

Approaching the challenges with the new data eco-system

Karin Blix, Quality Coordinator, kwb@dst.dk
Statistics Denmark



Challenges of the new data eco-system

Examples From Denmark

- Long tradition of using administrative data sources for official statistics.
- This presentation will elaborate on:
 - how the Danish statistical system relies on administrative sources and
 - how this is also a base for exploring new data sources and
 - what is done to ensure quality in official statistics in the new data eco-system.

Quality in statistics

- ISO, an international body for formulating standards, has defined quality as: “Degree to which a set of inherent characteristics fulfils requirements”.
- Quality in statistics is the set of properties the statistics needs to have to fulfil the users’ needs.
- Our starting point is the users and their needs.

Fit for purpose

Administrative data



Data which are originally collected by **public authorities** for their own purposes

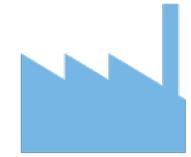
A blurred image of a data table with columns of numbers. The numbers are arranged in a grid-like structure, typical of a spreadsheet or database table. The text is out of focus, but some numbers are visible.

Usually organised and structured in **administrative registers**



Administrative registers are a **valuable asset** also for producing statistics

Data flow from data owners to end-users



4

Data are collected in registers as by-products of administrative routines



5 & 6

Producing statistics

Creating knowledge

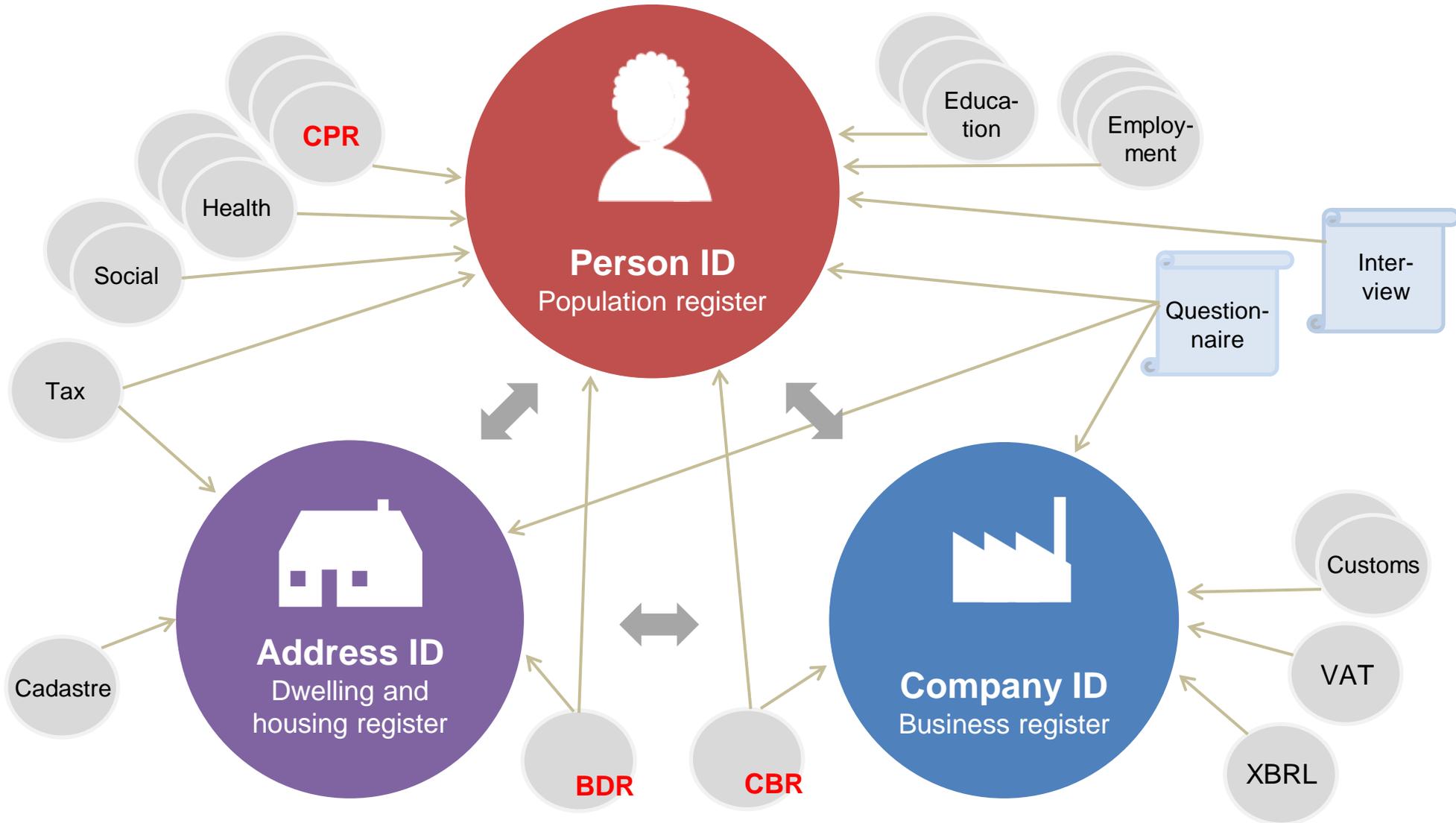


7

Disseminating statistics to users in Society



Statistical information system



Examples of exploring data sources

Administrative data

- Life Lines is an example of exploring existing administrative sources for new uses
- Networks of opportunities another example exploring existing administrative sources

