

# Semantic annotation pipeline & data publication

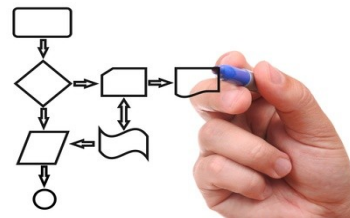
Functional and conceptual aspects regarding our approach to automate the production of semantic data from relational Databases



Functional



Algorithm



Conceptual



Code

Ghislaine Monet

Damien Maurice

Antoine Schellenberger

Yahiaoui Rachid

# INTRODUCING THE PROBLEM

What we are talking about ..

**To have a tool to ensure the production of semantic data as generic and automated as possible.**



What ?

**Simplify the production of semantic data from existing DBs of the AnaEE-F infra and more.**



WHY ?

**By using open source semantic tools + some specific Developments**

which I would talk about later



HOW

# INVENTORY OF THE SEMANTIC OPEN SOURCE PROJECT



Before Developing  
anything

## On-The-Fly Translation tools



- On-the-fly translation allows publishing of RDF from large live databases



- ~~Not intuitive mapping~~  
(~~Especially for those who handle SQL~~)



- ~~Fail last~~ Errors Mapping Raised at runTime

- ~~No Native GUI ! External Project (AuReli)~~



- On-the-fly Ontology-based Data Access



- Intuitive Mapping ( using SQL )

the most interesting  
feature provided by  
ONTOP



- ~~GUI integrated with Protege~~



- Support SPARQL 1.0



# INVENTORY OF THE SEMANTIC OPEN SOURCE PROJECT



Semantic DataBases..  
Two Kinds ...

## TripleStore

### \* Sesame



- Robustness : **K.O**
- Scaling out : **K.O**
- Performance : **Err**

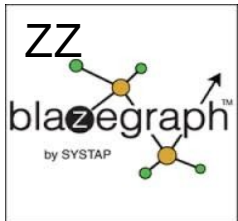
### \* Sol-RDF



- Robustness : **OK**
- Scaling out : **OK**
- Performance :

