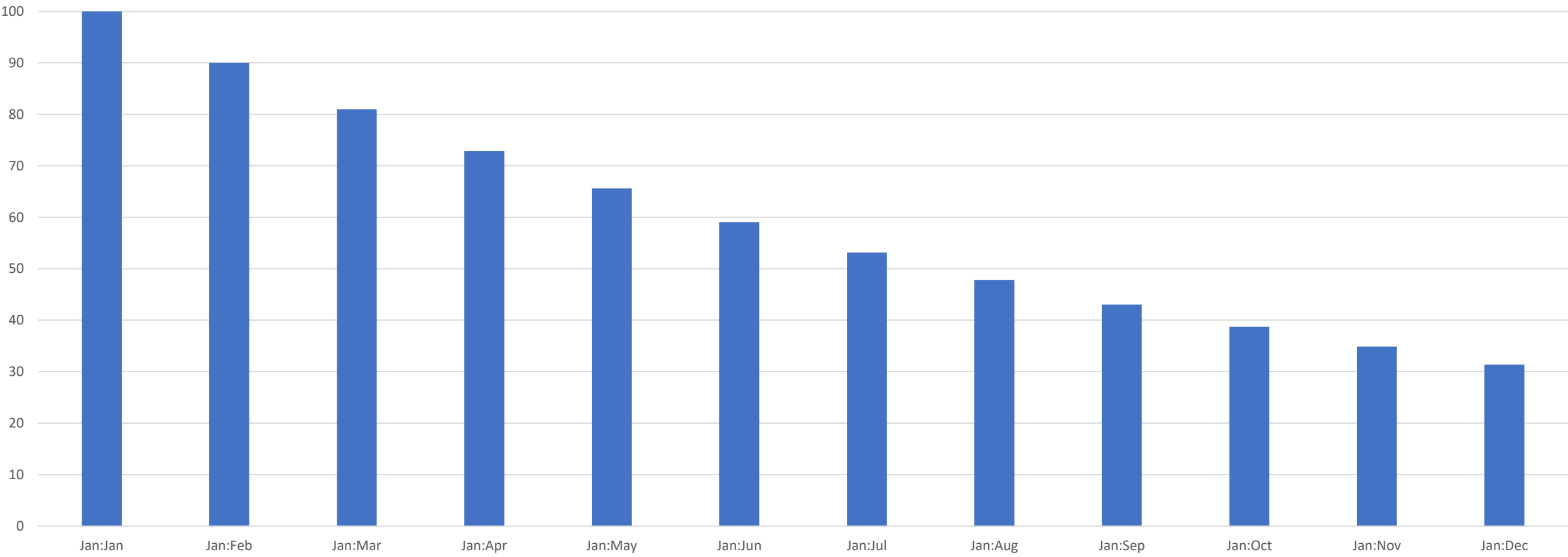
Some of the research

Scalability

Not possible to manually scrutinise big data, e.g. classification.



