# Future Advanced Data Collection

## The Future is NOW
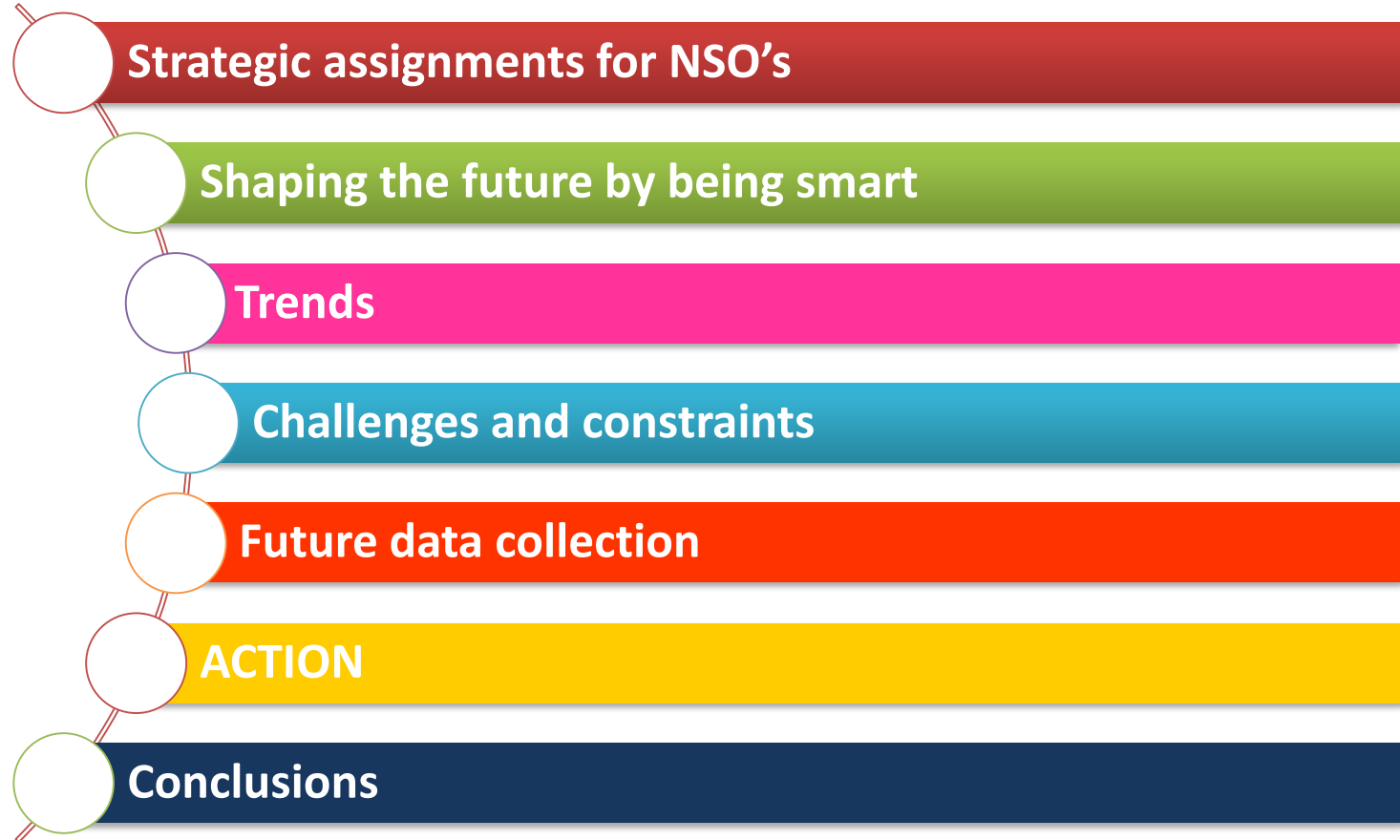
Irene Salemink

NPSO Innovatie dag 26 November 2019

"Statistical output is generated to a maximum extent using non-primary data sources. Searching for available and applicable data sources, data capture modes, and data sharing solutions is an essential part of the **collection strategy**, as is protecting confidentiality and privacy, with an appreciation for data suppliers and regard for social acceptance."

# Future data collection

- Strategic assignments for NSO's
- Shaping the future by being smart
- Trends
- Challenges and constraints
- Future data collection
- ACTION
- Conclusions

# Strategic assignments NSO's

# Strategic assignments NSO's

- Demand driven and user centric
  - Broad range statistics/information
  - Fit-for-purpose
  - Accurate and timely
  - Policymakers
    - Private sector
      - Society



