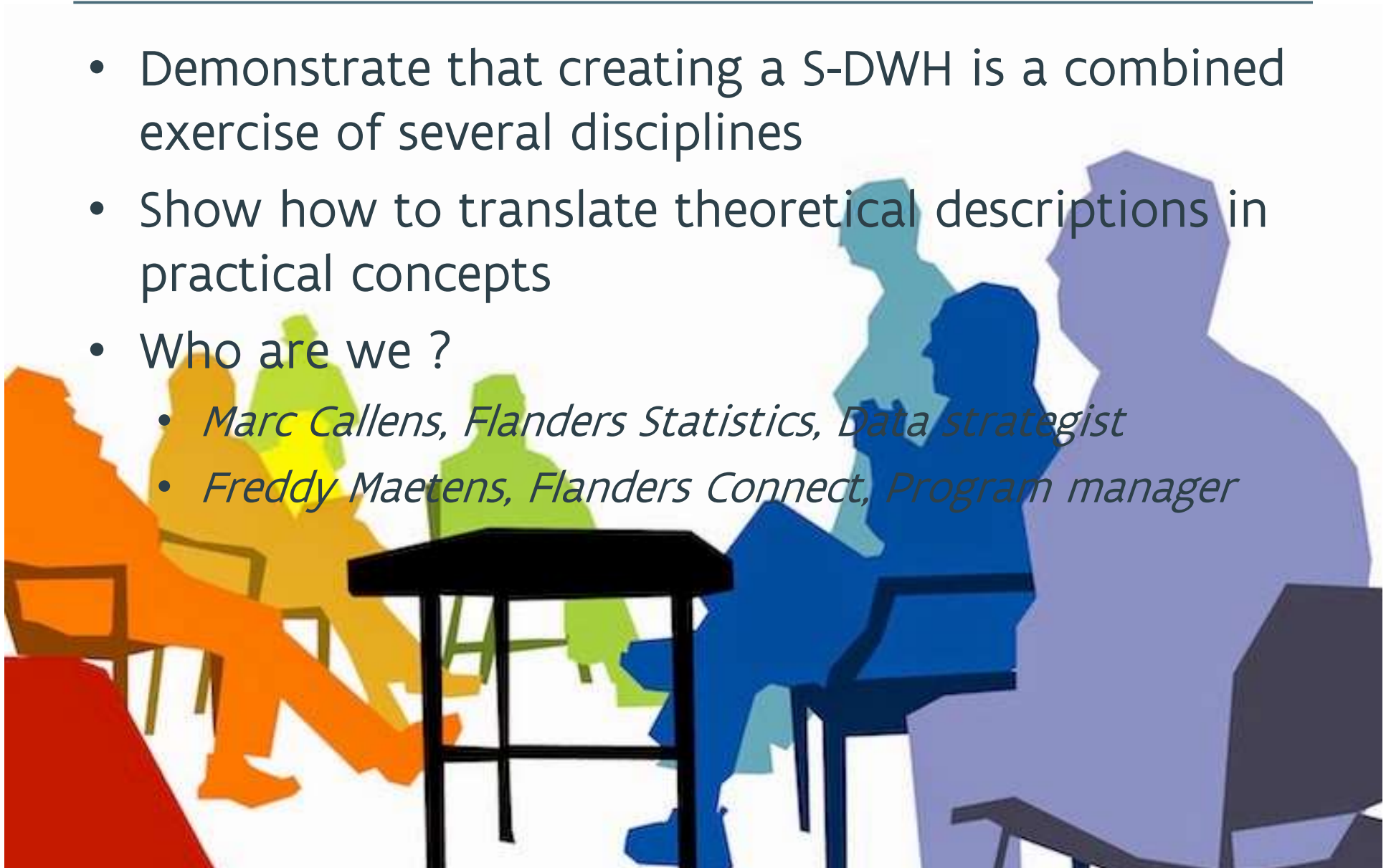


Defining
a practical approach to realize a
Statistical Datawarehouse platform
Using
GAMSO GSBPM GSIM

Objectives of this presentation

- Demonstrate that creating a S-DWH is a combined exercise of several disciplines
- Show how to translate theoretical descriptions in practical concepts
- Who are we ?
 - *Marc Callens, Flanders Statistics, Data strategist*
 - *Freddy Maetens, Flanders Connect, Program manager*