Product Churn – synthetic



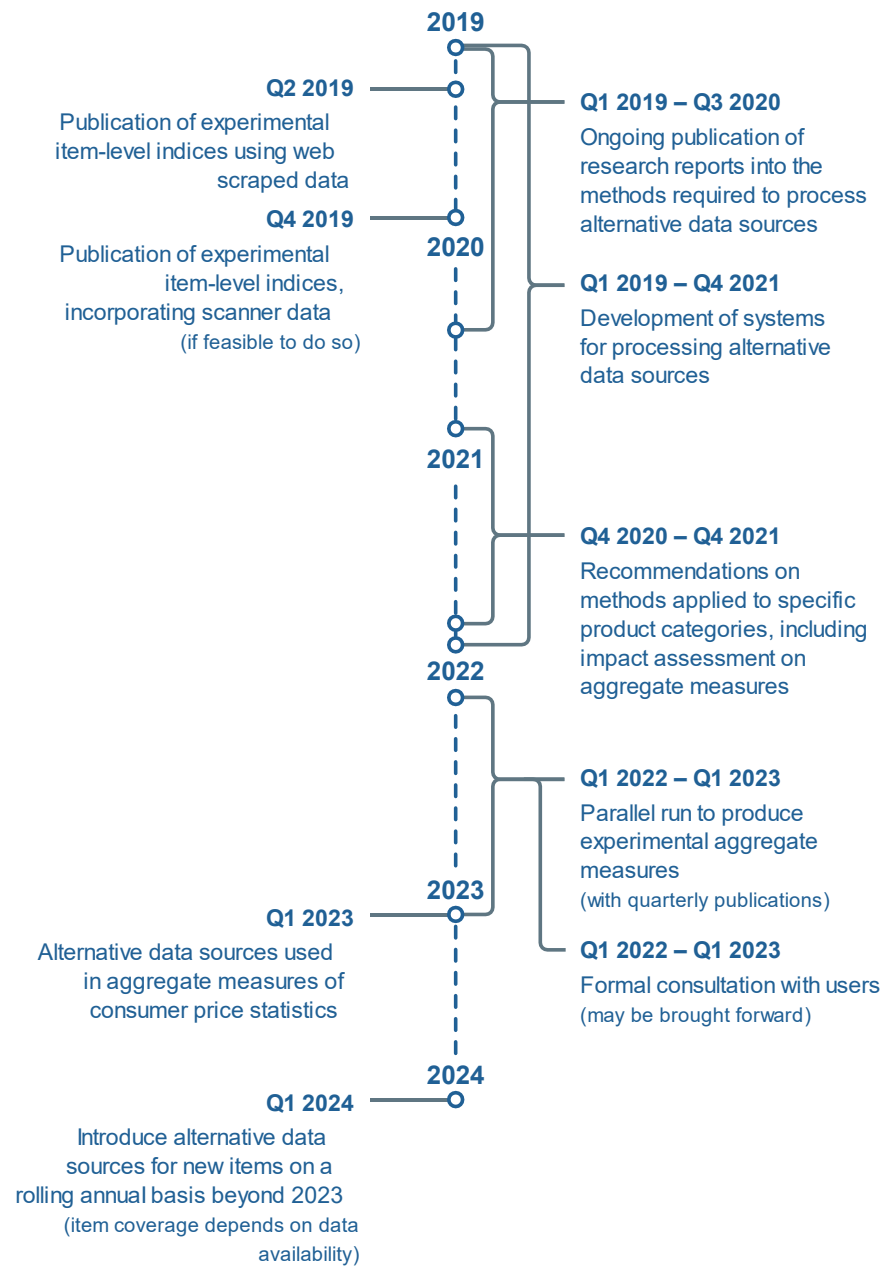
The combination problem

Lots of steps to calculate indices

Different methods at each step

Leads to many potential combinations!

