eosc | FAIR-EASE

# The Mythical Data Lake
## *From Data Lake to Data Spaces*

Sept 27, 2023 @OSFAIR
Katrina Exter (VLIZ), Marc Portier (VLIZ)

The services FAIR-EASE will build

- DDAS: Data Discovery and Access Service. Users and clients will be able to
  - Find and access their own data and associated/related data
  - Find and access subsets within datasets
- VRE: Virtual Research Environment. Users will be able to
  - Call up software and services to work on data
  - Save/export the outputs of those services
  - Build up workflows/notebooks

So the FE environment needs to

- know where the data are (DDAS)
- know how to interrogate the data so they can be searched, accessed, subsetted, quantified, … (DDAS, VRE)
- know how to move those data into the services (VRE)
- know how to incorporate the data output from the services back into the DDAS and VRE

The services FAIR-EASE will build

- DDAS: Data Discovery and Access Service. Users and clients will be able to
  - Find and access their own data and associated/related data
  - Find and access subsets within datasets
- VRE: Virtual Research Environment. Users will be able to
  - Call up software and services to work on data
  - Save/export the outputs of those services
  - Build up workflows/notebooks

So the FE environment needs to

- know where the data are (DDAS)
- know how to interrogate the data so they can be searched, accessed, subsetted, quantified, … (DDAS, VRE)
- know how to move those data into the services (VRE)
- know how to incorporate the data output from the services back into the DDAS and VRE

The services FAIR-EASE will build

- DDAS: Data Discovery and Access Service. Users and clients will be able to
  - Find and access their own data and associated/related data
  - Find and access s
- VRE: Virtual Research

**DATA LAKE**

  - Call up software a
  - Save/export the o
  - Build up workflows/notebooks

So the FE environment needs to

- know where the data are (DDAS)
- know how to interrogate the data so they can be searched, accessed, subsetted, quantified, … (DDAS, VRE)
- know how to move data into the services (VRE)
- know how to incorporate the data output from the services back into the DDAS and VRE

*An undefined "somewhere" that "somehow" would allow the DDAS and VRE to find, search, access, subset, and otherwise do whatever we want with FAIR-EASE data*
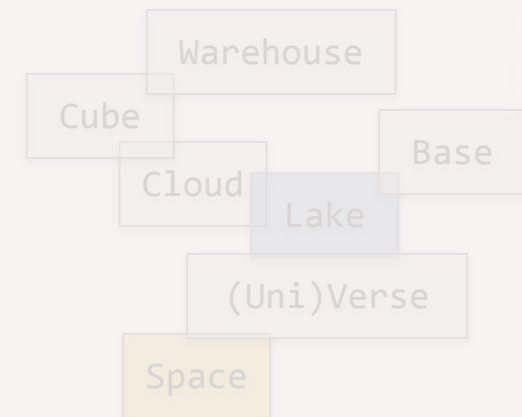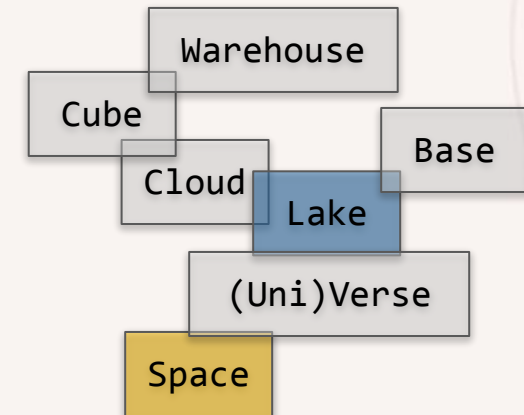
An amorphous blob with datasets swimming around in it?
A single, huge dataset that all individual datasets have "bled into" ?
A place where datasets bob about waiting to be fished out?

Data-

Soup

Warehouse

Cube

Base

Cloud

Lake

(Uni)Verse

Space

*An undefined "somewhere" that "somehow" would allow the DDAS and VRE to find, search, access, subset, and otherwise do whatever we want with FAIR-EASE data*

An amorphous blob with datasets swimming around in it?
A single, huge dataset that all individual datasets have "bled into" ?
A place where datasets bob about waiting to be fished out?