- Fact based policy making
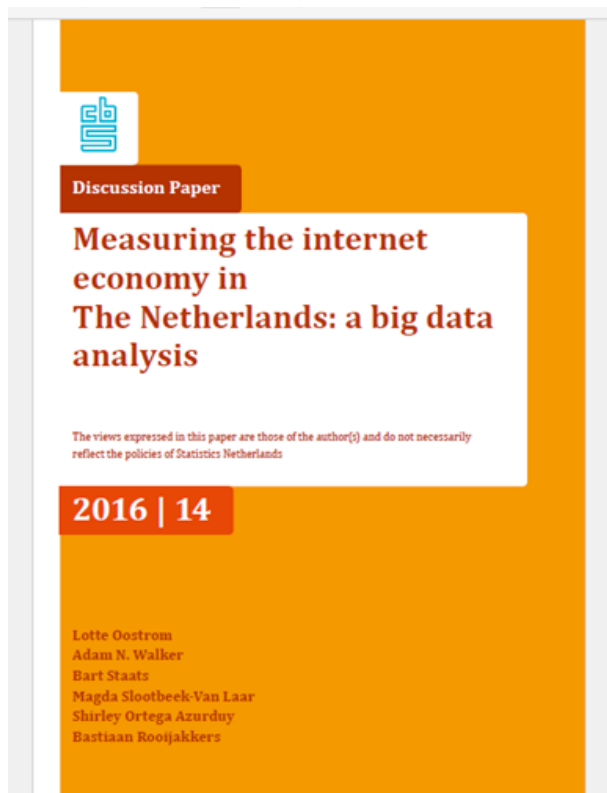  - Actionable Intelligence

# Strategic assignments NSO's

- Dealing with complexity
  - Complex societal and economic phenomena
  - Real time statistics
  - Tailored to all aggregation levels
  - More detail, regional, local
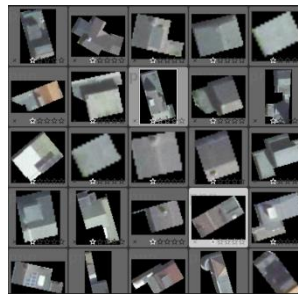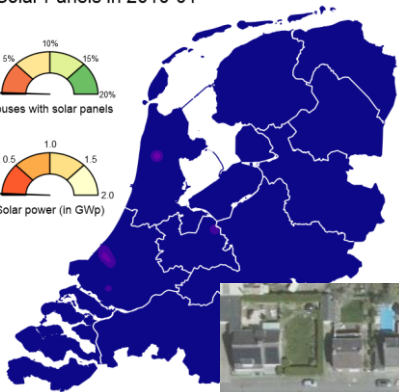  - Fit –for-purpose quality

# New information products



**Discussion Paper**

**Measuring the internet economy in The Netherlands: a big data analysis**

The views expressed in this paper are those of the author(s) and do not necessarily reflect the policies of Statistics Netherlands

**2016 | 14**

Lotte Oostrom
Adam N. Walker
Bart Staats
Magda Slootbeek-Van Laar
Shirley Ortega Azurduy
Bastiaan Rooijakkers

# Complex phenomena



## Energy transition

Support local government with insights to support sustainable development goals

Production Installation Register (solar panels)
VAT registers from Tax authority
Basic register Addresses and buildings

# Complex phenomena



## Innovative enterprises

Web scraping and text mining to identify small innovative enterprises

Classification
Linking to background characteristics

Center for Big Data Statistics

# Shaping the future by being smart

# Shaping the future by being smart

- Official statistics become "smart statistics"
    - Smart technologies
    - Smart data

- Guaranteed confidentiality

- Privacy by design

## Trusted Smart Statistics

# Trusted Smart Statistics

- Facts for evidence based policy making
- Quantitative monitoring of development and progress of policy
- Society oriented
- Reliable and innovative
- Protect confidentiality and privacy