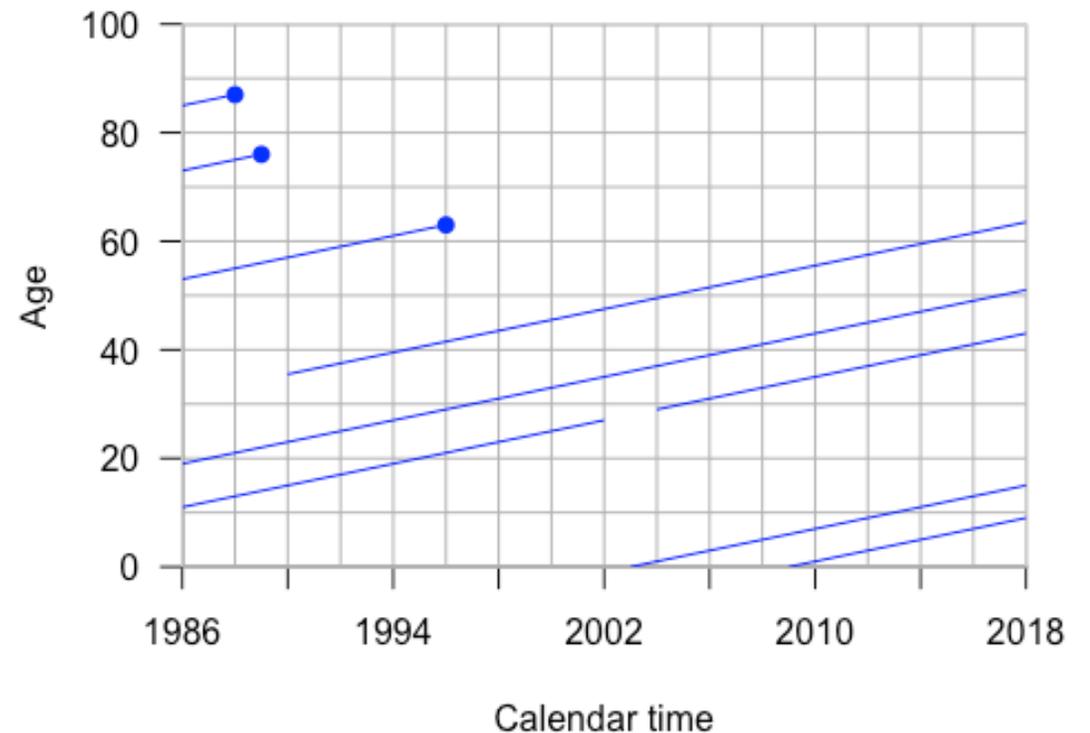
New data sources

- The use of scanner data for price statistics has been in regular use for many years
- AIS position data for ships is in use for experimental statistics on port calls
- Data from smart-meters on use of electricity is being explored for several uses.

Examples – use of admin data - population

Life-lines

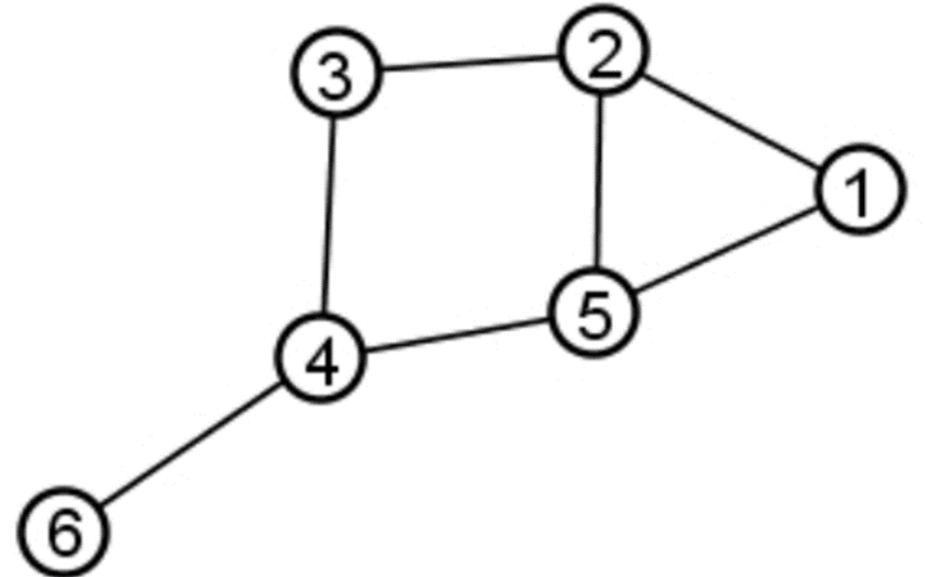
- The longitudinal Register (Life lines) is updated annually and shows the individuals that have been part of the Danish population through their life courses.
- Based on the population statistics register: all persons in Denmark from 1986 to 2021 plus an extraction from the CPR (Central Person register) from 1968.
- One line for each person's presence in the population, e.g. as a period from birth to death or from the person's immigration to emigration/death
- The purpose: form a population based on life lines, e.g. extraction of cohort, period, age, length of life line, etc.