Data-  Soup

Warehouse
Cube
Cloud    Base
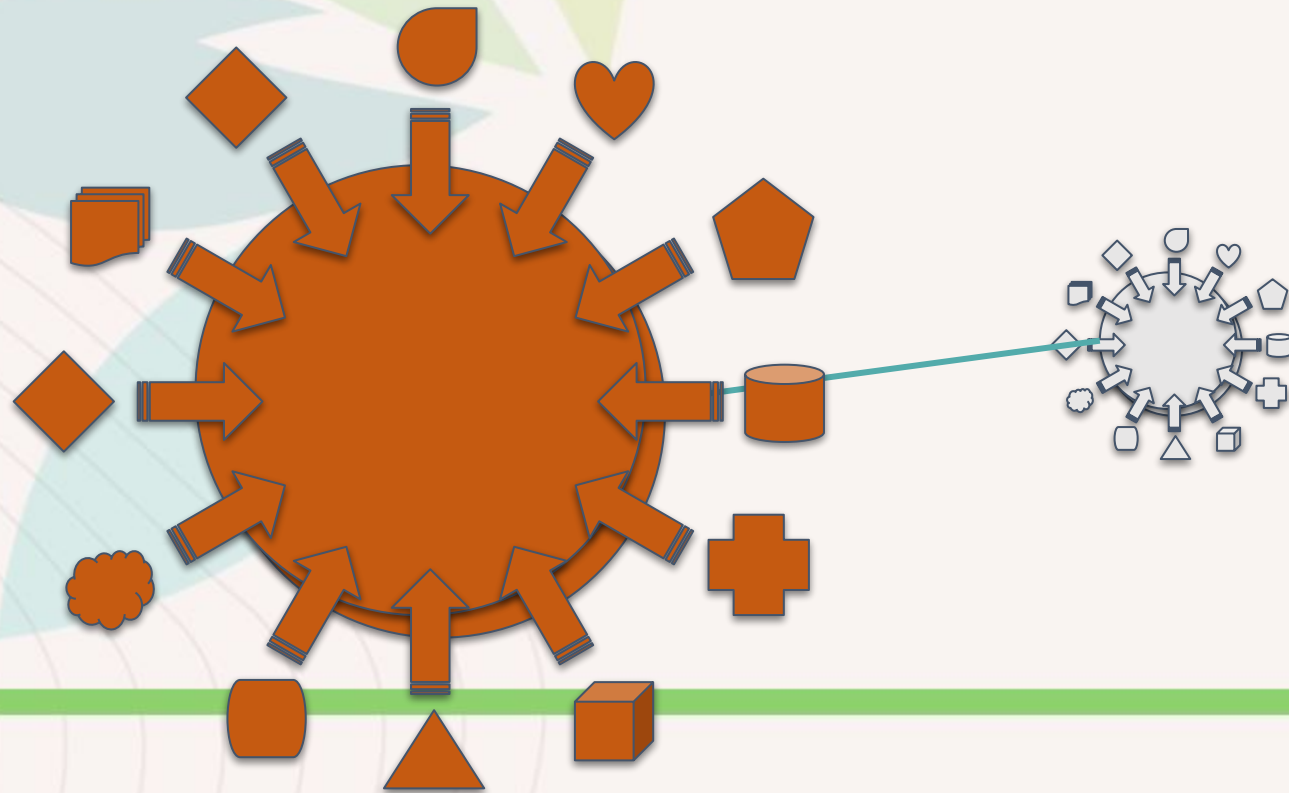Lake
(Uni)Verse
Space

Requirements

1. To allow all FAIR-EASE data to be findable and accessible
2. To allow sharing of data with the community
3. To allow for search and access to subsets and recombinations of data
4. To provide download management and manage user authentication
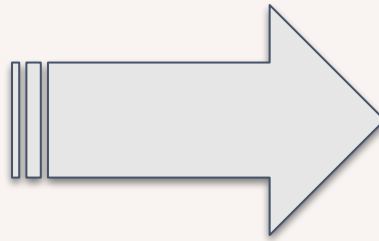5. To be flexible to evolving data sets and requirements

Restrictions

1. The diversity of data (types, sizes, accessibility, formats) is enormous and evolving
2. We need to be sustainable: whatever is developed, must continue to exist after FAIR-EASE ends
3. We want a solution that can be adopted by others for their own projects
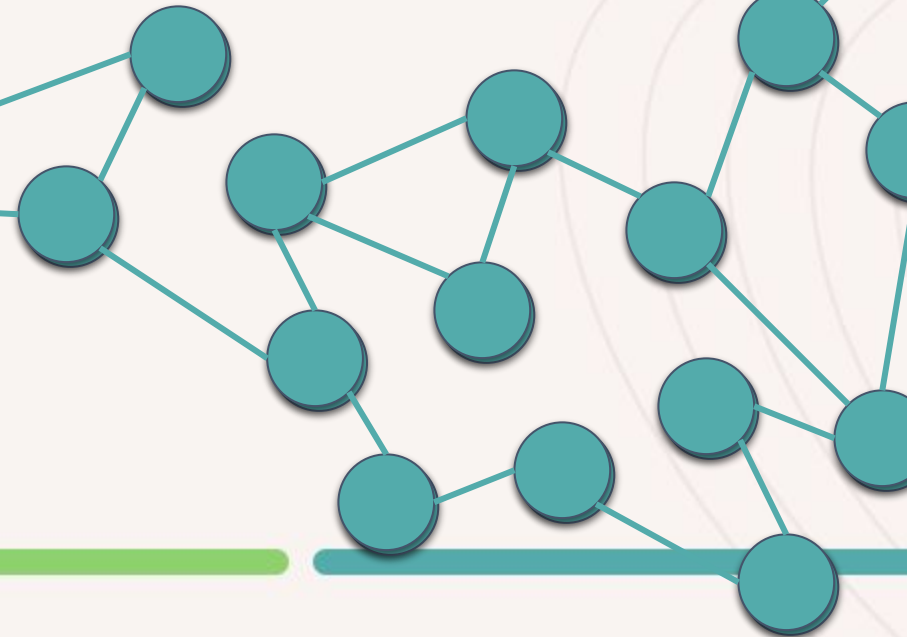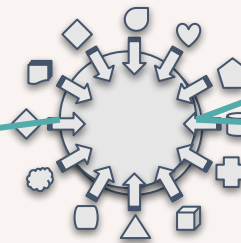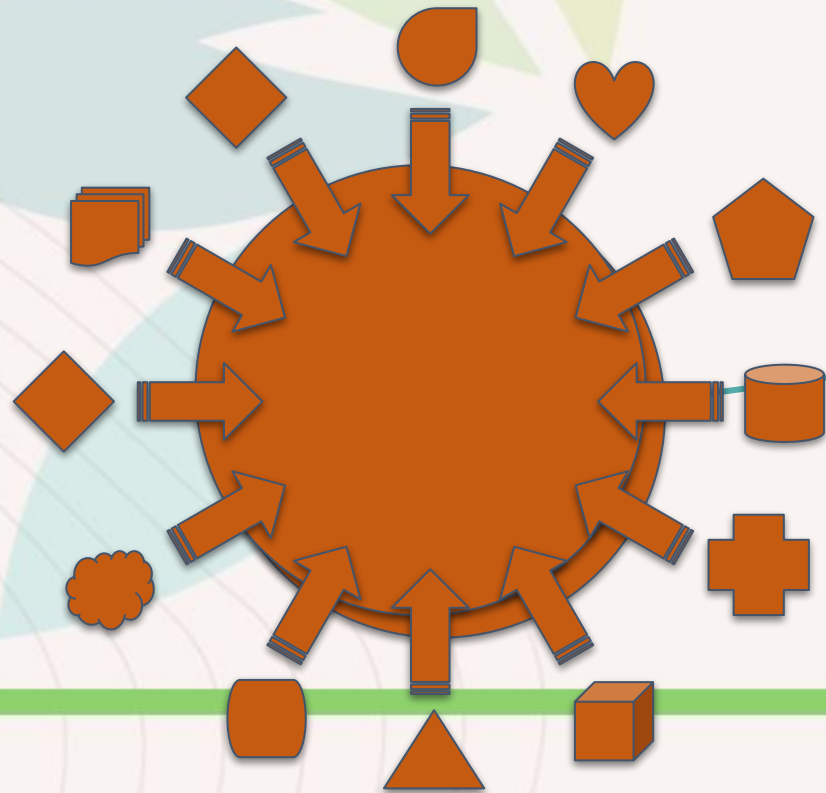4. We do not and cannot copy data from source to our data lake