## Nature of data collection is bound to change
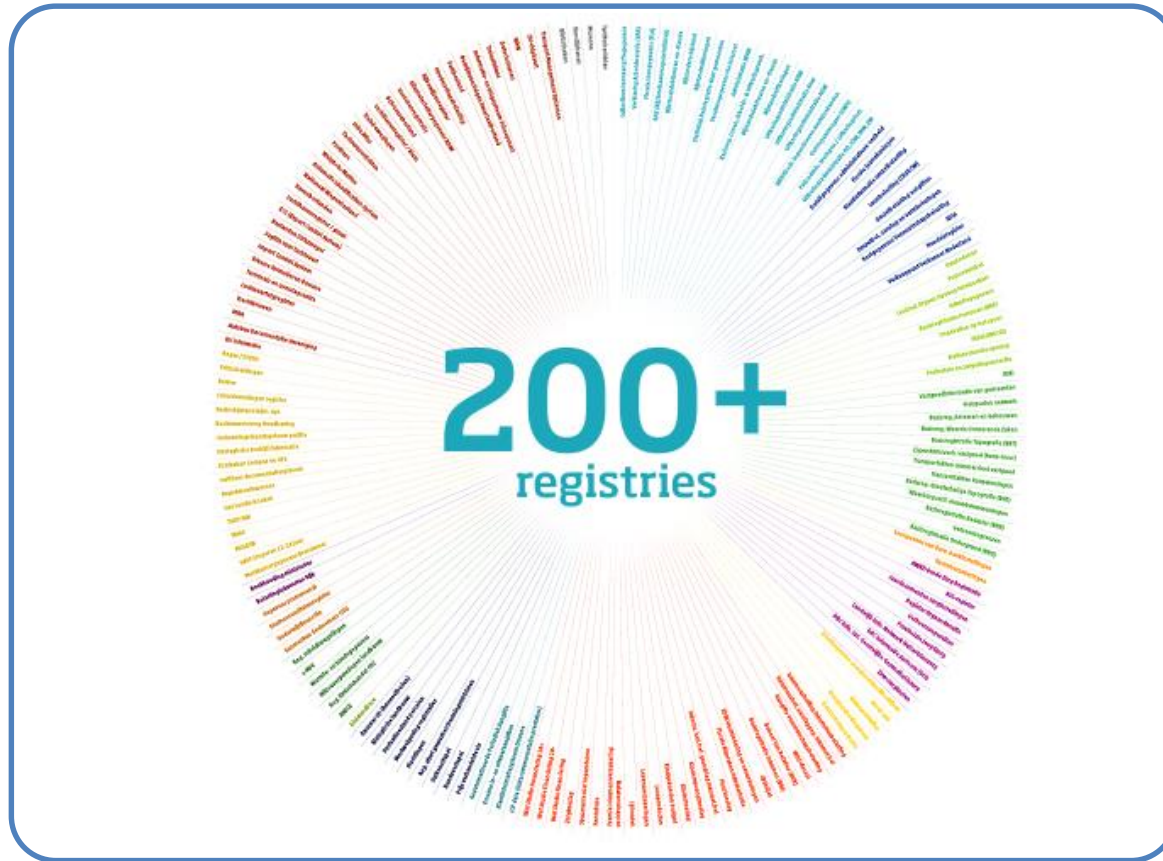
# Trends

# Trends



- Datafication



- Data storage and access



- Computing power and analytical capabilities

# Administrative data use at CBS



200+
registries

# Sensors and IoT



geek & poke

MY COFFEE
MACHINE HAS
UNFOLLOWED ME

THE INTERNET OF THINGS



INTERNET OF THINGS

We Live in Exponential Sensor Times

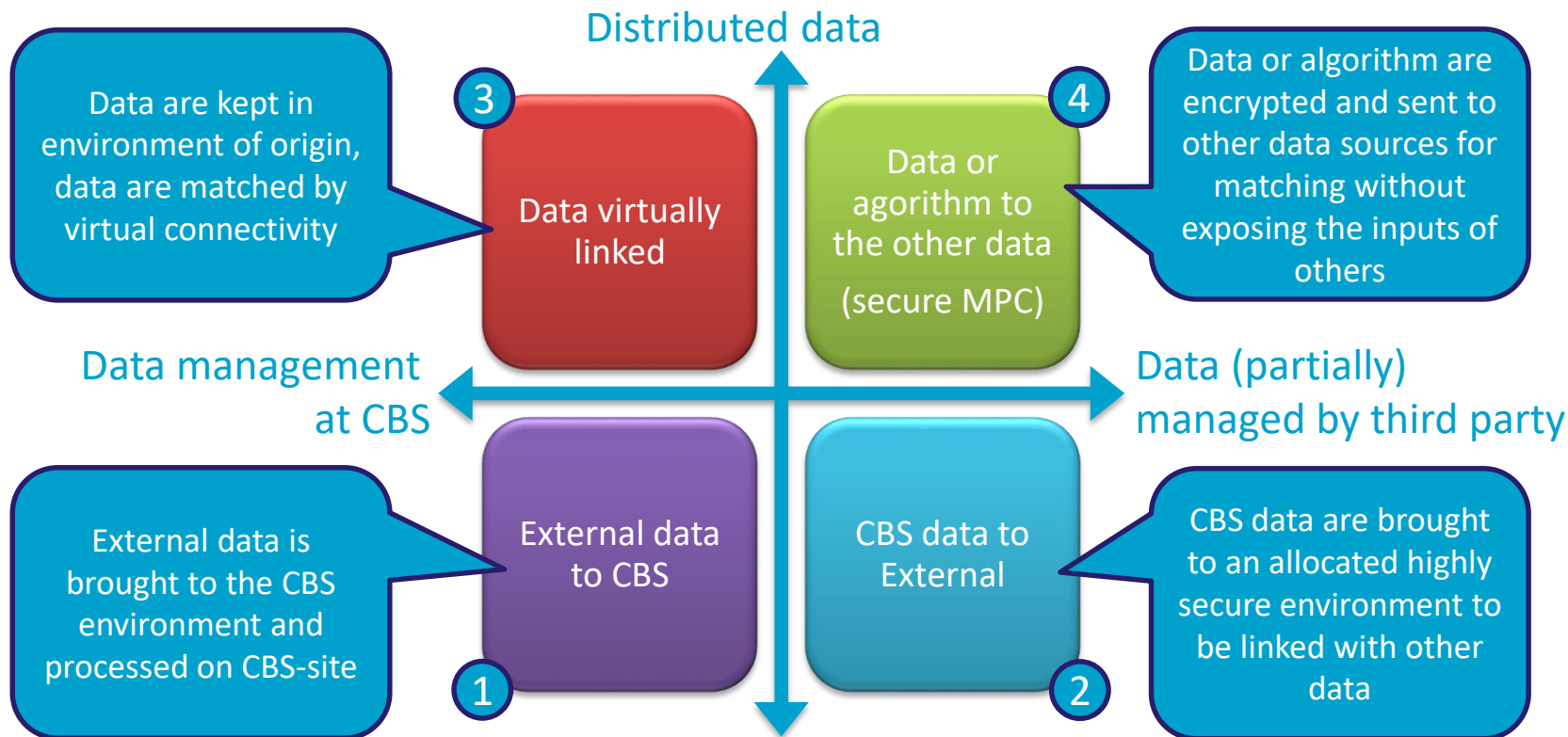| Period | Annual sensor volumes | New industry consuming sensors |
|---|---|---|
| 1960s | 10,000s | Defense, Avionics |
| 1970s | 100,000s | Process Control, Automotive |
| 1980s | 1,000,000s | Medical |
| 1990s | 10,000,000s | Consumer |
| 2000s | 100,000,000s | Mobile |
| 2010s | 10,000,000,000s | Wearables, mHealth, Internet of Things, Big Data |
| 2020s | 10,000,000,000,000s | Internet of Everything, Social Cloud |