Our plans



Consumer Prices Data Transformation: Development

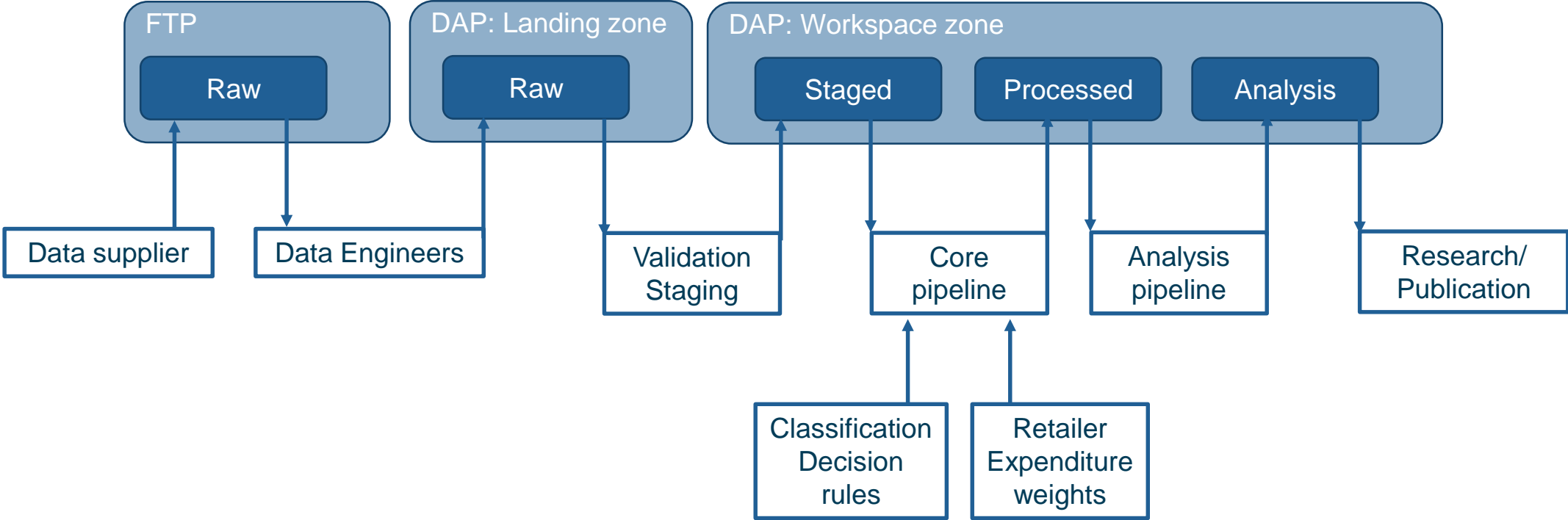
Lefteris Karachalias

Emerging Platforms Development and Support Team

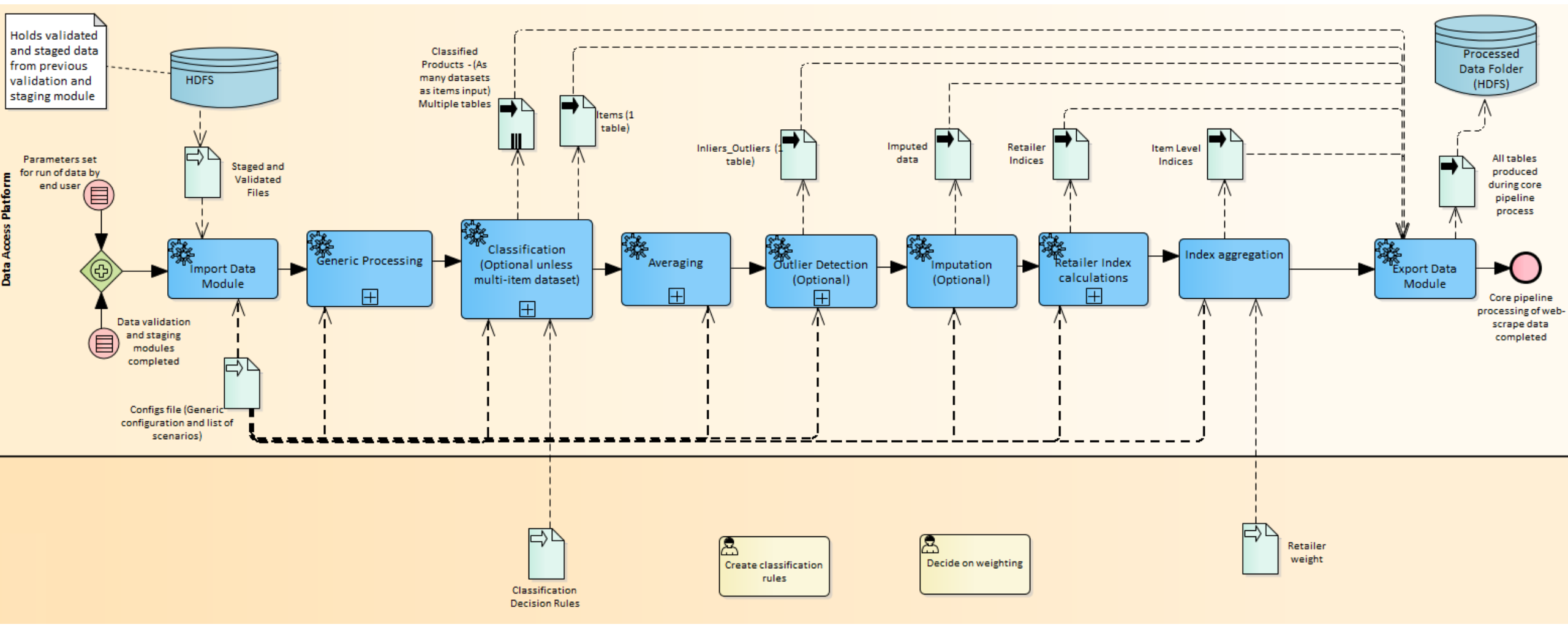
Overview

- The system
 - Overall system architecture
 - User interaction
- Development framework
 - Development project delivery team
 - Tools
 - Documentation
 - Dev&Test