Short introduction to Statistics Flanders



Statistics
Flanders

STATISTIEK
VLAANDEREN

 Vlaamse
overheid

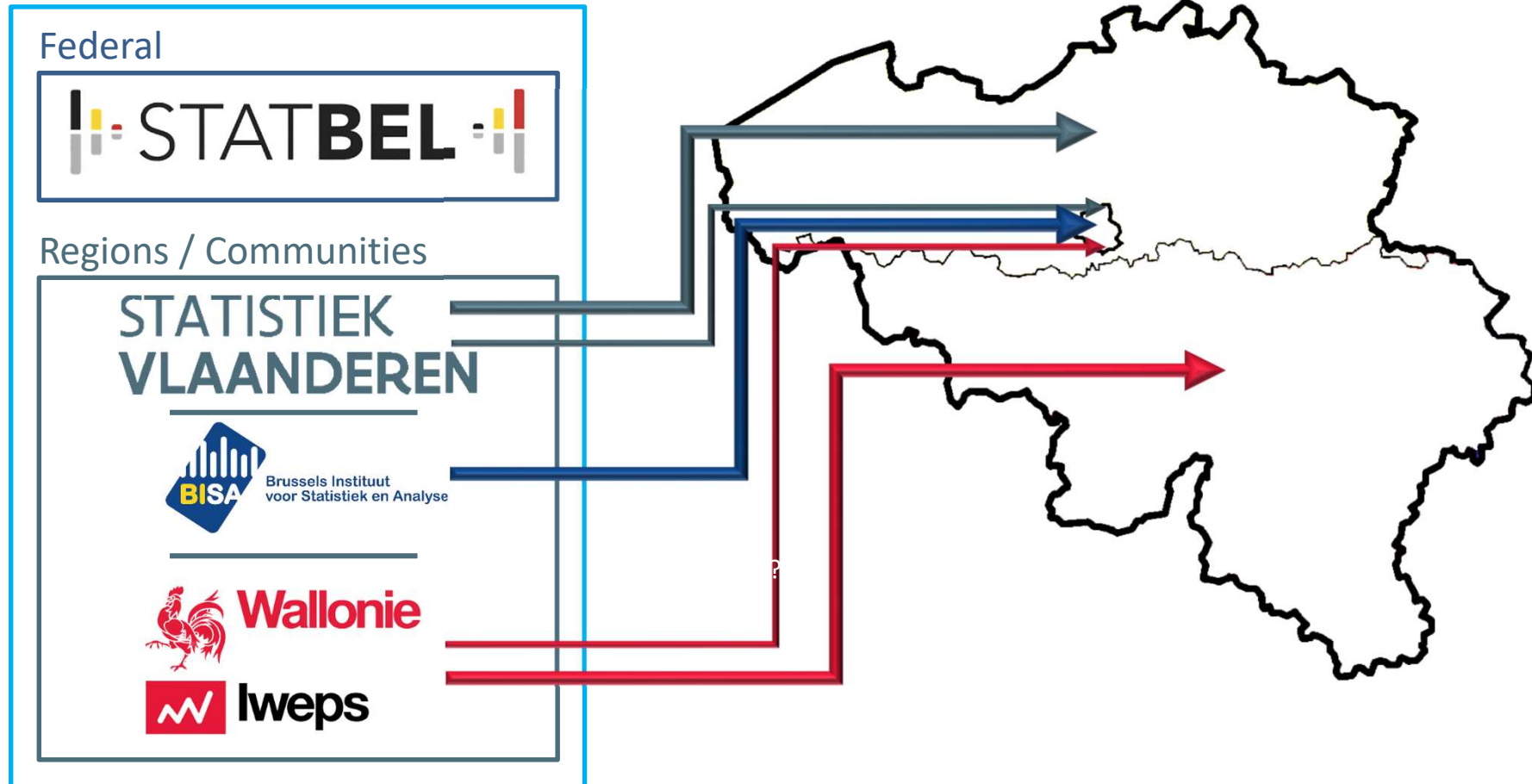
STATISTIEK
VLAANDEREN

 Vlaamse
overheid

Context

IIS

Interfederal Statistical Institute

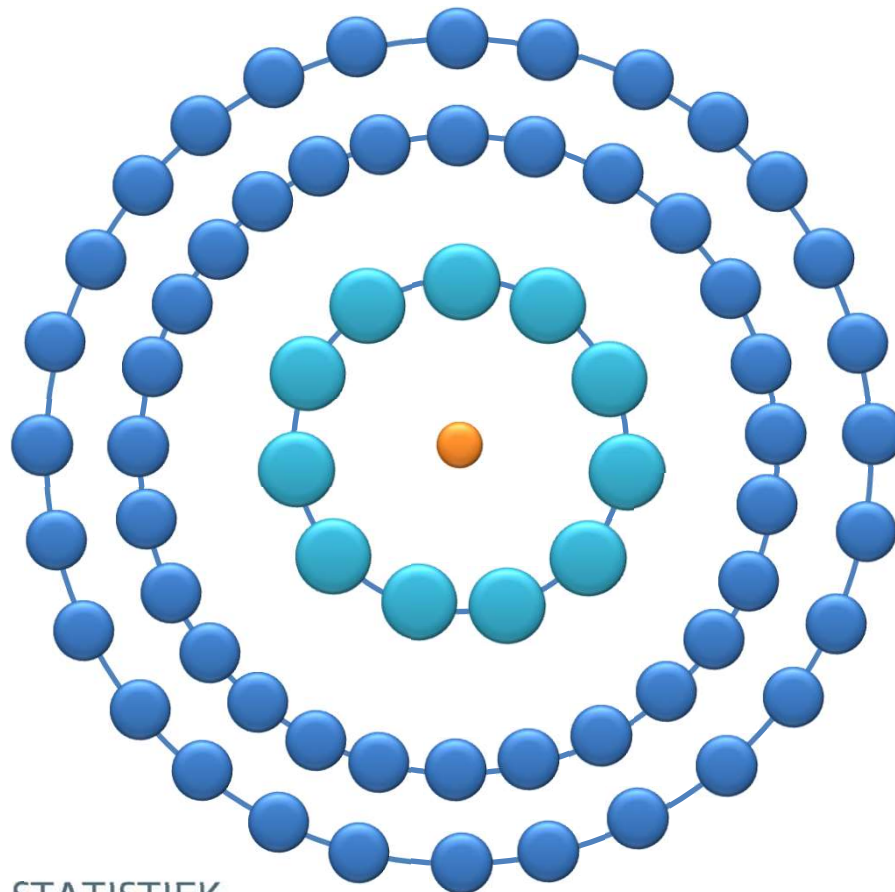


STATISTIEK
VLAANDEREN

Simplified representation of the Belgian statistical landscape.

Coordination of a decentralized network

Network Statistics Flanders



Decentralized production between:

- 11 policy areas
- 55 entities
- VGC: Flemish Community Commission
- Municipalities & provinces

Sources:

- Regional, local and federal sources and data owners
- Also sources from outside the government

Flemish statistical authority

- 3 Objectives:
 - Coordination in collaboration with
 - Council for Flemish Public Statistics (RVOS)
 - Operational forum for Flemish Public Statistics (CVOS)
 - Production:
 - Official Flemish statistics (VOS),
 - Other output
 - in depth analysis,
 - experimental statistics
 - methodology
 - Data & Methods:
 - Mainly supporting the other 2 objectives