Mobile market accelerated
the sensor growth by an
order of magnitude

# Data storage and access



- Need for :

    Fast, easy, and free access to all relevant data

- Reality:

    Datasets; too big to copy, not allowed legally to "leave the building", need matching between multiple (different) sources, require knowledge, only the proportion that is needed can be accessed…

- Solution:

    Dependent on type of data sharing

# Data architecture patterns



Distributed data

Data are kept in environment of origin, data are matched by virtual connectivity

**3** Data virtually linked

**4** Data or agorithm to the other data (secure MPC)

Data or algorithm are encrypted and sent to other data sources for matching without exposing the inputs of others

Data management at CBS

Data (partially) managed by third party

External data is brought to the CBS environment and processed on CBS-site

**1** External data to CBS

**2** CBS data to External

CBS data are brought to an allocated highly secure environment to be linked with other data
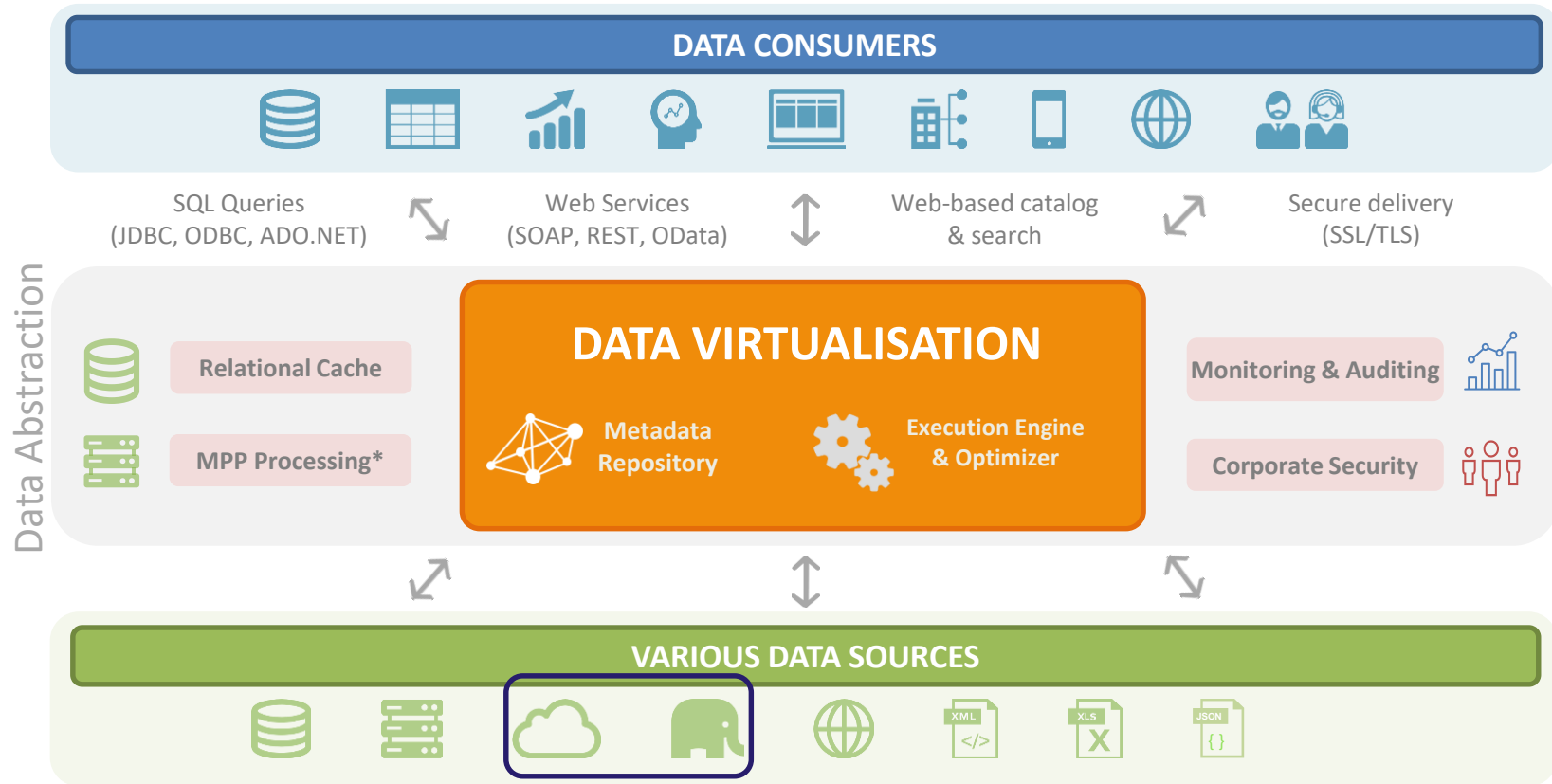
18