Examples - use of Life Lines for networks

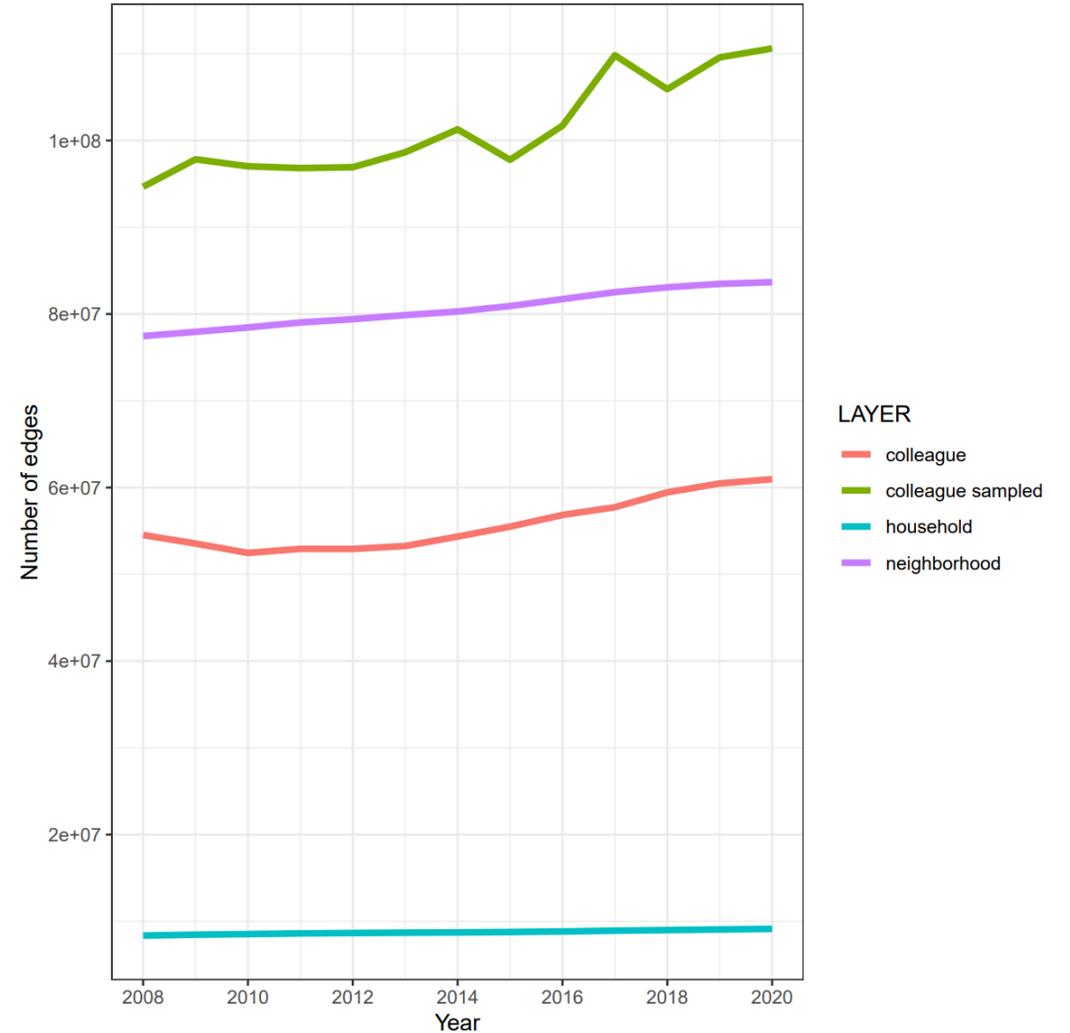
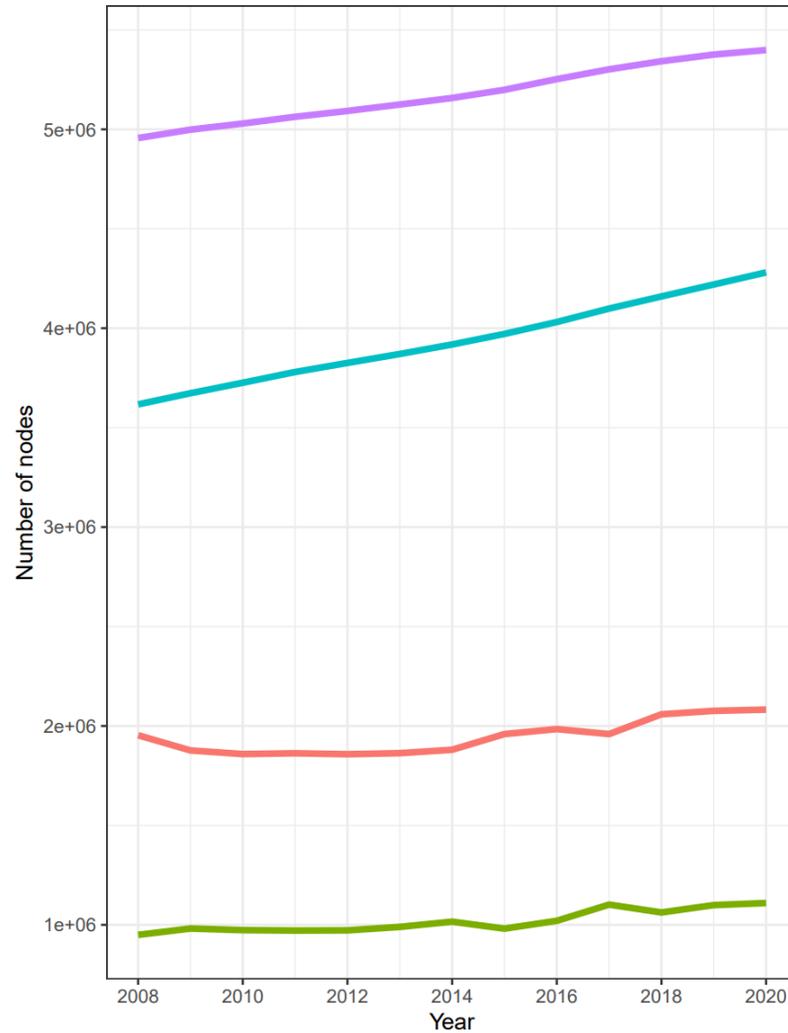
Network of opportunities

- Based on Life-Lines
- Networks composed of
- Nodes (individuals)
- Edges (relations between individuals)
- Degree (how many relations one individual has – edges per node)
- Statistics on number of family members, colleagues, class mates etc.
- Segregation (how many one potentially knows that has a university degree, two vacation homes, a mother from Norway etc.)
- How something can spread through a network (sickness, education, employment etc.)