Introduction to Statistical Data Warehouse

Statistical
Datawarehouse

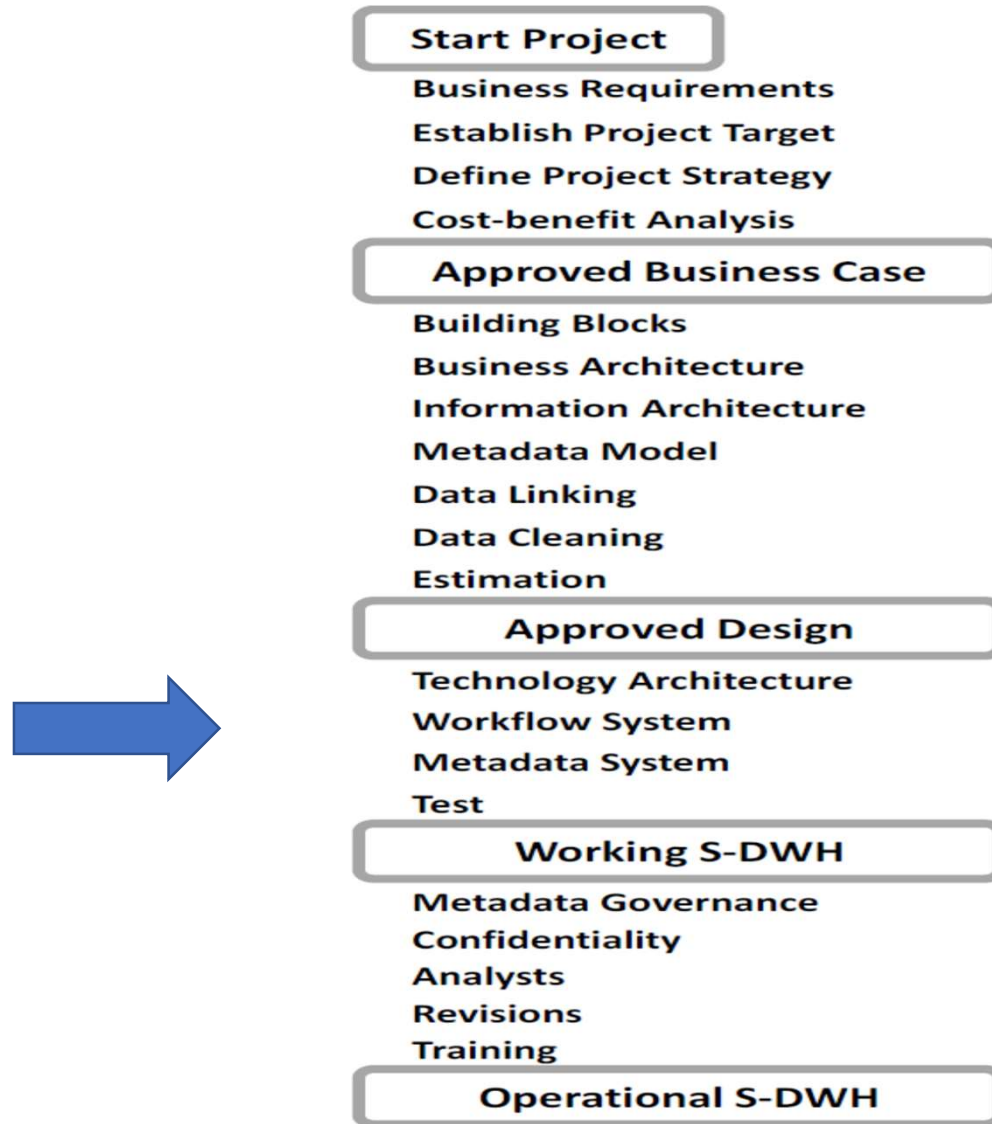
Introduction to the S-DWH Project - Objectives

- Modernisation of the statistical methodology by introducing new methods, techniques and tools
- ICT modernization of the statistical production process; collect, store, analyse and disseminate statistical data
- Increase efficiency by
 - Linking of data
 - Re-utilization of data
- Create a platform to organize data exchanges with the network and with federal statistical institute

S-DWH Organisational constraints - challenges

- As-Is situation
 - Due to historical reasons know how and approach is very dispersed
- Limited resources
- Infrastructural landscape at Flemish Government
 - IT infrastructure and access management is managed centrally => Has to be brought in line with statistical legislation and approach