# Data storage and access

- These 4 patterns come with capabilities that need further investigation

  - Privacy preserving analytic techniques

    - Secure multi party computation

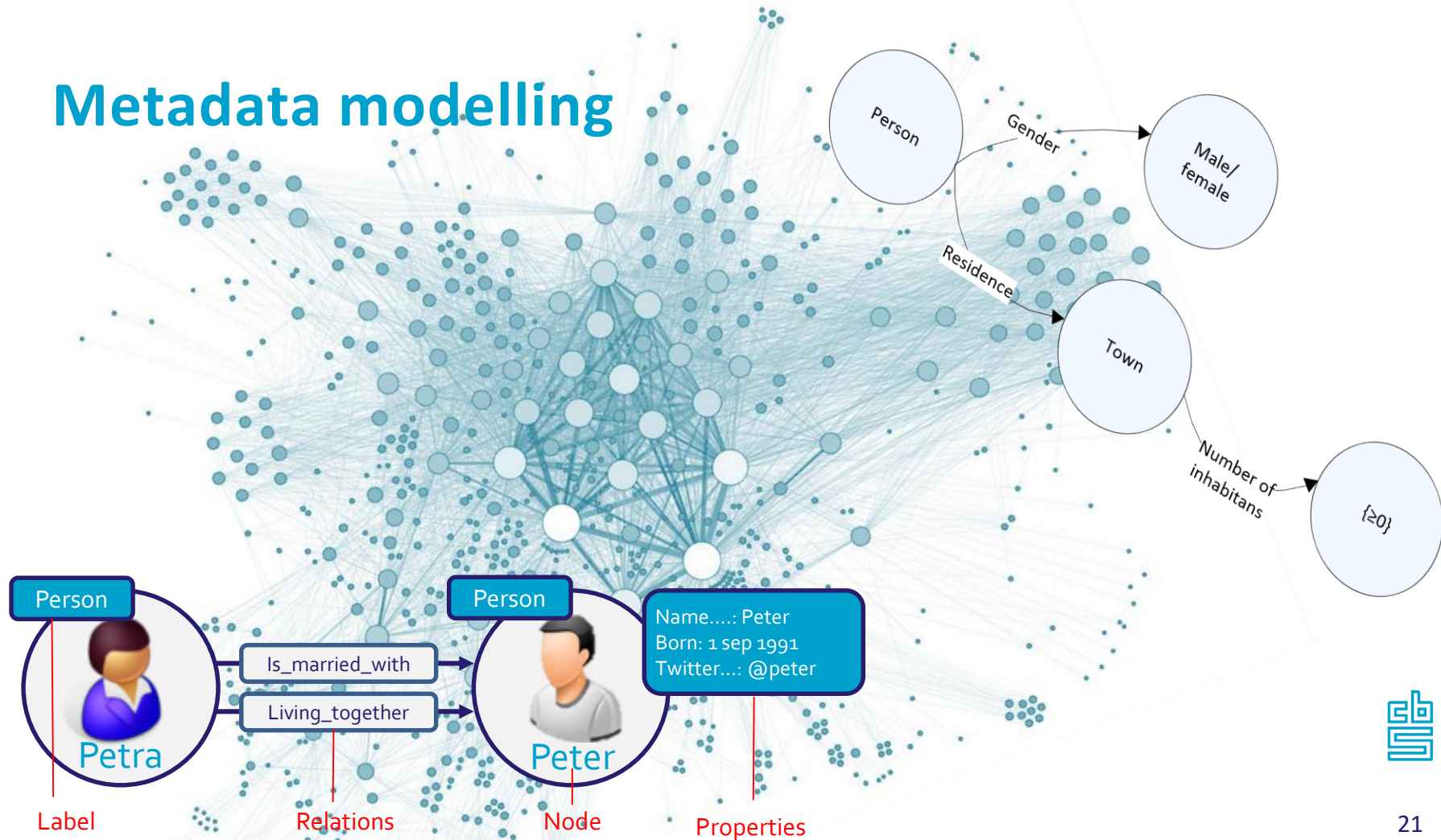  - Data virtualisation and data abstraction