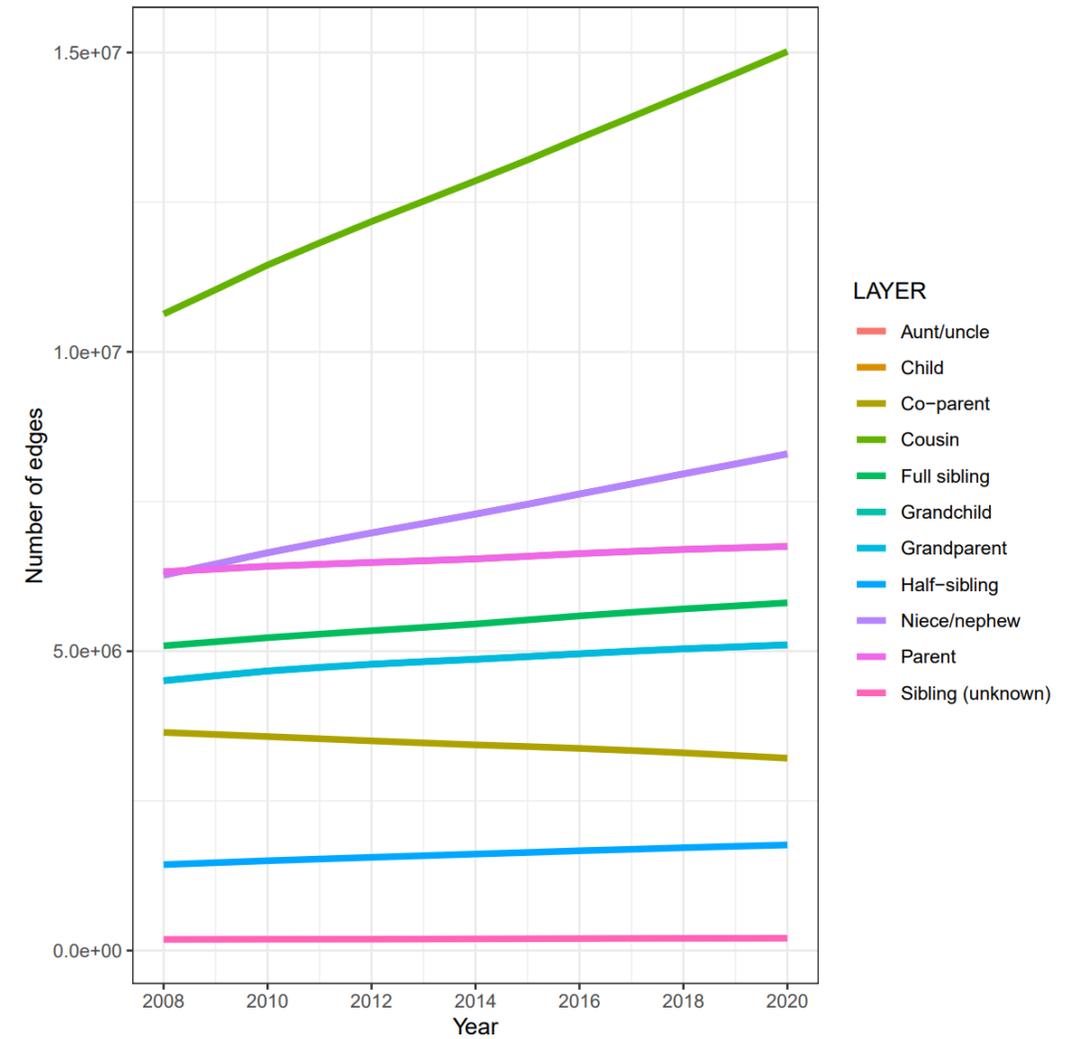
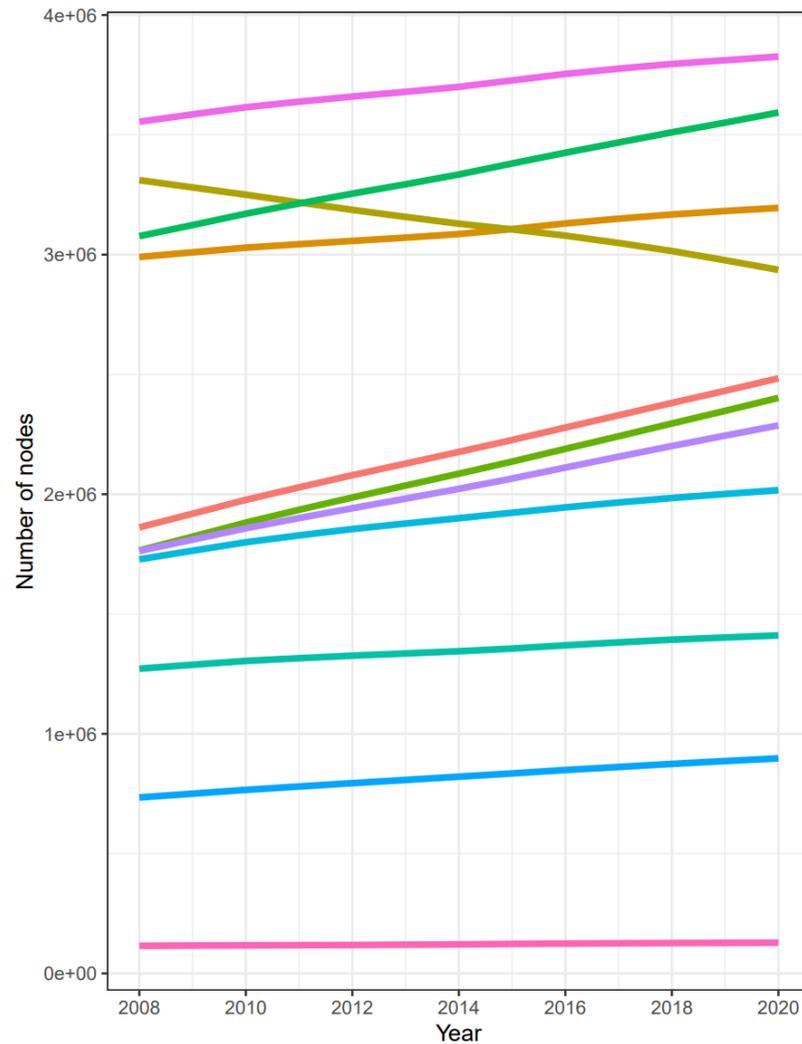
Example – networks - descriptives