Requirements

1. To allow all FAIR-EASE data to be findable and accessible
2. To allow sharing of data with the community
3. To allow for search and access to subsets and recombinations of data
4. To provide download management and manage user authentication
5. To be flexible to evolving data sets and requirements

Restrictions

1. The diversity of data (types, sizes, accessibility, formats) is enormous and evolving
2. We need to be sustainable: whatever is developed, must continue to exist after FAIR-EASE ends
3. We want a solution that can be adopted by others for their own projects
4. We do not and cannot copy data from source to our data lake

Centralisation

Centralisation

Distribution
+
Uniformisation
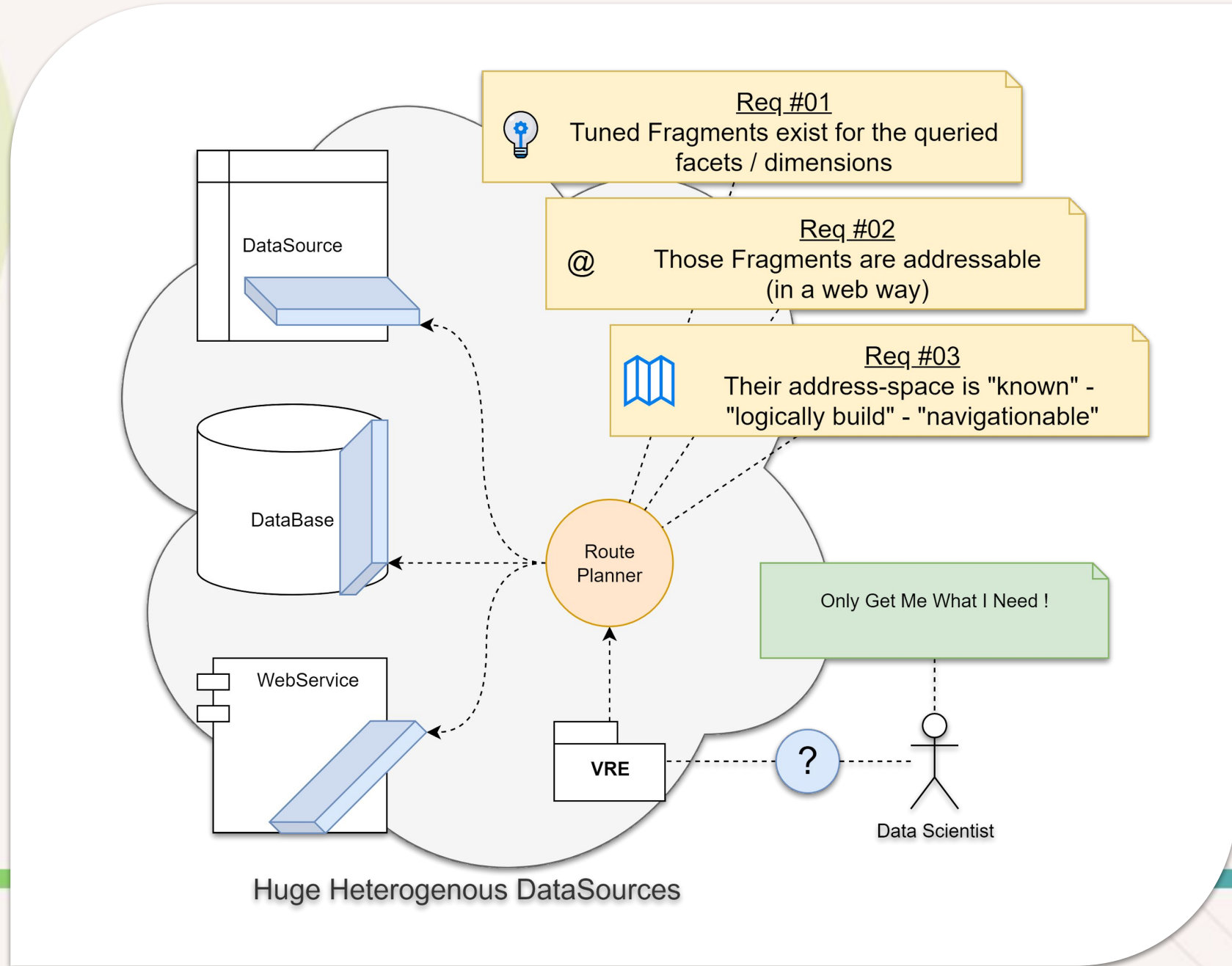
The data can remain *where* they are