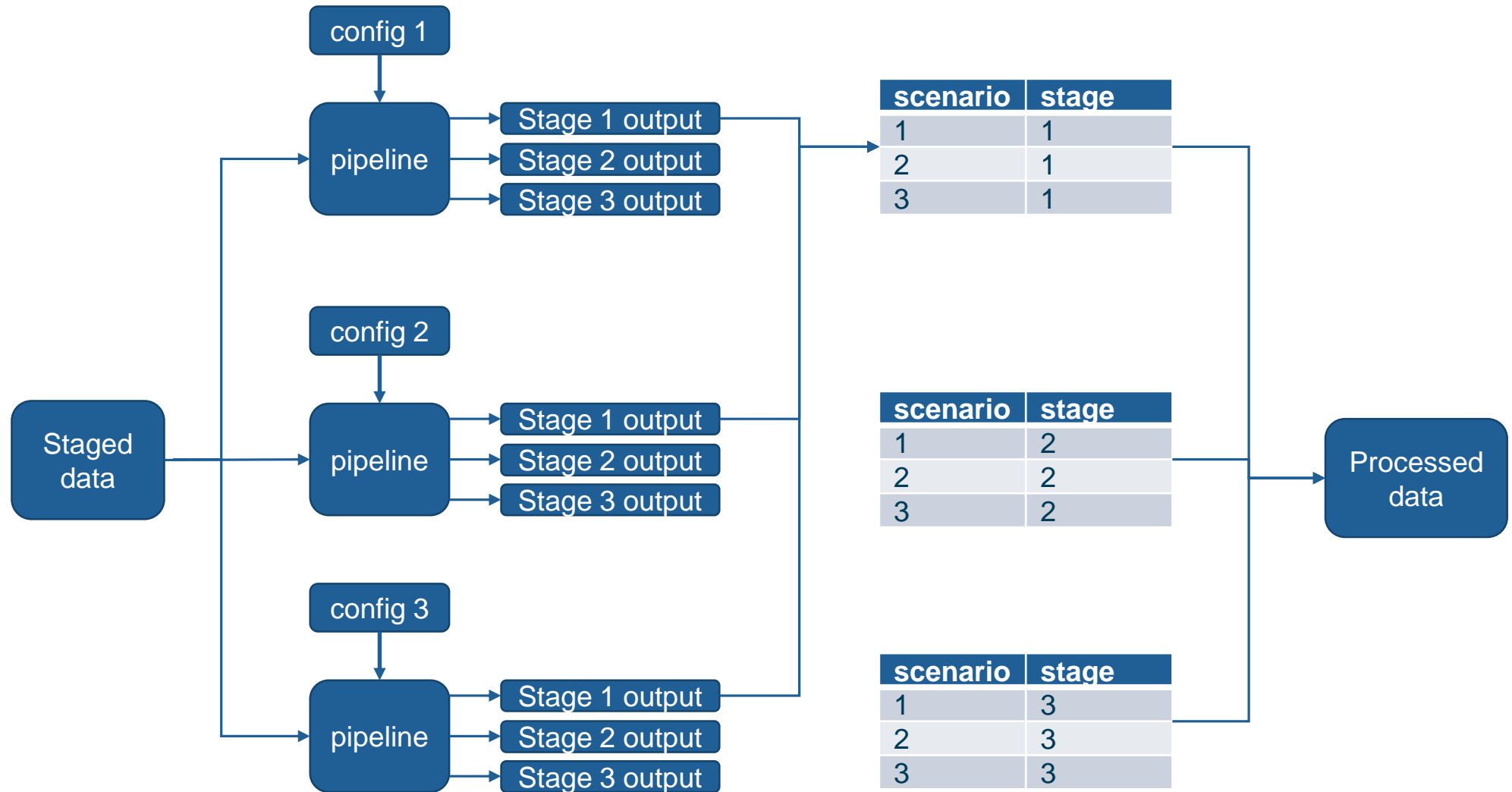
Overall system architecture



Core pipeline



Multiple configuration scenarios



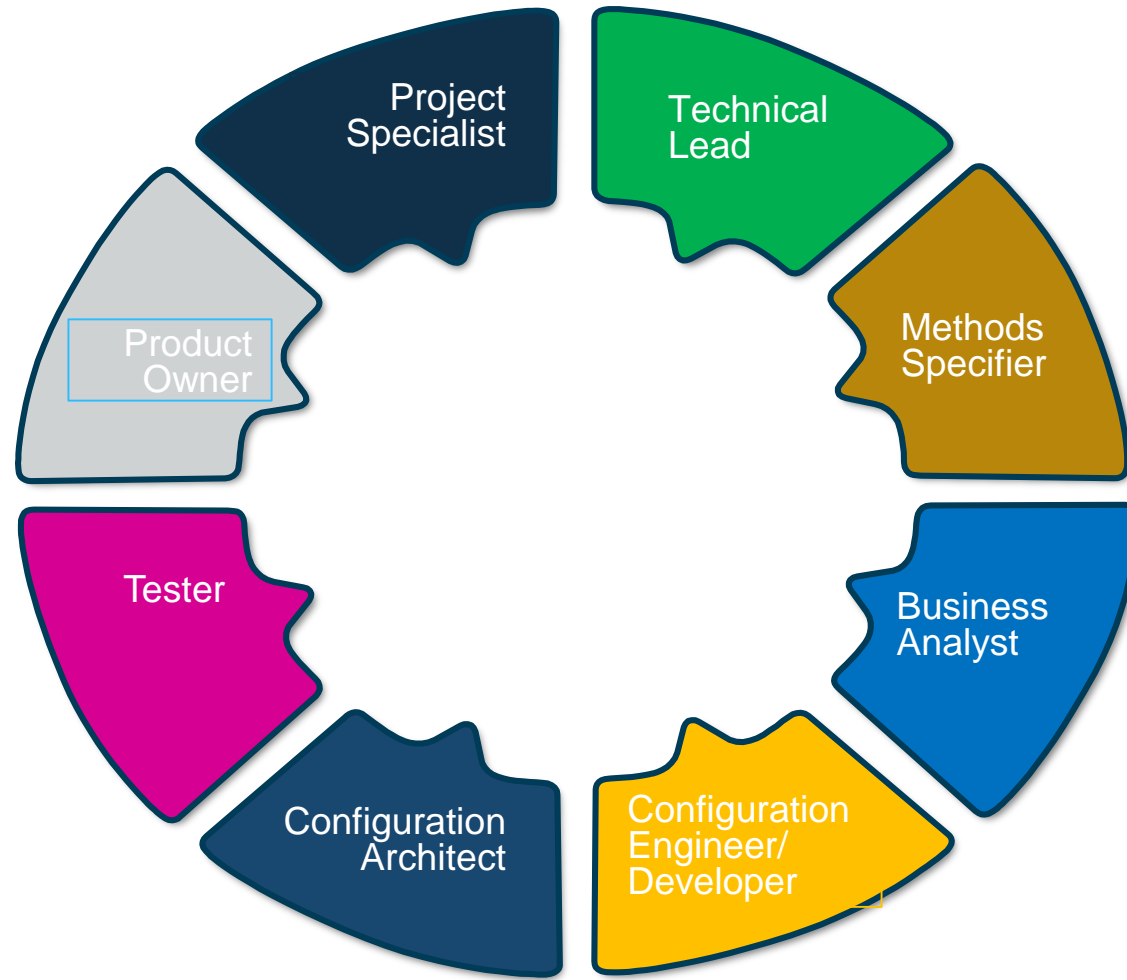
User interaction

- UI: CDSW / HUE
- Manual
- Configuration
- Mappers (BAU)
- Dashboard (cannot share VDI)
- Output tables

Development framework (1)

- Project delivery team
- Development phase: Between Discovery and Alpha phase
- Agile, Jira
- DAP, PySpark, HDFS, HIVE (sensitivity)
- Git, GitLab