Descriptives



Example – networks

Family



Example – use of electricity data

What we have

- Database registry of construction projects
- Access to smart-meter electricity data on address level
- Machine learning estimation model correcting delay
- Quarterly statistics
- Issues
 - Late and incomplete registry regarding projects

New approach

- Combining electricity data with data about known construction projects
- Exploring patterns of interest e.g. identifying construction phases
- Pilot project
 - Find a sample project and follow the path through the registry

Example – use of electricity data

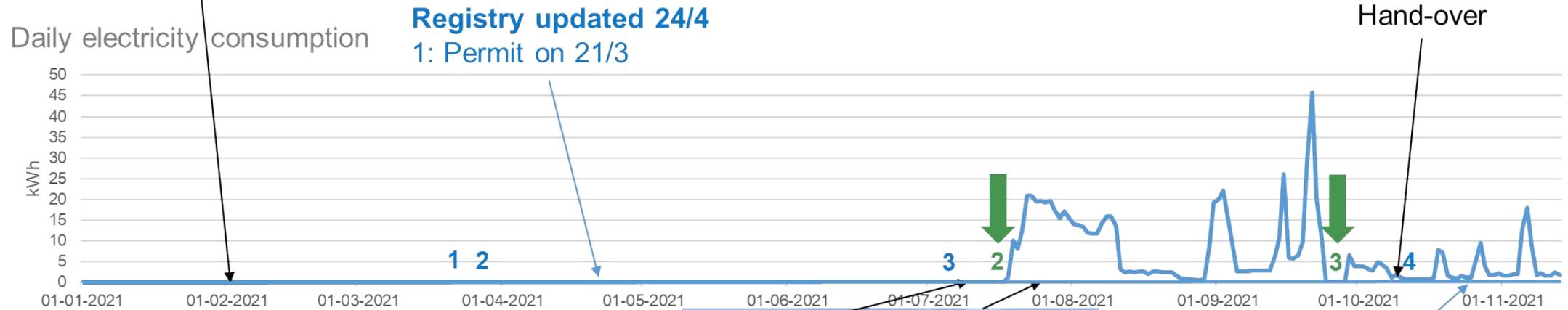


Purchase

Smart-meter data: available with 8 days delay

New agreement from 15/06 appears in July

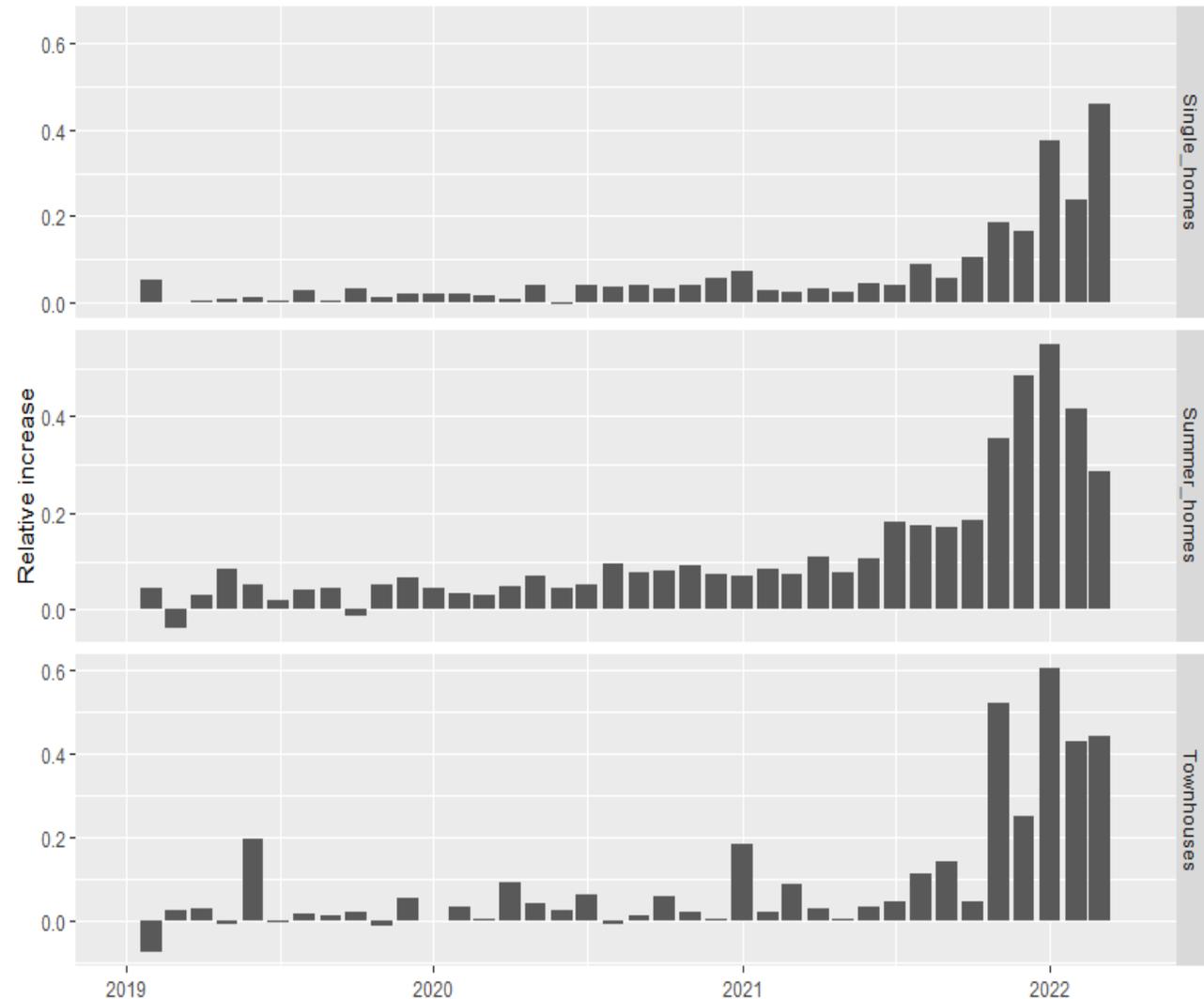
Use of electricity begins on 18/7



Example - Future work - electricity data

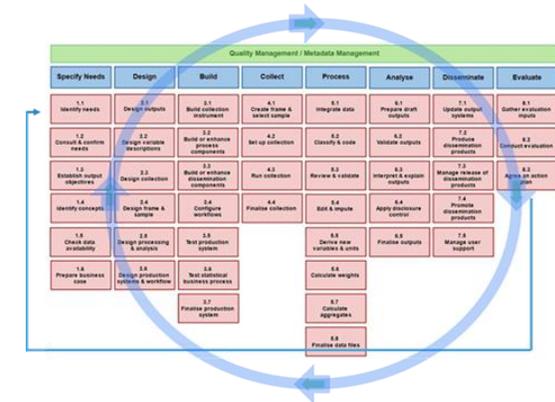
- Immediate potential
 - Using electricity data as proxy
 - Development of new indicator
 - Now-casting statistic
 - Pattern recognition for identifying more projects – data science lab

- Other uses of smart-meter data
 - Charging of electric cars
 - Identifying empty dwelling
 - Assessing use of summer homes
 - Forecasting bankruptcies