The data can remain *as* they are

*But*, a layer of uniformisation sits on top of the data that allows the same findabiliity, accessibility, interoperability, and re-usability for all the data
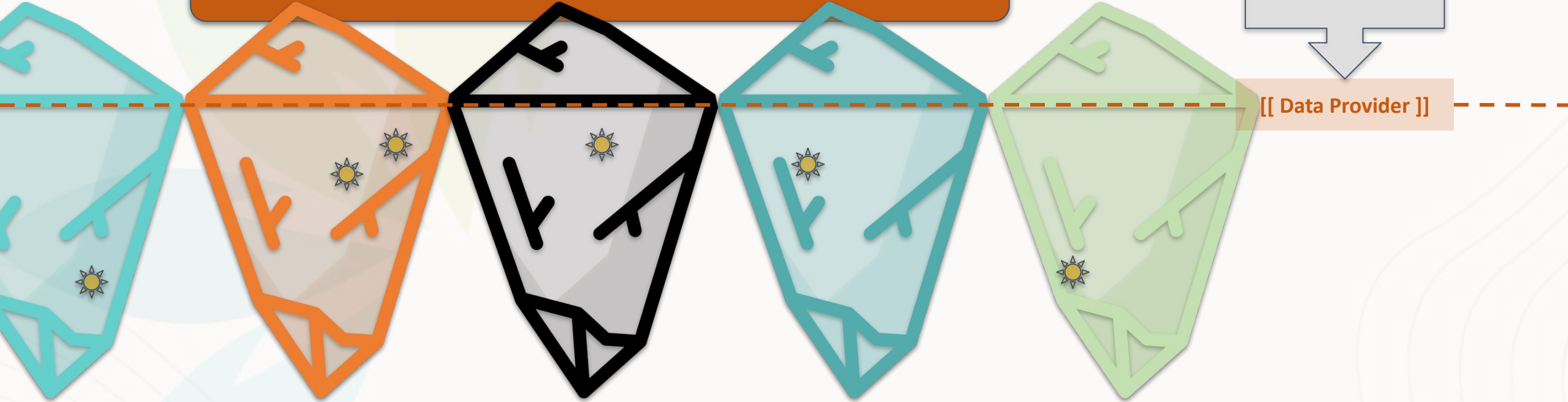
The Data Lake is not a single place where data are stored

The Data Lake is instead the specifications necessary to add this uniformisation layer

In these Large & Diverse Data Sources …

Measurements
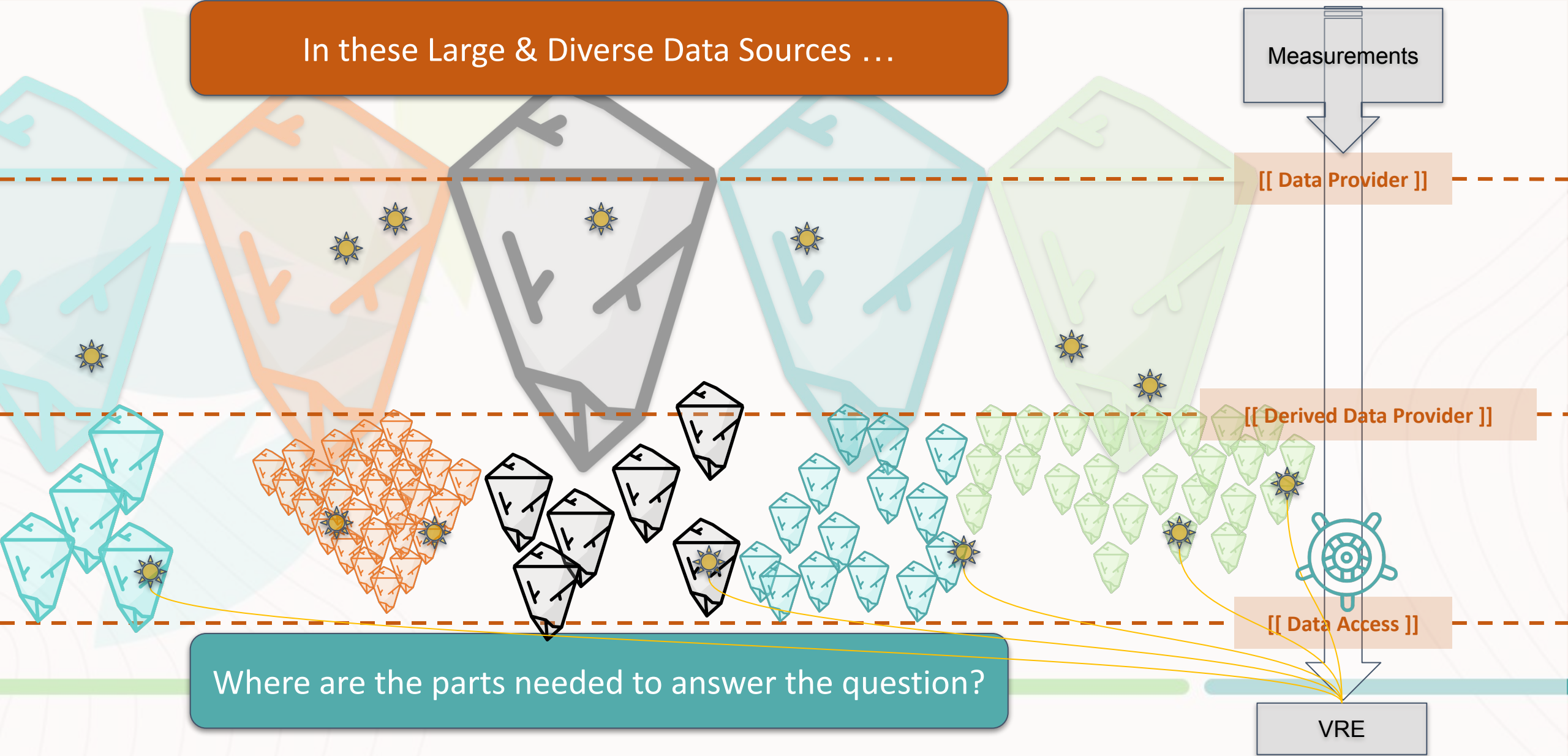
[[ Data Provider ]]