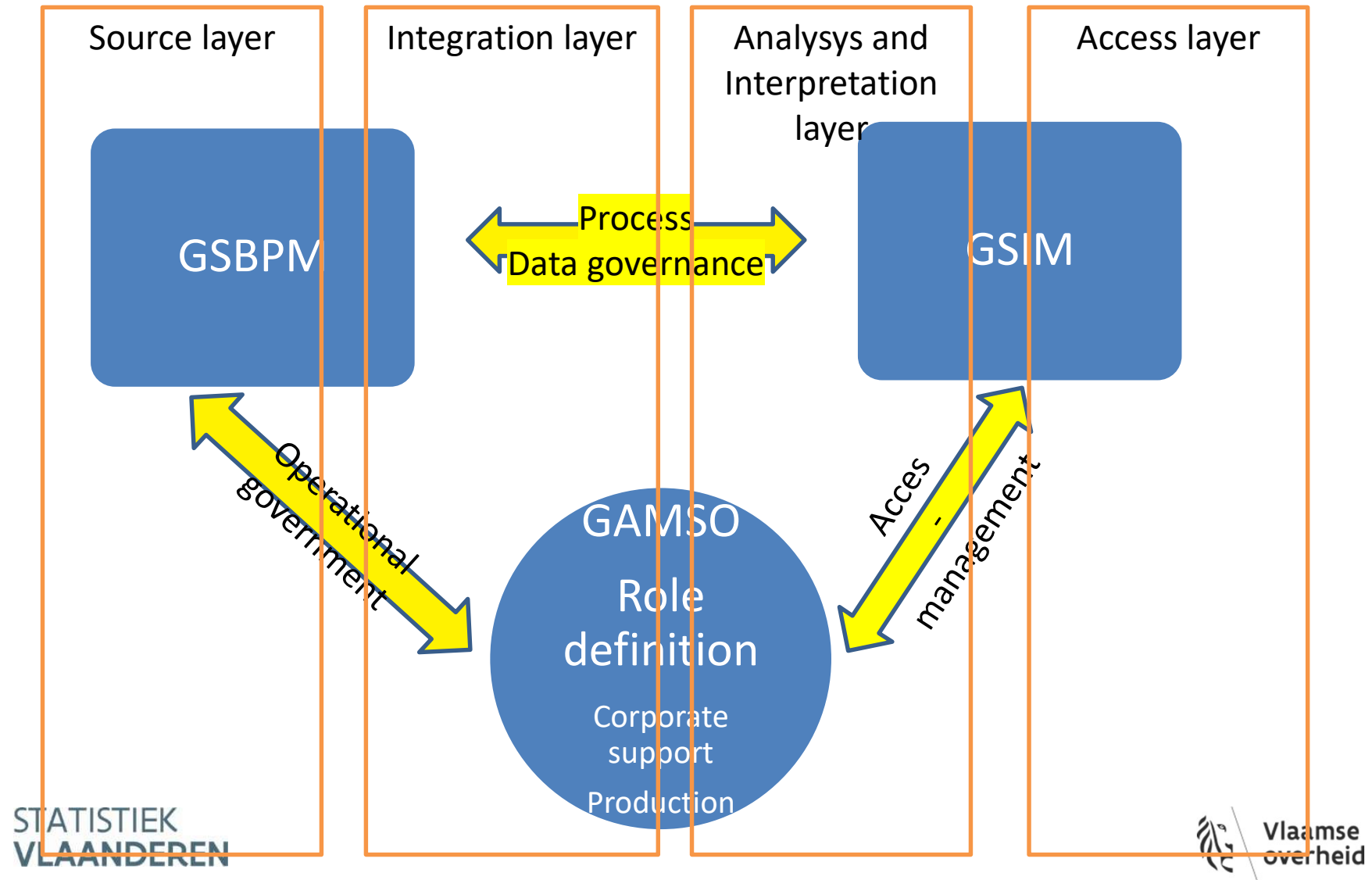
Starting point: COE Roadmap to build S-DWH



Theoretical phase

- How will we link the different methodologies to each other?

Theoretical phase : link the methodologies



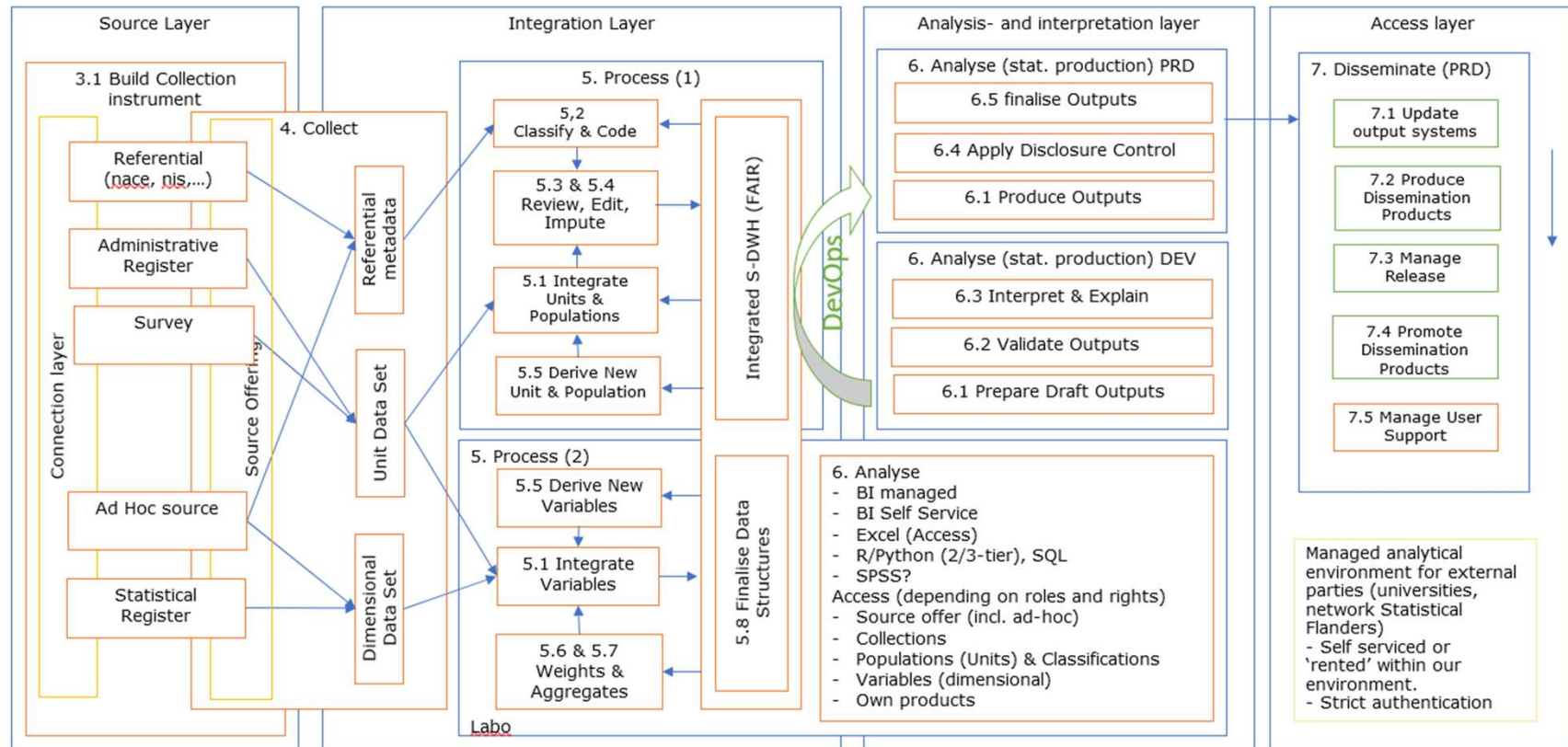
GAMSO provides a framework

- Production roles
 - Linked to GSBPM phases 4 to 7
- Corporate support
 - Linked to all GSBPM phases : methodology, broader support
- Strategy & leadership
 - Linked to a 5year strategy plan which gave us long term objectives
- Capability management
 - Linked to different work groups defining and evaluation competences, training and education