  - Metadata management

# Data virtualisation architecture



**DATA CONSUMERS**

SQL Queries
(JDBC, ODBC, ADO.NET)

Web Services
(SOAP, REST, OData)

Web-based catalog
& search

Secure delivery
(SSL/TLS)

Data Abstraction

Relational Cache

MPP Processing*

**DATA VIRTUALISATION**

Metadata
Repository

Execution Engine
& Optimizer

Monitoring & Auditing

Corporate Security

**VARIOUS DATA SOURCES**

XML
</>

XLS
X

JSON
{}

* Massive Parallel Processing

# Metadata modelling



Person — Gender → Male/female

Person — Residence → Town

Town — Number of inhabitans → {≥0}

Person
Petra

Is_married_with →
Living_together →

Person
Peter

Name....: Peter
Born: 1 sep 1991
Twitter...: @peter

Label        Relations        Node        Properties

# Computing power and analytics

- Explore and Confirm
  - Deductive; data analysis to explain, check or validate ideas
  - Inductive; data analysis to generate new ideas

- Edge analytics
  - Analysis and data quality framework are brought to the data gathering devices instead of moving the data to the (centralized) analytics and quality frameworks

# Challenges and constraints

# Challenges

- Information Gap due to a Data gap

- Burden to society and response rates

- Privacy protection and difficult acces tot data  tot data
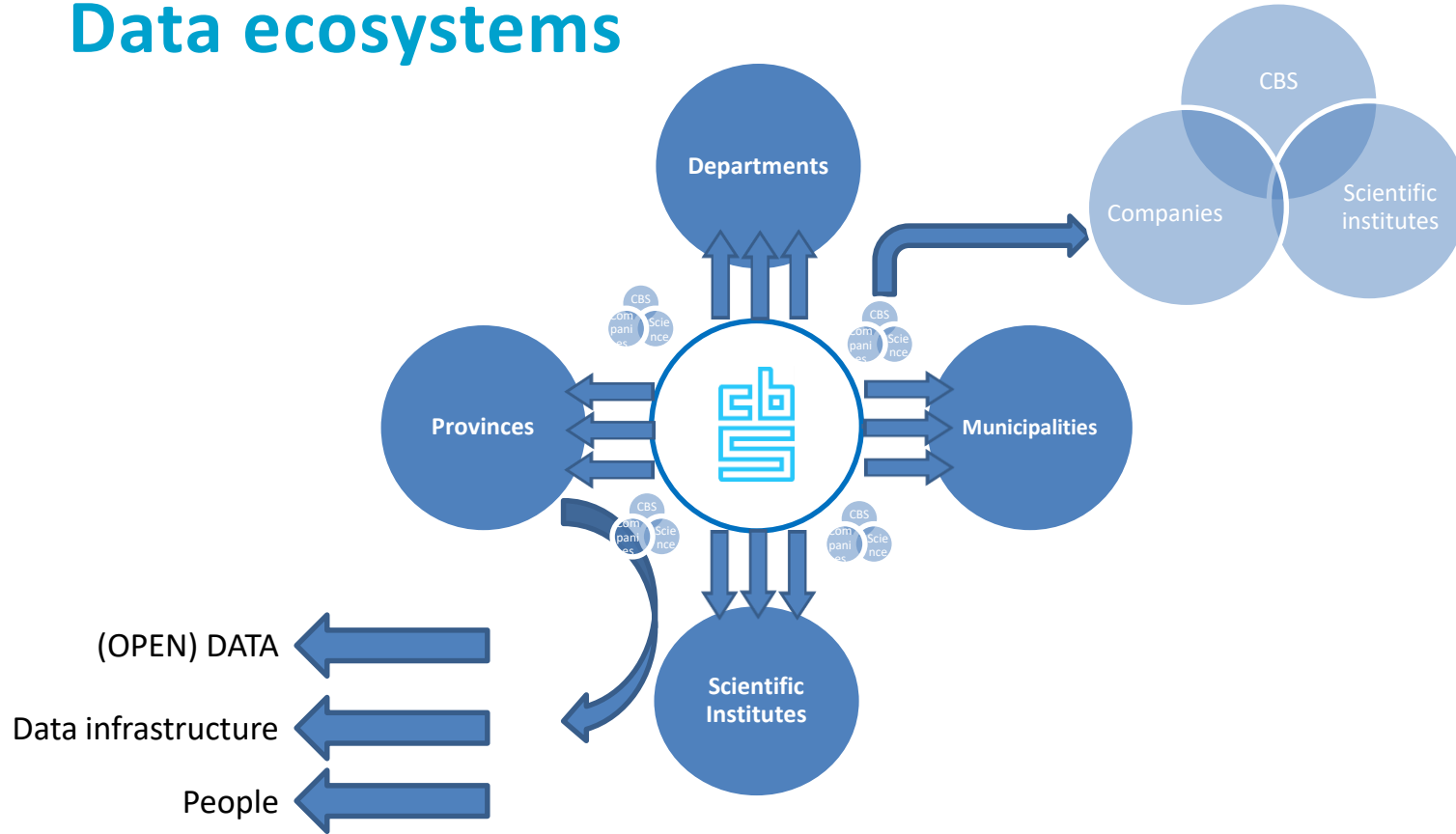
- Methodology