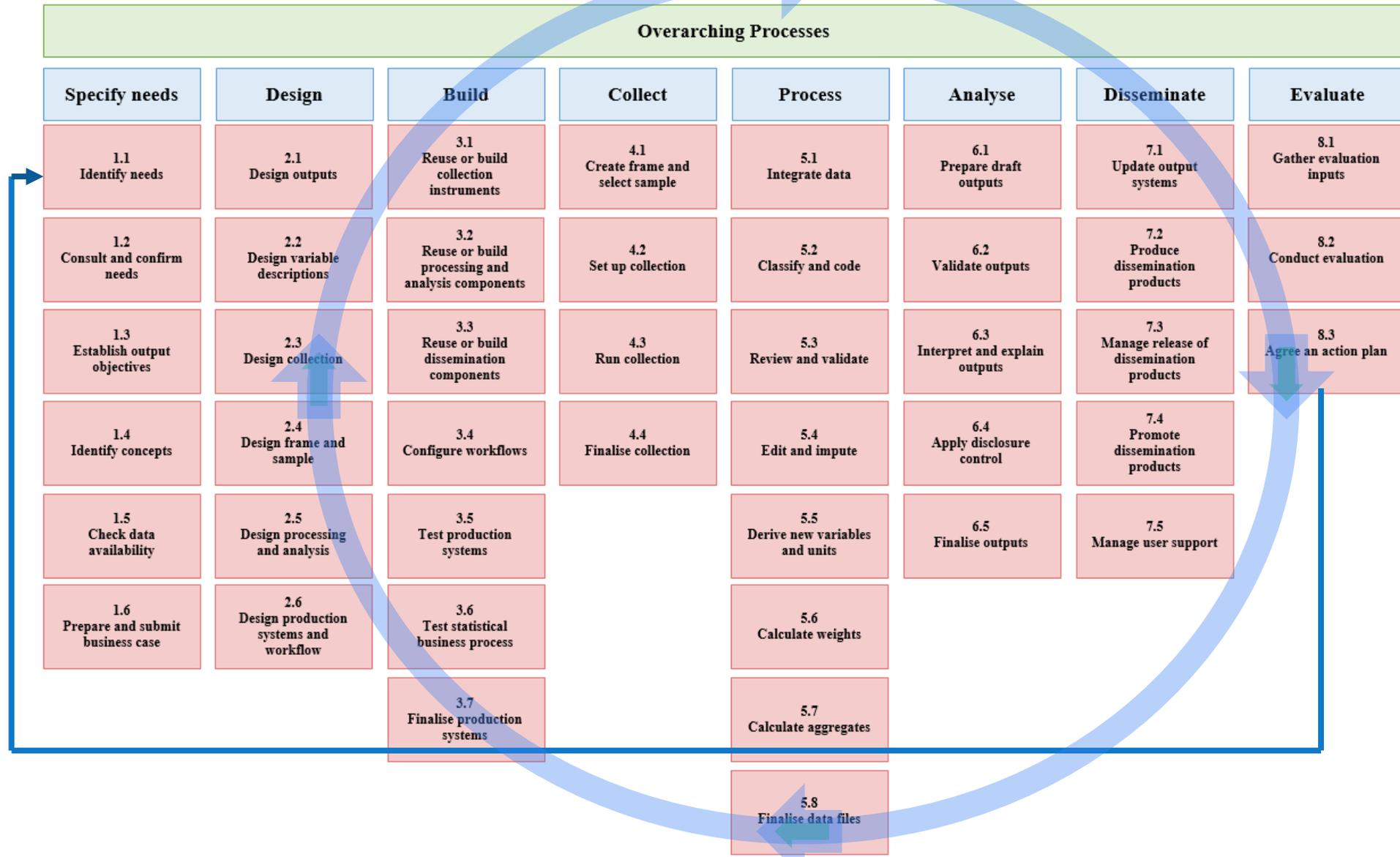


How we ensure quality

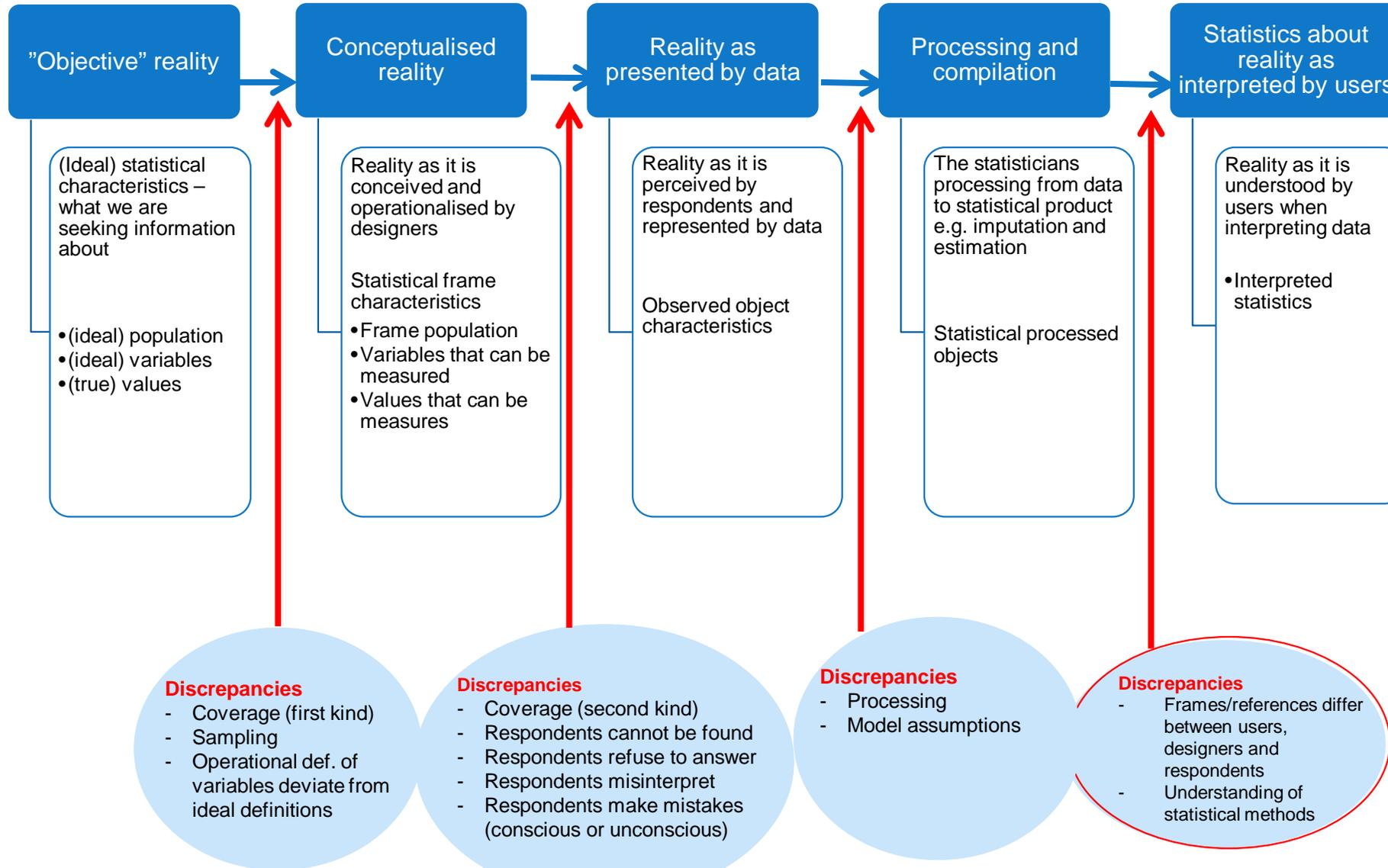
- Cooperation with data owners
- Cooperation with data providers (e.g. government agencies, municipalities, educational institutions etc.)
- Methods unit in Statistics Denmark
 - Sample selection
 - Standardised and modernised error detection
 - Other methodological issues
- Quality unit in Statistics Denmark
 - Quality assurance of documentation of statistics (quality reports)
 - Quality reviews
 - Standardisation – e.g. the process model (GSBPM)
 - Coherent metadata system
- Data science lab in statistics Denmark
 - Exploring existing data sources
 - Exploring new data sources
 - Exploring new methodology



Statistical processes (GSBPM 5.1)



Production of statistics



Two types of metadata

Structural metadata

- Used to identify statistical data
- Headlines, variable names, unit of measure, reference time etc.
- Must go together with statistical data
- Impossible to interpret statistics without it

Reference metadata

- Describes content, statistical processing, relevance etc.
- Can be detached from the statistical output
- Quality Reports is a type of reference metadata
- ...so is methodological metadata