Business requirements: GSBPM (processes)

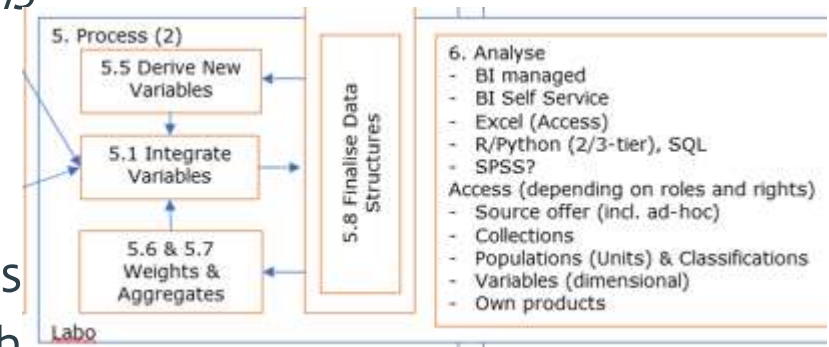
Quality Management / Metadata Management							
Specify Needs	Design	Build	Collect	Process	Analyse	Disseminate	Evaluate
1.1 Identify needs	2.1 Design outputs	3.1 Build collection instrument	4.1 Create frame & select sample	5.1 Integrate data	6.1 Prepare draft outputs	7.1 Update output systems	8.1 Gather evaluation inputs
1.2 Consult & confirm needs	2.2 Design variable descriptions	3.2 Build or enhance process components	4.2 Set up collection	5.2 Classify & code	6.2 Validate outputs	7.2 Produce dissemination products	8.2 Conduct evaluation
1.3 Establish output objectives	2.3 Design collection	3.3 Build or enhance dissemination components	4.3 Run collection	5.3 Review & validate	6.3 Interpret & explain outputs	7.3 Manage release of dissemination products	8.3 Agree an action plan
1.4 Identify concepts	2.4 Design frame & sample	3.4 Configure workflows	4.4 Finalise collection	5.4 Edit & impute	6.4 Apply disclosure control	7.4 Promote dissemination products	
1.5 Check data availability	2.5 Design processing & analysis	3.5 Test production system		5.5 Derive new variables & units	6.5 Finalise outputs	7.5 Manage user support	
1.6 Prepare business case	2.6 Design production systems & workflow	3.6 Test statistical business process		5.6 Calculate weights			
		3.7 Finalise production system		5.7 Calculate aggregates			
				5.8 Finalise data files			

GSBPM and the Layered CoE model



GSBPM and the CoE Layers

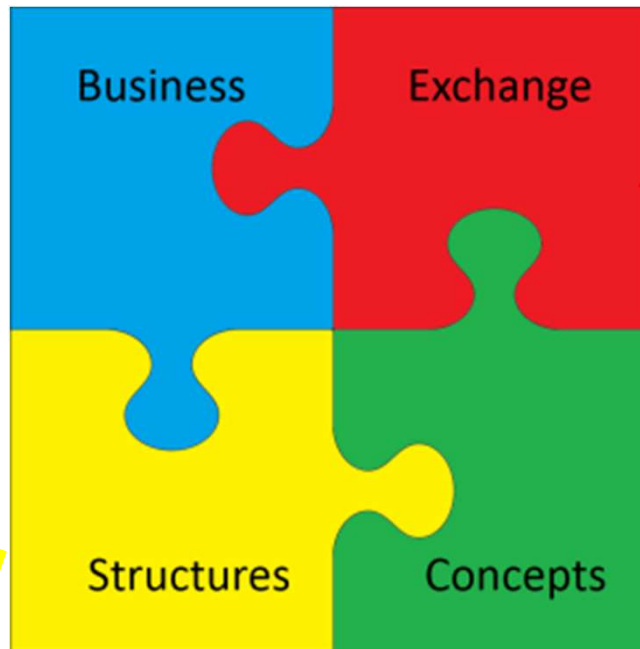
- Results:
 - 4 work phases of the GSBPM onto the 4 layers of CoE
 - 1 source connection layer for the collection instruments
 - Process phase
 - Management of *CategorySets, UnitTypes, Populations*
 - Integrating, transforming, deriving, analysing *Variables*
 - Laboratory environment
 - (partial) Process and Analysis phases
 - Emphasis on extreme flexibility with tooling (R, Python, SPSS, SAS,...)
 - Statistics production (with DevOps)



GSIM Groups

Capture designs & plans of statistical programs.
Identification of a Statistical Need.

Describe and define the terms used in relation to information and its structure.



Catalogue information that comes in and out the statistical organization.
Describe the collection and dissemination of information.