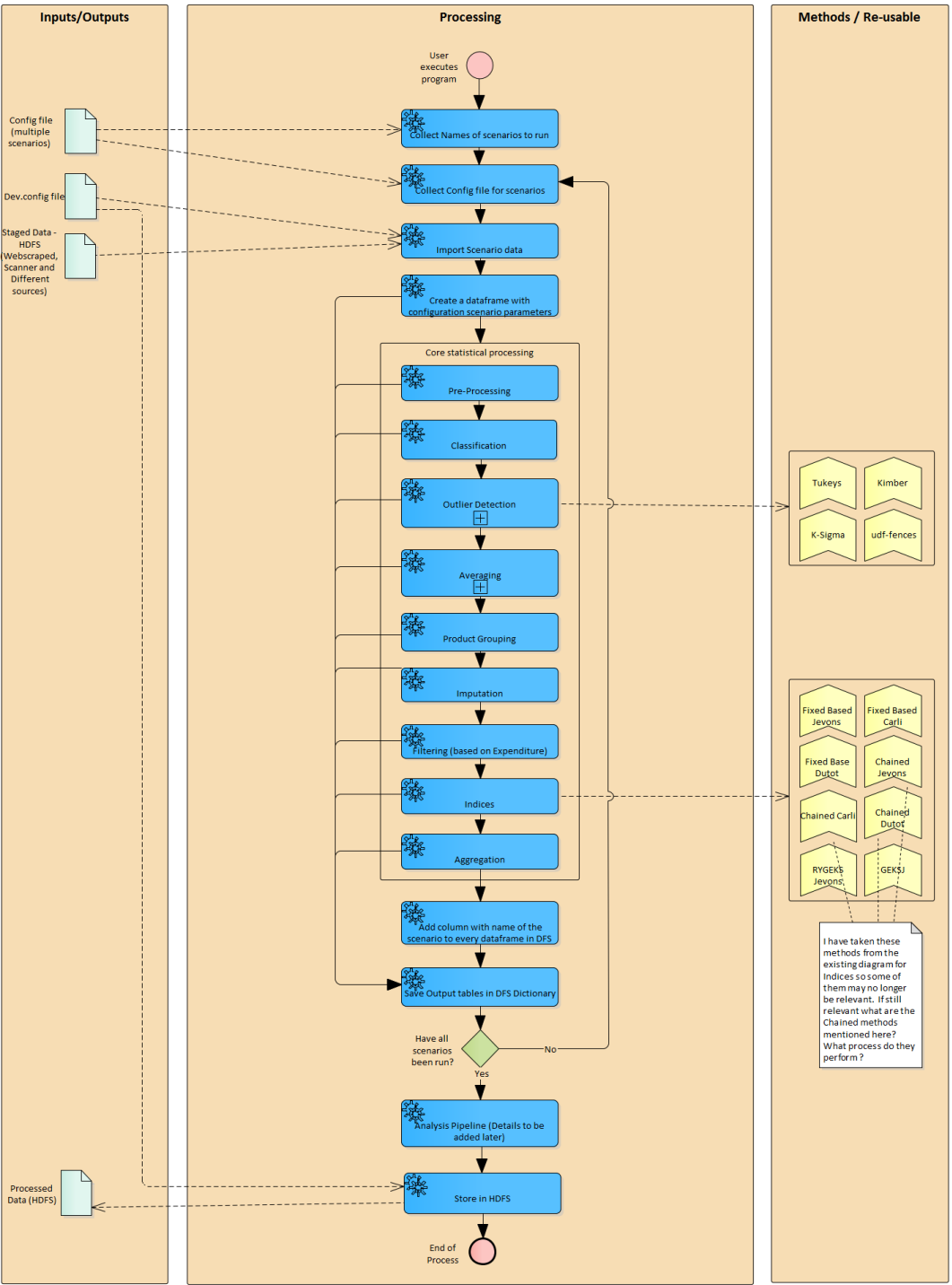
Project delivery team



Development framework (2)

- Unit testing, CI with Jenkins, UAC
- Documentation Sphinx, user manuals
- Business Analysis models, Sparx
- Business Architecture: pushing to the SML
- Synthetic data, Dev&Test, packaging