Statistics without metadata

	2 881 620
	2 908 337
<hr/>	
	2 868 172
	2 976 785

...with structural metadata

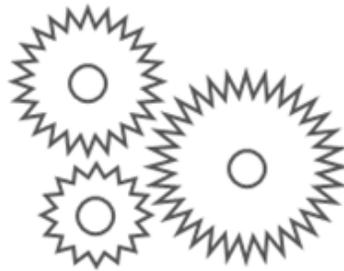
Population	
All Denmark	2018Q3
Men	2 881 620
Women	2 908 337
Unit : number	
<hr/>	
Real estate market value	
One-family houses	2016
Brøndby	2 868 172
Vallensbæk	2 976 785
Unit : Average Market value (DKK)	

...and reference metadata

Population	
All Denmark	2018Q3
Men	2 881 620
Women	2 908 337
Unit : number	
<hr/>	
Real estate market va	
One-family houses	
Brøndby	
Vallensbæk	
Unit : Average Market value	



Documentation of statistics - content

Introduction 	Statistical presentation 	Statistical processing 	Relevance 
Accuracy and reliability 	Timeliness and punctuality 	Comparability 	Accessibility and clarity 

User profiles

■ General population

- Everyone in contact with SD through the flow of news
- Mr and Mrs Smith
- People interested in social affairs

■ Specially interested parties

- Actively searching for facts
- No special qualification
- E.g. journalists, students and politicians

■ Professional users

- Systematically use figures from SD
- Can combine and extract data
- E.g. specialists, trade and business press and public servants

■ Analysts

- Awareness of statistics
- Can extract and process complex data
- E.g. researchers, large-scale consumers and data analysis units



[← Documentation of statistics](#)

- Climate footprint (experimental statistics)** ^
- Statistical presentation
- Statistical processing
- Relevance
- Accuracy and reliability
- Timeliness and punctuality
- Comparability
- Accessibility and clarity

SHARE THIS PAGE






Climate footprint (experimental statistics)

The purpose of the statistics is to show the investment and emissions of greenhouse gases from the supply chains for Danish final use. Global emission constitutes Denmark's share of the global emission which uses it for the production of goods and services. Global Afrappro

Statistical presentation

The statistics show the amount of greenhouse gas that has been emitted in the supply chains for Danish final use. The emissions are distributed by type of final use, emitting industries and countries, as well as by supplying industries.

[Read more about](#)

Statistical processing

The climate footprint statistics are distributed internationally by supplying industries.

[Read more about](#)

Relevance

The climate footprint statistics show global emission

[← Climate footprint \(experimental statistics\)](#)

- Statistical presentation**
- Statistical processing
- Relevance
- Accuracy and reliability
- Timeliness and punctuality
- Comparability
- Accessibility and clarity

SHARE THIS PAGE






Statistical presentation

The statistics show the amount of greenhouse gas that has been emitted in the supply chains for Danish final use annually from 1990 onwards. The emissions are distributed by type of final use, emitting industries and countries, as well as by supplying industries.

Data description

The statistics show the amount of greenhouse gas that has been emitted in the supply chains for Danish final use. The emissions are distributed by type of final use, emitting industries and countries, as well as by supplying industries.

The calculation of the climate footprint uses 100-year Global Warming Potentials from the IPCC's fourth assessment report (AR4) to convert tonnes of a given greenhouse gas into tonnes of CO2 equivalents.

The supply chain for a type of final use is defined in these statistics as *all the production activities in Denmark and the rest of the world that have been necessary to produce the products for final use*. The supply chain behind e.g. milk includes both raw milk production and further processing, the production of dairy cows and feed for them, the production of electricity to run the stables and dairies, as well as steel and wood to build the stables and dairies, etc.

The emissions are calculated in tonnes of CO2e (CO2 equivalents) and include the greenhouse gases CO2, CH4 (Methane), N2O (Nitrous oxide) and F-gases (SF6, HFC-gases and PFC-gases). In relation to LULUCF (unfccc.int) the climate footprint only includes emissions from land use in the agricultural sector.

The statistic has five variables: - Types of use: The type of Danish final use whose supply chain led to the greenhouse gas emissions. - Supplying industry: The Danish industry that formed the last link in the supply chain for the final use. Imports from foreign industries directly for final use are entered under the item "Imports for final use" so that the statistics are fully comprehensive. (Supply industry is only included in the AFTRYK2 table) - Emitting industry: The industry where the production that emitted the greenhouse gas took place. - Emitting country: The country where the production that emitted the greenhouse gas took place. - Year: The year of final use.

Contact info

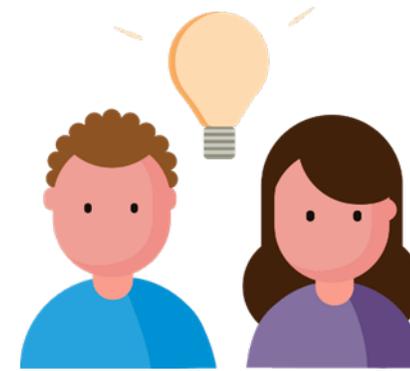
National Accounts, Economic Statistics
 Peter Rørmose Jensen
 +45 3917 3862
 prj@dst.dk

Get as PDF

[Climate footprint \(experimental statistics\)](#)

Conclusion

- At Statistics Denmark, we are constantly discovering new applications and gains from having a single system based approach to metadata reporting and storage.
- In combination with a dedicated and expanding team of quality practitioners principles of...
 - relevance,
 - accuracy & reliability,
 - timeliness & punctuality,
 - coherence & comparability and
 - accessibility & clarity
- ...are becoming routine and synonymous with the way we do business on a daily basis.



Thank you for your attention