# Future data collection

# CBS in relation to data owners & end users

- Signalling trends and potential future needs
- Meet the public demand for information
- Data hub
- CBS fulfils role of a platform
- Enabling broad cooperation between
    - Governmental bodies
    - Municipalities
    - Companies
    - Scientific institutes

# Data ecosystems



(OPEN) DATA

Data infrastructure

People

# Tapping into and unlocking new data sources

- Datascouting

  - Awareness of data sources being there

  - Usability and availability of the data source

  - Organization and facilitation of the acquisition, tapping into and unlocking of new data sources


- Data scout en community

  - Link between internal and external stakeholders

# Surveys – primary data collection

- CAWI hybrid mode

- Multi mode and multi source

- Custom-fit data collection*

- Experimenting

  - Crowdsourcing

  - Interactive Voice Response

- Validation of alternative sources / data

# Call to ACTION

# Call to action

- Methodology
- Collect – Connect – Link
- Metadata
- Proprietary sensor networks
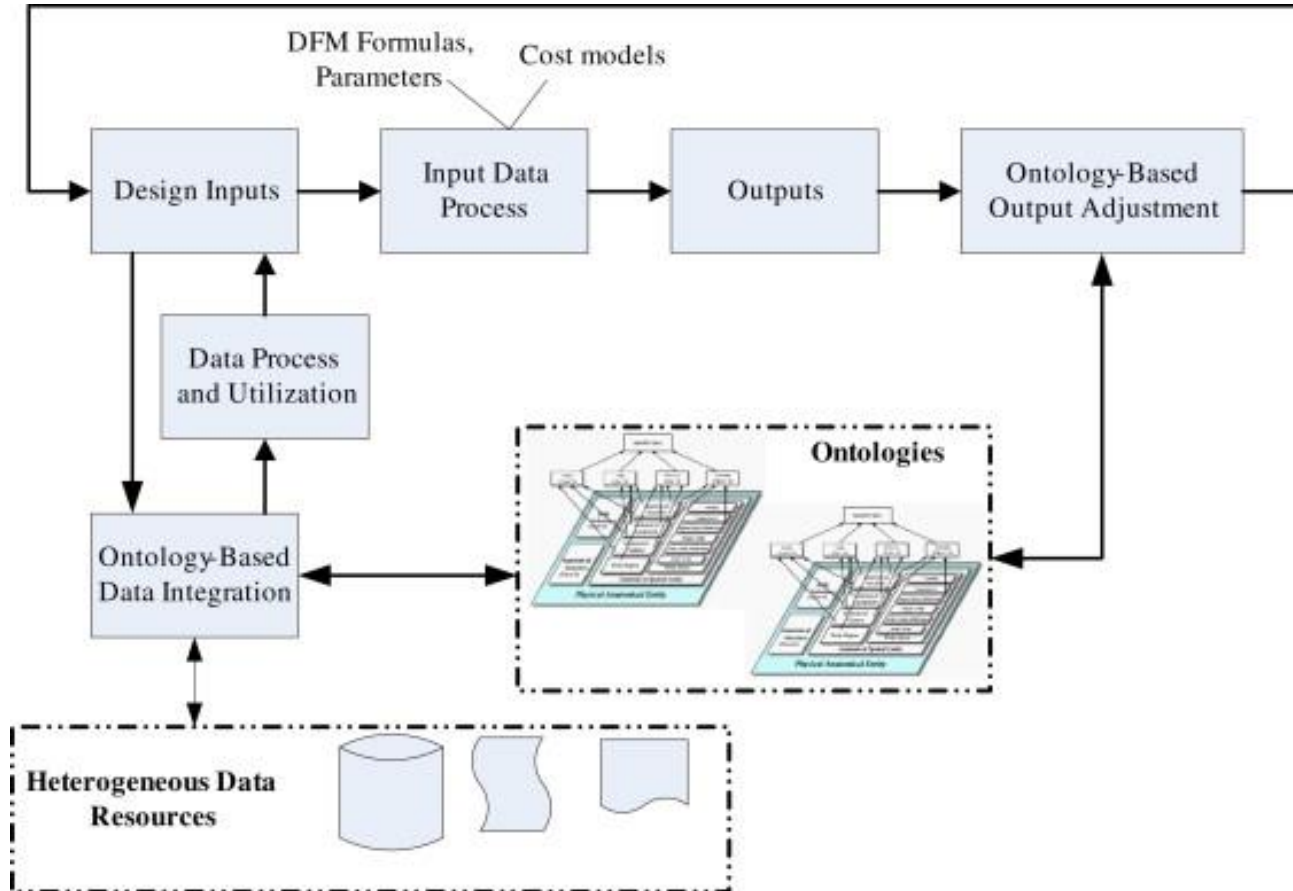- Legal frameworks and social acceptability
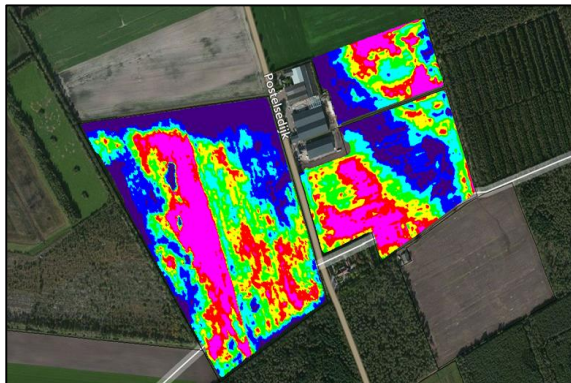
# Methodology on combining data