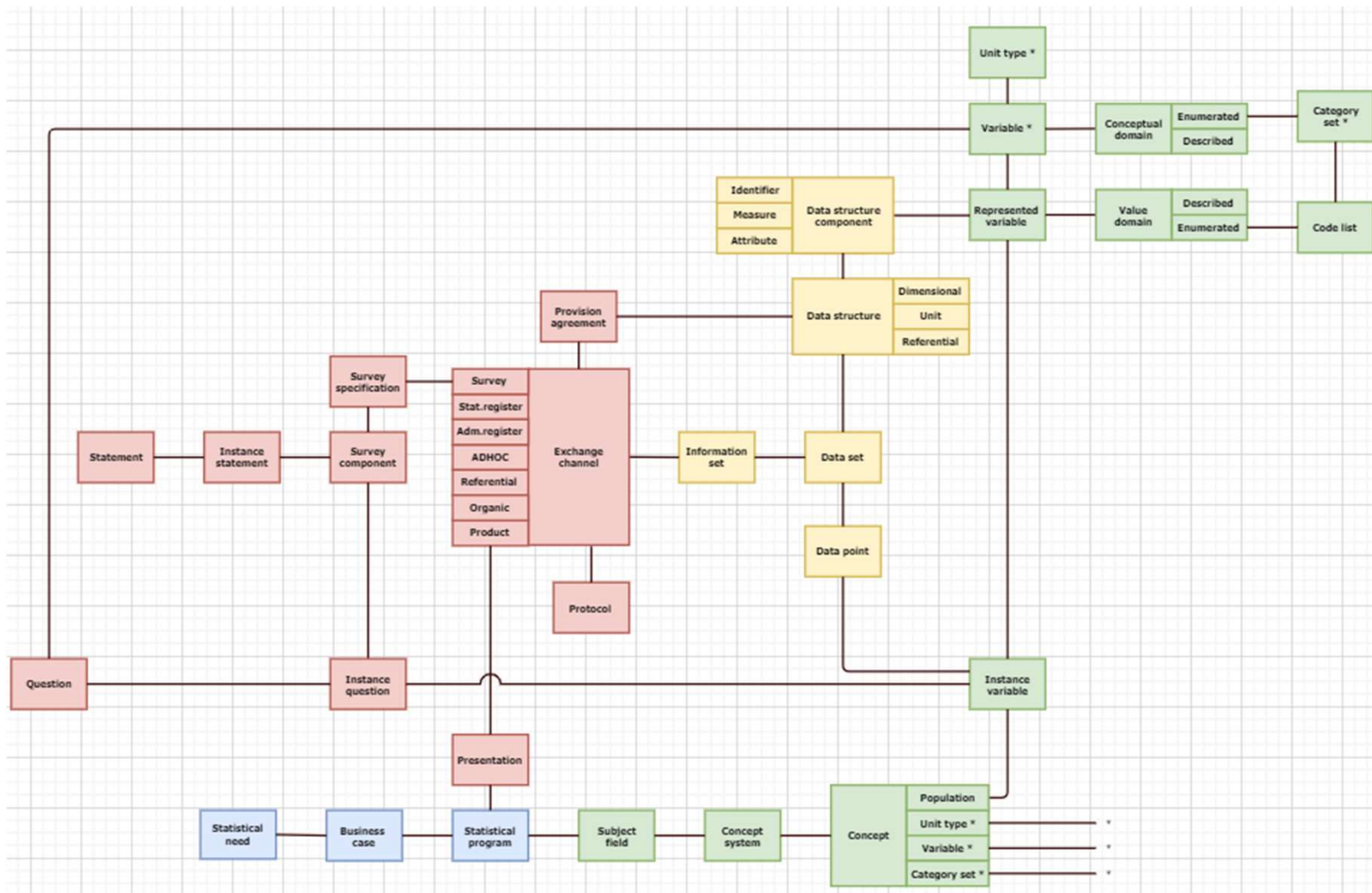
Define the meaning of data, providing an understanding of what the data are measuring.

GSIM Workshops

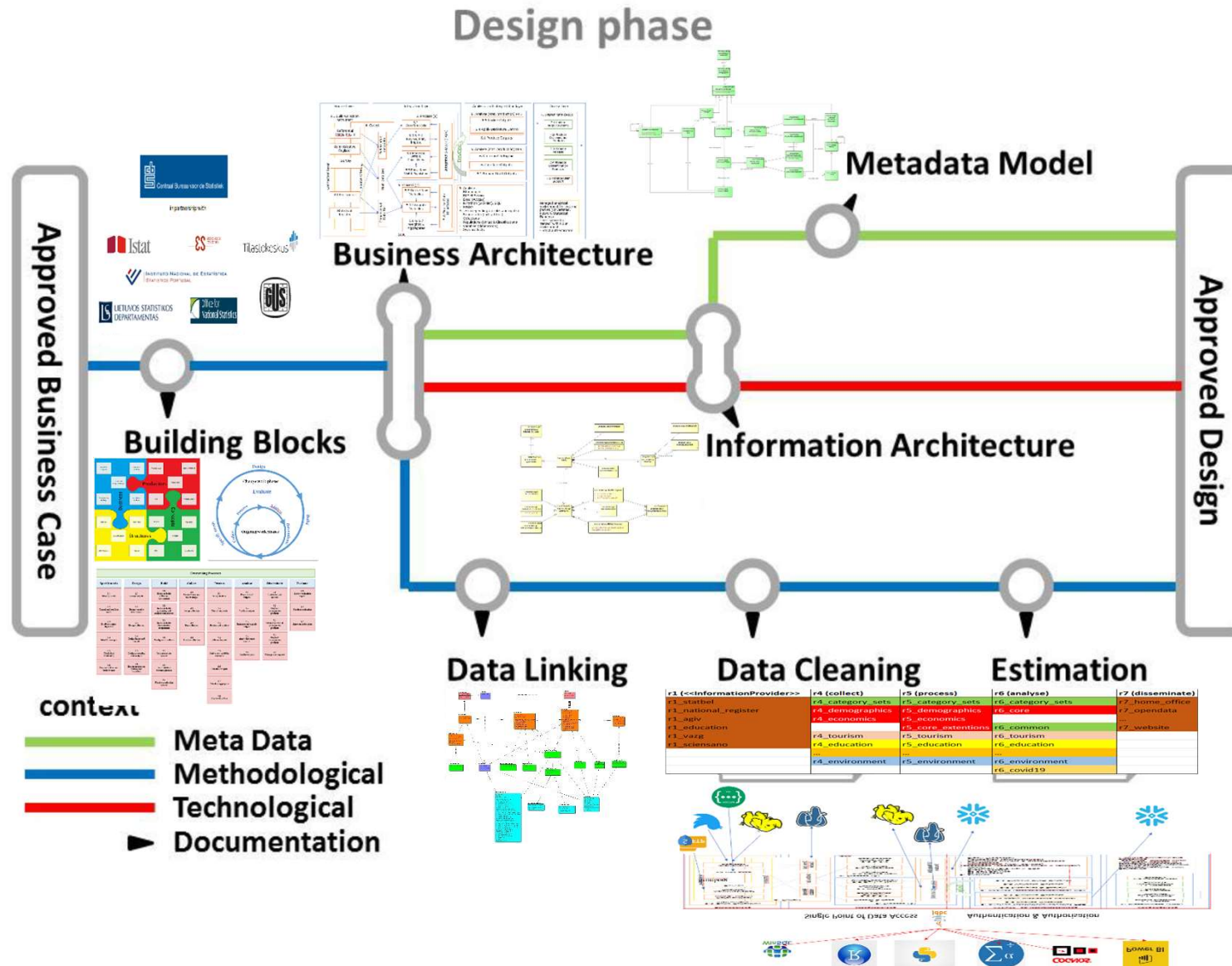
- GSIM has a very complete, detailed and large set of entities (+-120)
- Distilled to a **base model** after in-depth study (+- 40)
- Base model is discussed and reviewed by VSA
 - SME's
 - Data scientists
 - Data Strategist
 - S-DWH team
- This was done during 3 workshops