Statistical process model

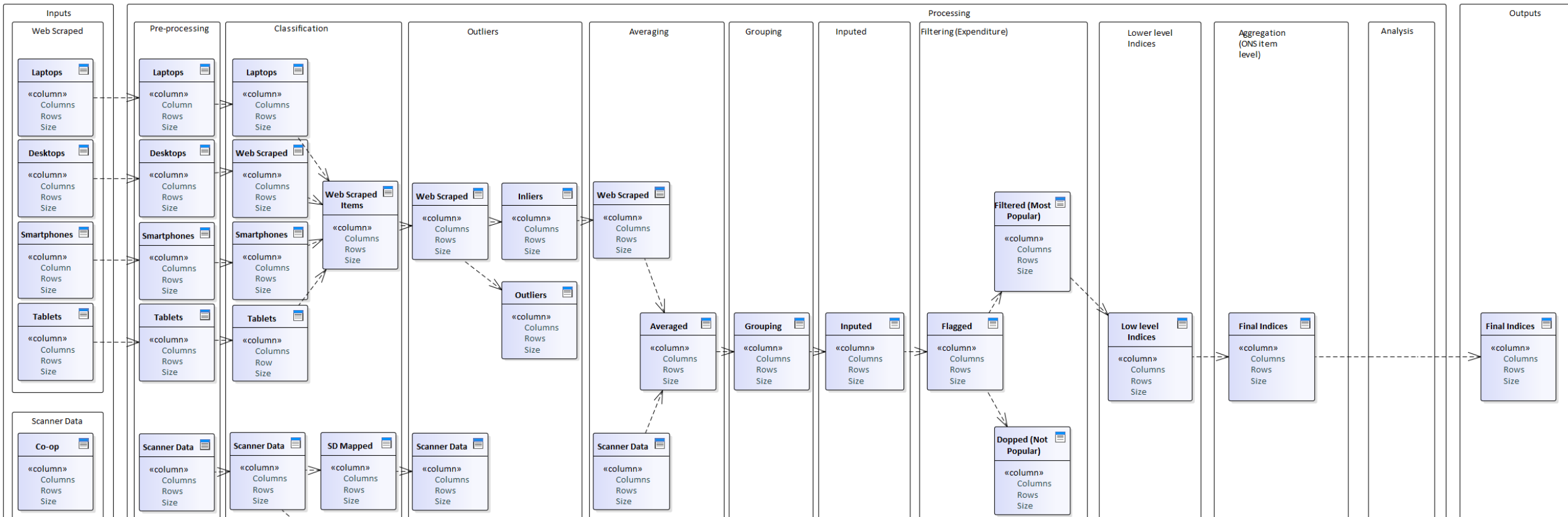


Data (journey) model