eosc | FAIR-EASE

In these Large & Diverse Data Sources …
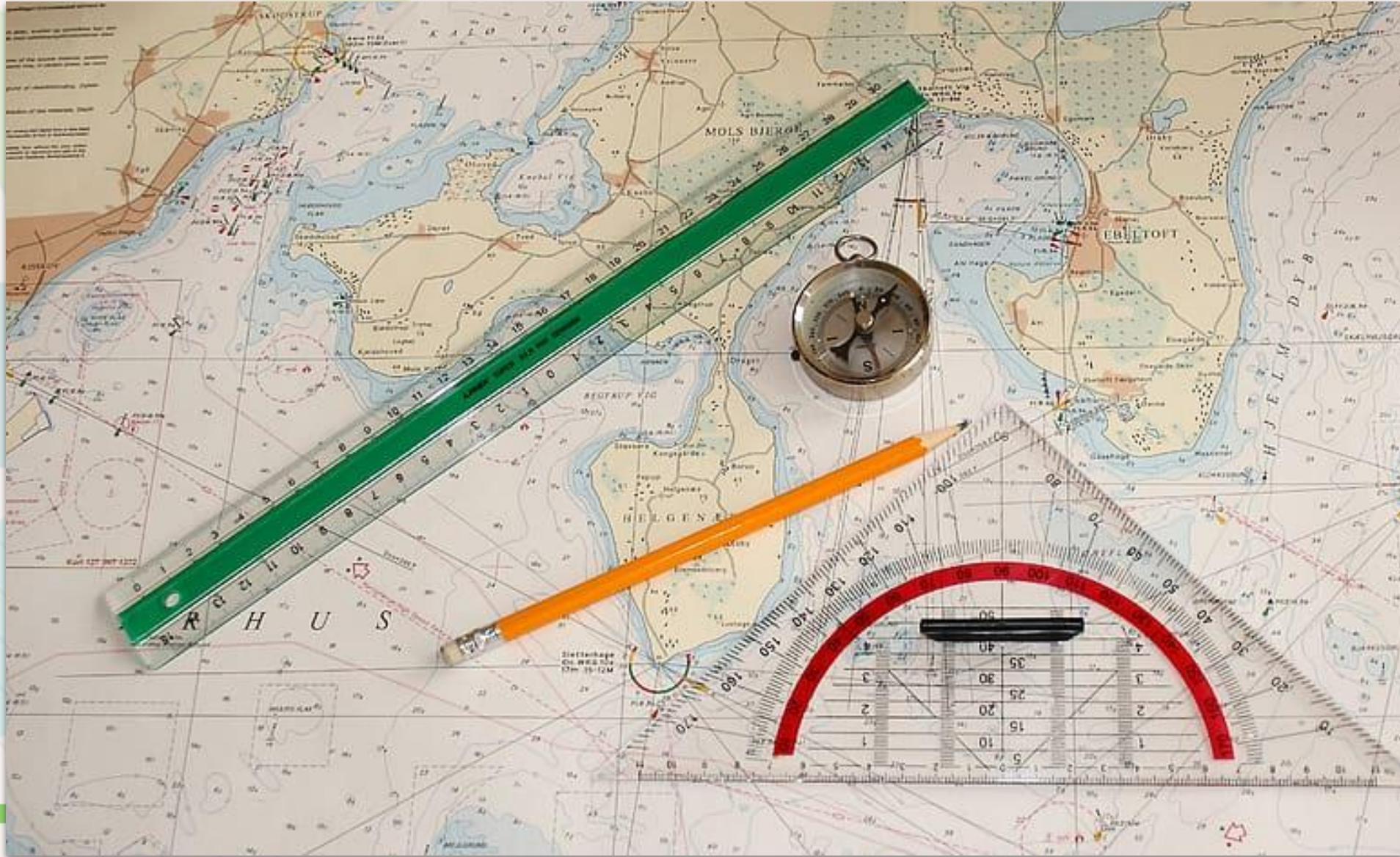
Measurements

[[ Data Provider ]]