GSIM model @ VSA



COE Roadmap - Design phase components



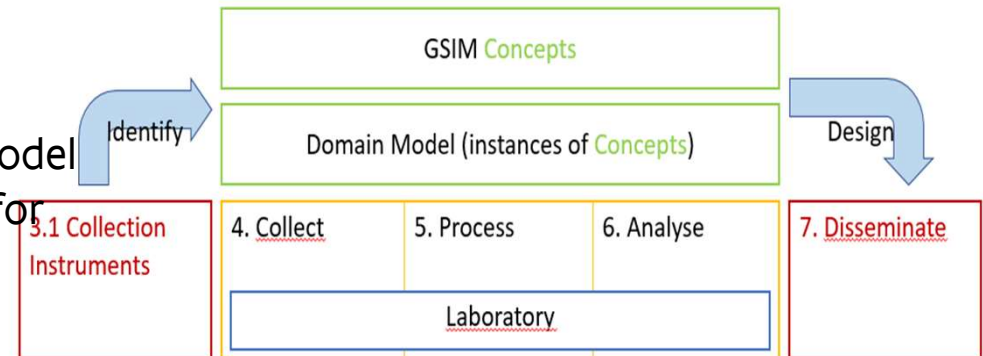
Practical phase: Matching the model with the reality

- Goal: organise data & processes around 2 core's:
Persons and Organisations
- Organize workshops based on our incoming data files
 - exchange group : Why, when, from where, for what, outcome
 - concept & structure group : Data definitions, data format throughout the process

Integrating GSIM in GSBPM

- **Concepts group:**

- GSIM level (M2)
 - Building blocks for Domain model
 - Search and navigation space for metadata.
- Domain Model (M1)
 - Instances of **Concepts**
 - Navigation and access space for data.



- **Structures group:**

- Relates DataStructures to Represented Variables
- Currently 3 types of structures (dimensional, unit, referential) coming from 4 sources: micro data; aggregated data; survey data; big data
- Allows for grouping into InformationResources
 - Covering internal work phases (collect, process, analyse) and metdata (M2)
 - Further refinement according to SubjectField & ConceptGroup
 - Centralized authorization.
- Independent of actual physical storage.
 - Questionable Data Structures
 - Service Level Agreement

- **Exchange group**

- Manages input and output

Framework Templates – Microdata

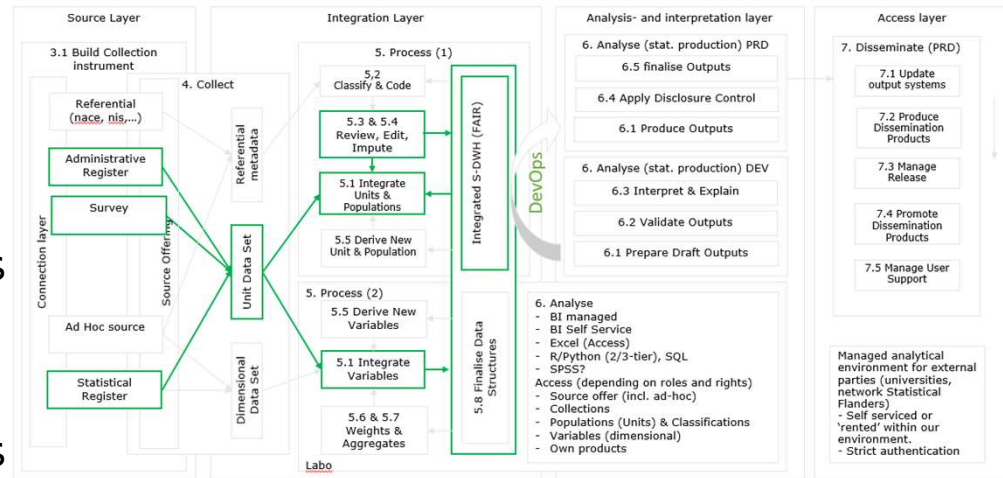
- Units (microdata)

- Collect phase

- Identifying Variables
 - Identify key spaces
 - Populations
 - Variables
 - Identify key spaces
 - Time dependence

- Process phase

- Units (with mapping/matching structures where required/possible)
 - Populations
 - Basic time independent variables
 - Template structures for time dependent variables