- Multiple collection modes and different data sources
    - units to be matched do not equal source units
    - sources do not contain overlapping units
    - matching errors
    - variables in multiple sources with different measurement errors

- Complex mixed mode designs

- Extent, combine and/or renew existing techniques
    - Probabilistic matching
    - Matching with supervised machine learning
    - Synthetic matching

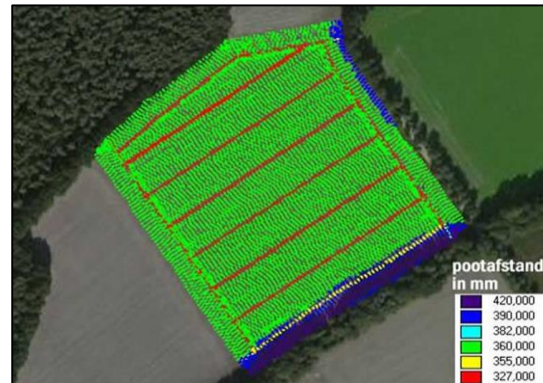- Integration by design

# Integration by design

# Sensor networks - precision agriculture cycle



## Winter
- Draw parcels
- Yield potential
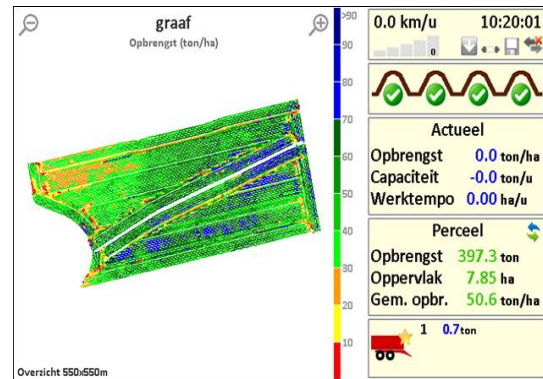- Tractor lanes

## Spring
- Fertilization
- Variable planting

## Summer
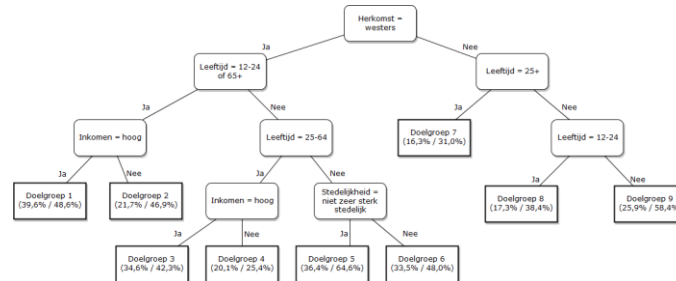- Additional fertilization, pesticides & water
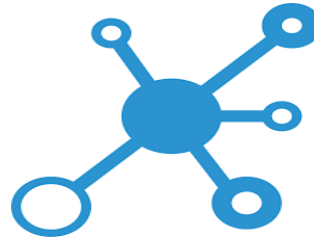- Based on sensordata

## Autumn
- Harvest
- Storage

# …and also….



Doelgroep 1 (39,6% / 48,6%)
Doelgroep 2 (21,7% / 46,9%)
Doelgroep 3 (34,6% / 42,3%)
Doelgroep 4 (20,1% / 25,4%)
Doelgroep 5 (36,4% / 64,6%)
Doelgroep 6 (33,5% / 48,0%)
Doelgroep 7 (16,3% / 31,0%)
Doelgroep 8 (17,3% / 38,4%)
Doelgroep 9 (25,9% / 58,4%)

# Conclusions

- Influence of NSI on content of captured data diminishes

- Availability and technology greatly determine the data to be collected

- New sources increase complexity

  - Validation, measuring errors & biases, mixed modes and multiple sources are a challenge

- AND possibilities

  - New – More – Cheaper –  Faster

# Conclusions

- Data collection becomes Advanced

- Data collection through data connection

- Statistical production process has to follow

- Integration by design

# Vision paper can be found at

https://www.cbs.nl/en-gb/uitgelicht/statistics-netherlands-at-isi-wsc-2019

Facts that matter