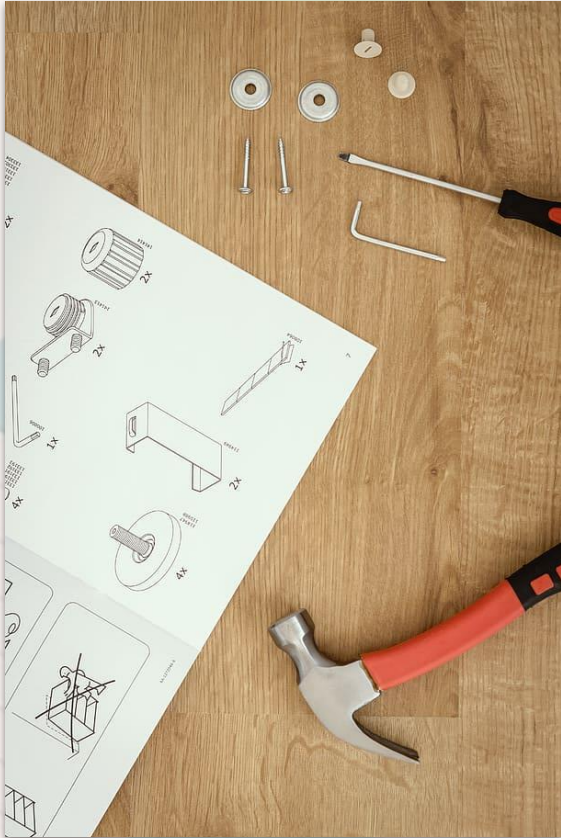[[ Derived Data Provider ]]

[[ Data Access ]]

Where are the parts needed to answer the question?