				Information Provider	Statbel			
				Exchange Channel	Jaarlijkse uitwisseling van demografische gegevens			
				Information Resource	Demografische datasets			
				Information Set	statbel_bevolking			
Unit Type	Variable	Variabele Omschrijving	Conceptual Domain	Tijdsaspecten? t0 geen historiek mogelijk ts Statistische waarnemingen te Event-stijl historiek	Population	Represented Variable	Instance Variable (voorlopig kolomnaam, suggesties welkom)	Value Domain
Persoon								
	vu_persoon	Identificatie van Persoon	Personen	t0	Bevolking	vu_persoon_rnrn_pseudo	ID_REC_C	Gepseudonymiseerd Rijksregisternummer
	vu_moeder	Identificatie van de Moeder	Personen	t0				
	vu_vader	Identificatie van de Vader	Personen	t0				
	vd_geboorte	Geboortedatum	Dag	t0	Bevolking	vd_geboorte	concat(geb_maand, geb_jaar, ge	Datum (yyyyMMdd')
	ve_geslacht	Geslacht	Geslacht	t0	Bevolking	ve_geslacht_statbel	geslacht_code	Geslacht-> sleutelruimte 1(statbel)
	ve_geboorteplaats	Geboorteplaats	Gemeentes - Subdomein van Bestuurlijke Indeling	t0				
	ve_geboorteland	Geboorteland ("moederland")	Land	t0				
	vd_overlijden	Sterfdatum	Dag	t0				
	ve_gemeente_overlijden	Gemeente van overlijden	Gemeentes - Subdomein van Bestuurlijke Indeling	t0				
	ve_hoofdverblijfplaats_overlijden	Hoofdverblijfplaats bij overlijden.	Gemeentes - Subdomein van Bestuurlijke Indeling	t0				
	vu_partner_overlijden	Identificatie van Partner bij overlijden	Personen	t0				
	ve_eerste_nationaliteit	Eerste Nationaliteit	Nationaliteit (of Land?)	t0	Bevolking	ve_eerste_nationaliteit_num3	fst_nat_code	Land-> sleutelruimte 1(num3)
	vd_eerste_nationaliteit	Datum eerste nationaliteit	Dag	t0	Bevolking	vd_eerste_nationaliteit	concat(fst_nat_jaar, fst_nat_maar	Datum (yyyyMMdd')
	ve_nationaliteit_overlijden	Nationaliteit	Nationaliteit (of Land?)	t0				
	ve_burgstat_overlijden	Burgerlijke Staat bij overlijden	Burgerlijke Staat	t0				
	vd_waarneming	Datum waarneming voor tijdsafhankelijke variabelen	Dag	ts	Bevolking	vd_waarneming	<jaar>'-'12'31'	Datum yyyy'-12-31
	vu_huishouden	Huishouden	Huishoudens	ts	Bevolking	vu_huishouden	id_huish	Huishouden id's bepaald door statbel
	ve_nationaliteit	Nationaliteit	Nationaliteit (of Land?)	ts	Bevolking	ve_nationaliteit_num3	nat_code nat_code_vj	Land-> sleutelruimte 1(num3)
	ve_hoofdverblijfplaats	Hoofdverblijfplaats	Gemeentes - Subdomein van Bestuurlijke Indeling - Subdomein van Adressen Postcodes	ts	Bevolking	ve_hoofdverblijfplaats_nis4 ve_hoofdverblijfplaats_nis4 ve_hoofdverblijfplaats_postcode	gem_code gem_code_vj post_code	Bestuurlijke Indeling-> sleutelruimte 1(NIS)- > level 4 (gemeentes) Postcodes
	ve_verwantschap	Verwantschap ten opzichte van referentiepersoon van het Huishouden	Verwantschappen	ts	Bevolking	ve_verwantschap_statbel)	verwant_code	Verwantschap-> sleutelruimte 1 (verwantschapcodes volgens Statbel)
	vu_hoofdverblijfplaats	Hoofdverblijfplaats	Gemeentes - Subdomein van Bestuurlijke Indeling	ts	Bevolking	vu_hoofdverblijfplaats_statbel	vu_pos_1000	Hoofdverblijfplaatscodes volgens Statbel

Framework Templates – Dimensional data

- Dimensional Data

- Collect phase

- Identify Variables Identifiable

- Primary Key (if exist)

- Identify Variables

- Time dependence

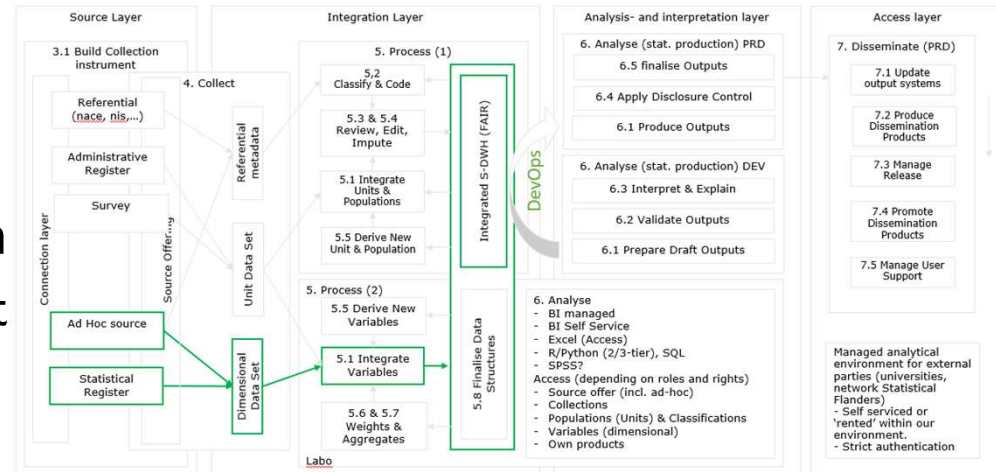
- Aggregation candidate

- Process phase (1)

- Match (and expand) CategorySets & Codelists (see referential metadata)

- Process phase (2)

- Link (CategorySets) and Integrate (Metadata)



MVP : S-DWH Approach

- Build, Buy, Meet in the middle?
 - Build from scratch?
 - Complex and elaborate GSIM & GSBPM.
 - Limited in-house resources (duration)
 - Waterfall effect (esp. with required flexibility)
 - Buy?
 - Comm. Products not aligned with GSIM/GSBPM (MOF, CWM)
 - Limited flexibility
 - Steep learning curve
 - Meet in the middle?
 - Technologies & Techniques from classical DWH? (useable?)
 - What functional components do we actually need?

Wrap up

- It is good to combine the different methodologies from the beginning
 - Stimulate the interaction
 - Evaluate the limits of each methodology
- Use the methodology as a model
 - Don't be afraid of making a simplified operational model of your ideal theoretical model

You can reach us at

- Marc.Callens@Vlaanderen.be
- Freddy.Maetens@Vlaanderen.be