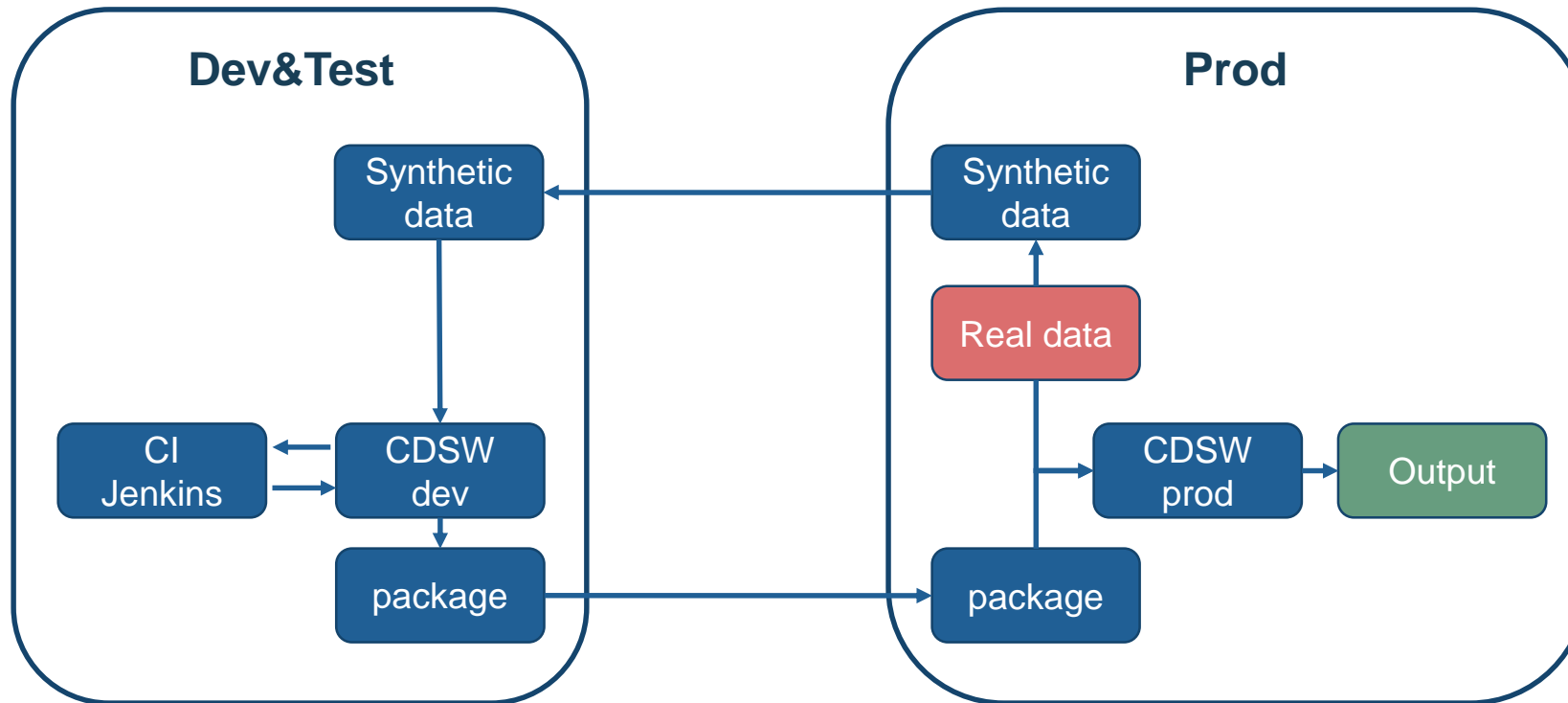
Input

Processing

Output



Dev and Test environment



Thank you!

Any questions?