VRE

- data = the stuff
- **metadata = the stuff about the stuff**
  - provenance (who, where, when, what, how)
  - citation & attribution
  - general findability (keywords, some links)
- **semantics = stuff about what the stuff is about**
  - making it understandable = interoperable
  - making it machine-actionable
  - is very often overlooked
    - provided for humans only
    - understandable to domain experts only
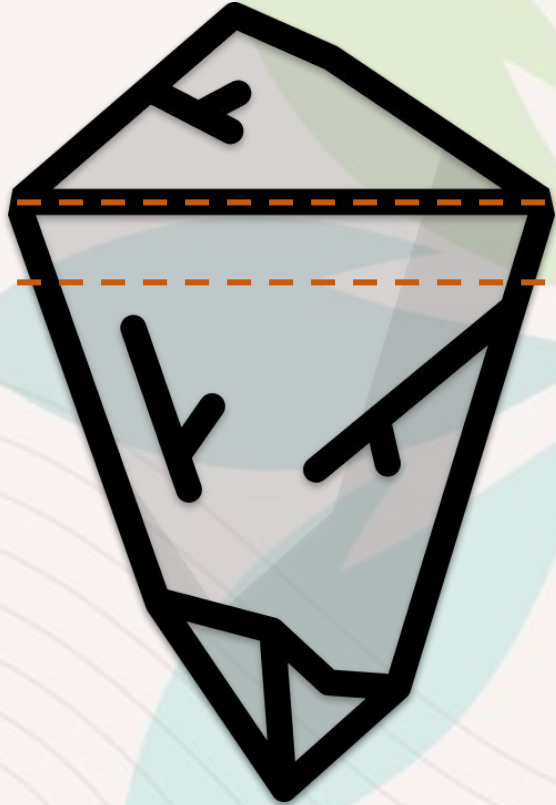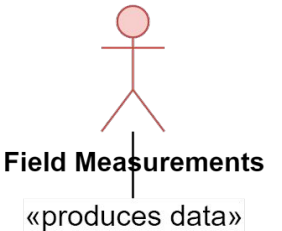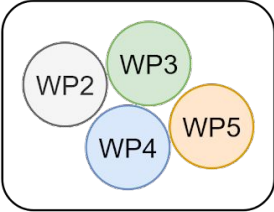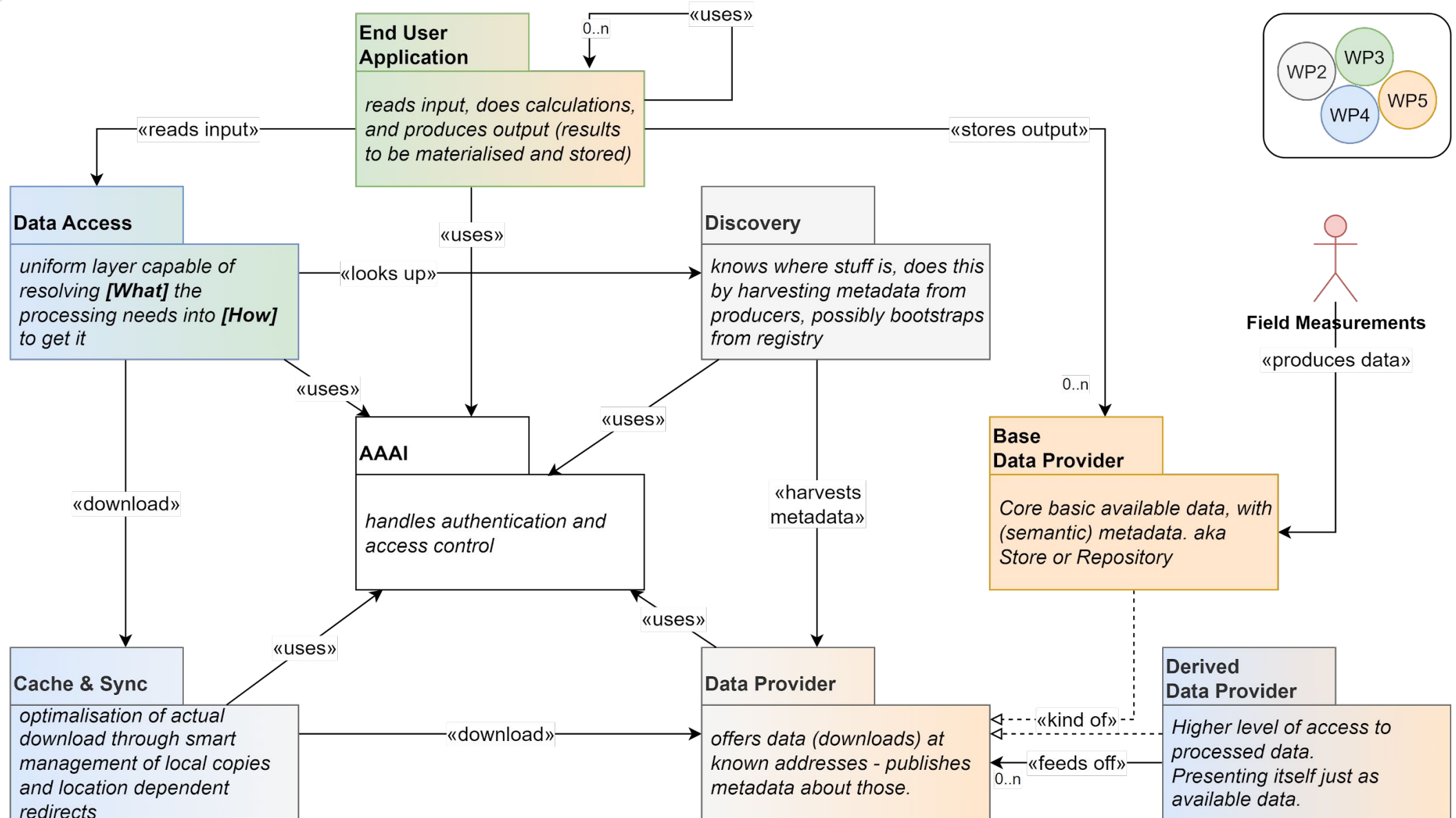    - incomplete
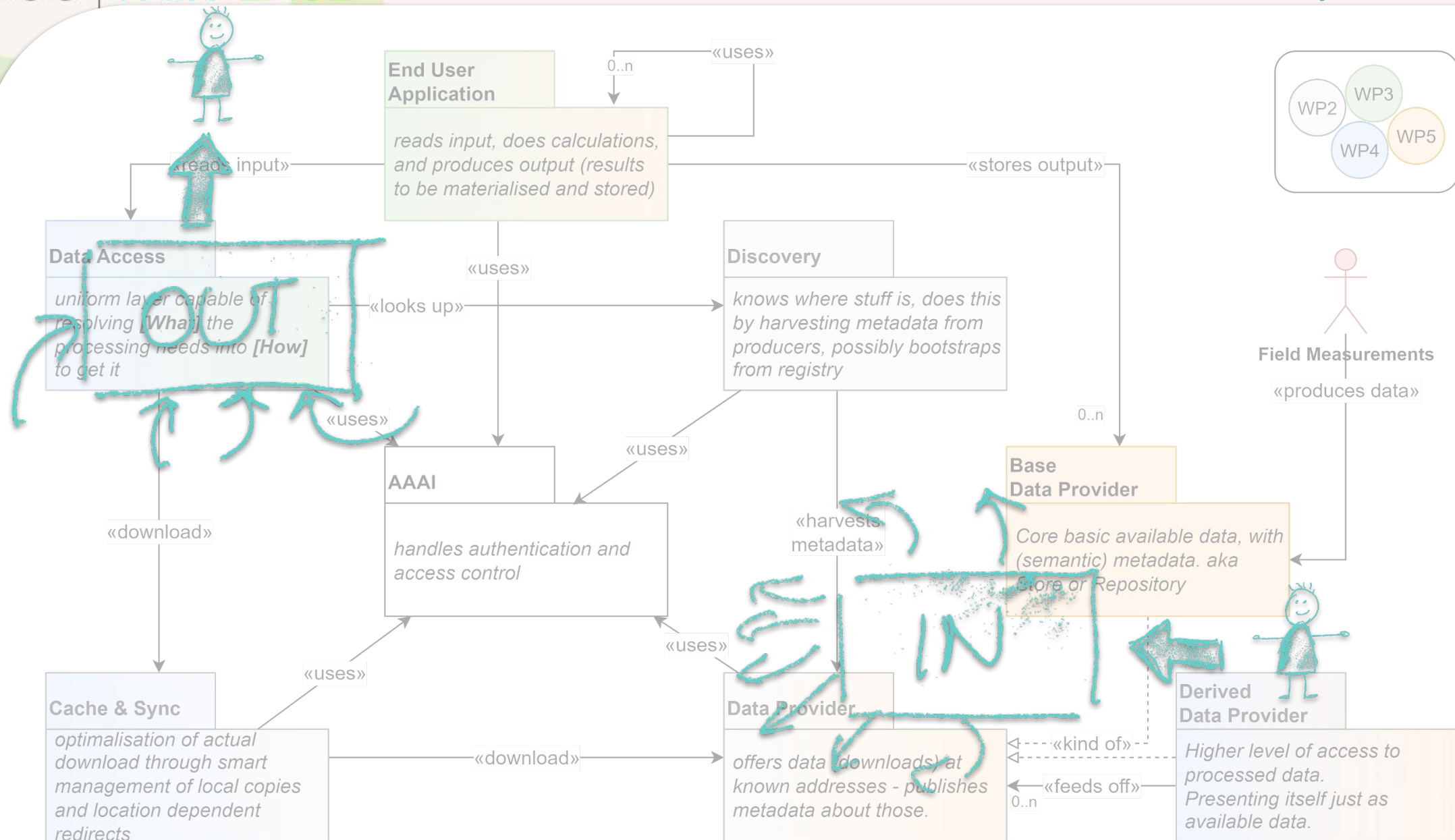
«your new bed»

MALM

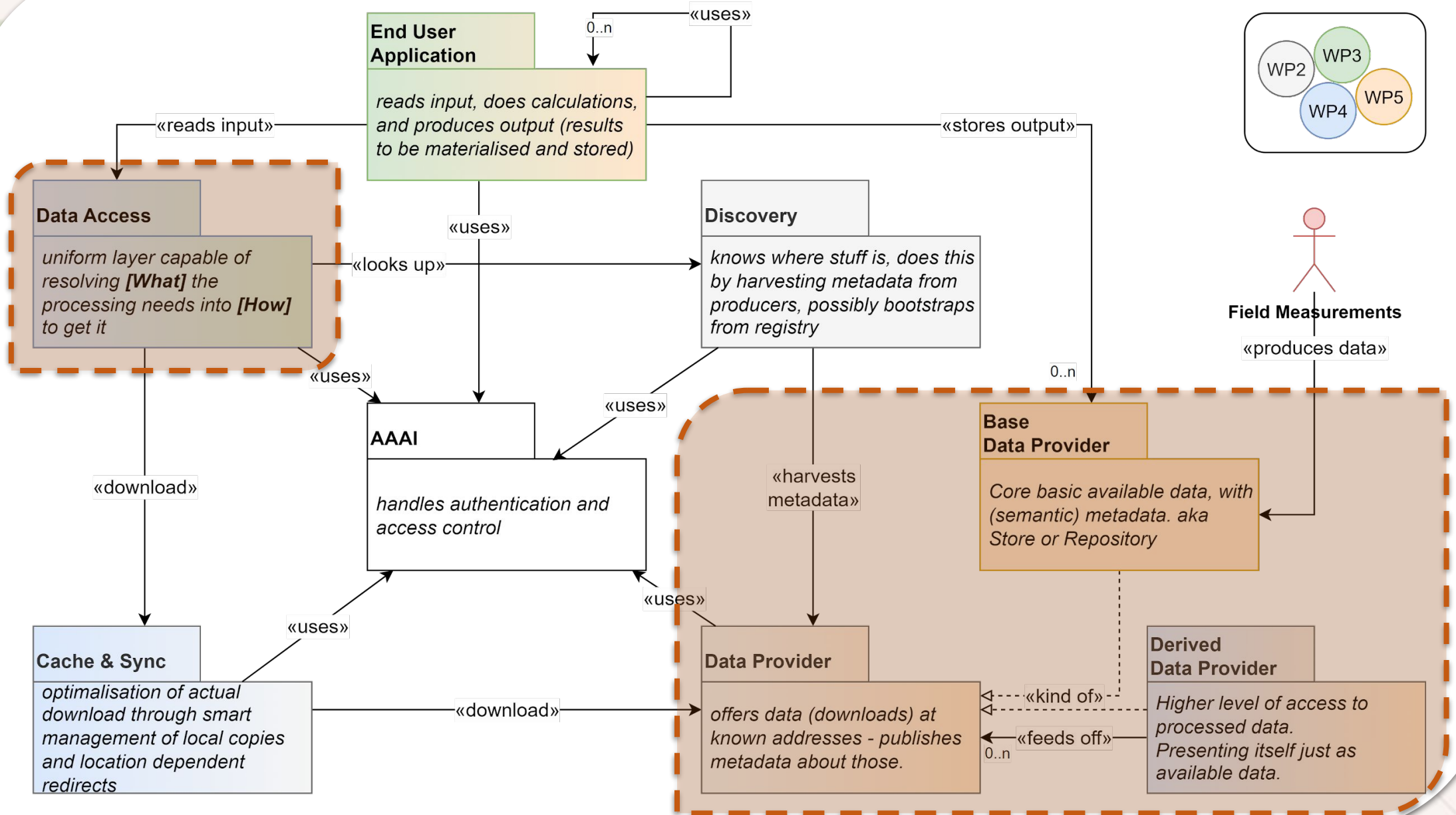299,- €

190 x 210

THE INSTRUCTION KIT

- **the «where» and «when»**
  - temporal coverage
  - spatial coverage  (both emerging via DCAT)
- **the «what» is in there and how is it connected**
  - content coverage (not just title and keywords))
  - shapes (rows, columns)
  - histograms / distributions of available values
  - quality coverage
- **the «how» for optimal retrieval**
  - format
  - size / delivery speed
  - available cache / copies / updates

Questions?

# Data Lake summary

Data Lake will

- Allow data to be *where* they are, *as* they are
- Will be the *specifications* to build an *interoperability layer* over the distributed FE data
- Will allow the DDAS and VRE to plot the routes from the user/client to the requested dataset/data subset/combined data
- Will allow data (sub)sets to be identified and retrieved via URLs, built on-the-fly
- Will have its work cut out for it in figuring out
  - how to obtain all the machine-actionable metadata to allow data to be fully understood
  - how to subset data in an efficient way