REST

### \* BlazeGraph



- Robustness : **OK**
- Scaling out : **OK \***
- Performance : **OK**

### \* Corese



- Robustness : **OK**
- Scaling out :
- Performance : **OK**

\* : Version 1.5.3

## Graph Databases

Has more generalized  
structure than a triplestore

11-13 Octobre 2017

PARIS

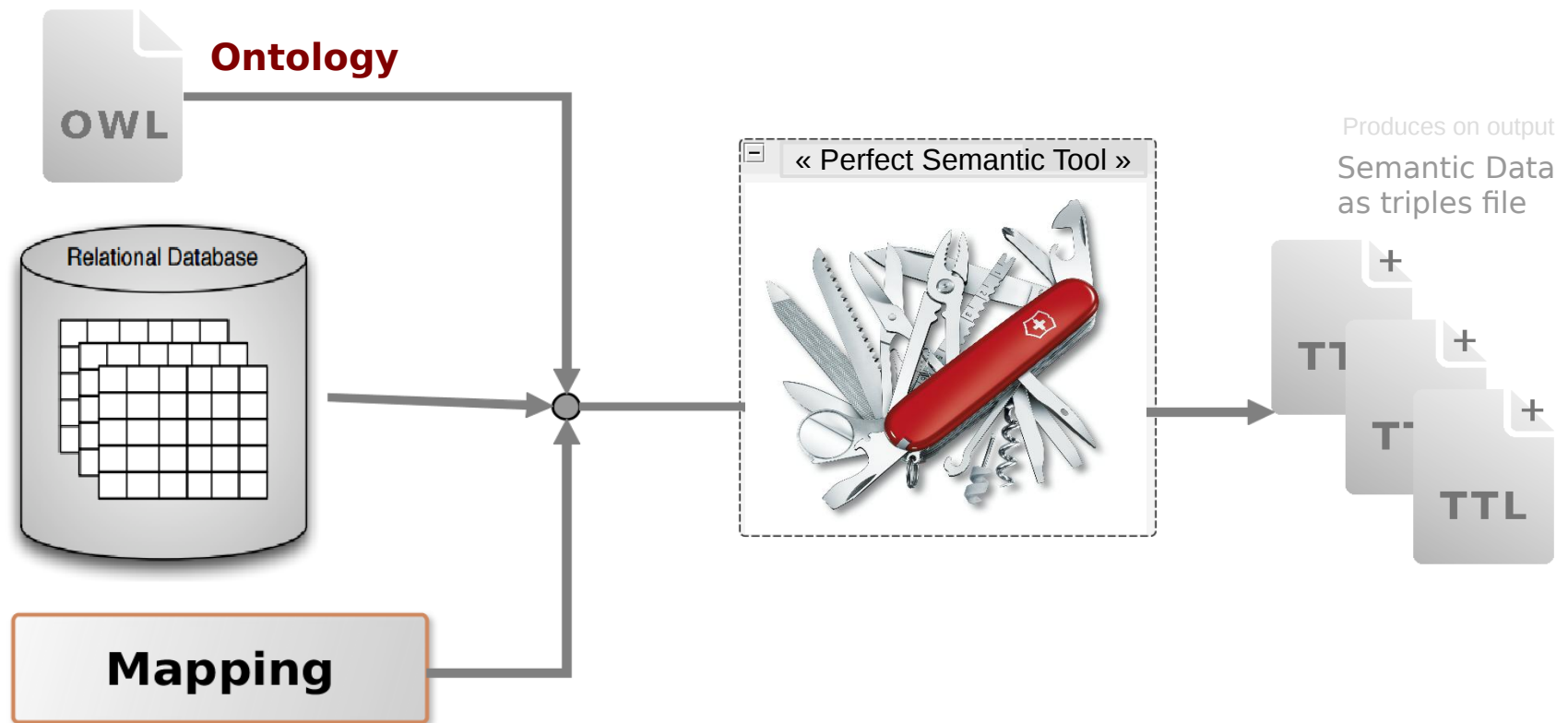
AnaEE - Envri+ f2f meeting

4

# AUTOMATION APPROACH



what's the perfect semantic tool in our case ?

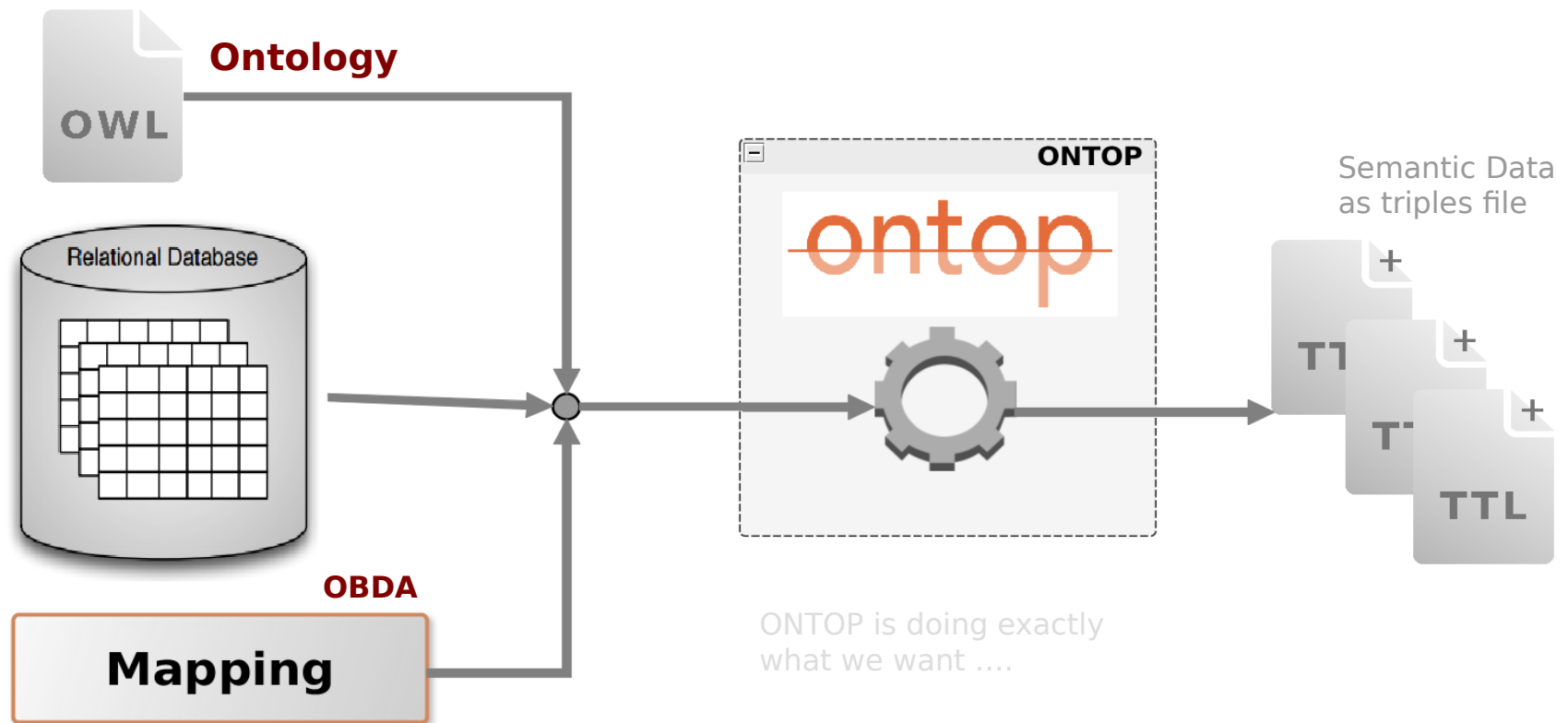


Something which says how relational data are transformed to semantic data graph

# AUTOMATION APPROACH



Good news :-)



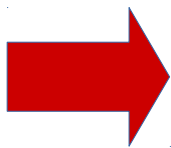
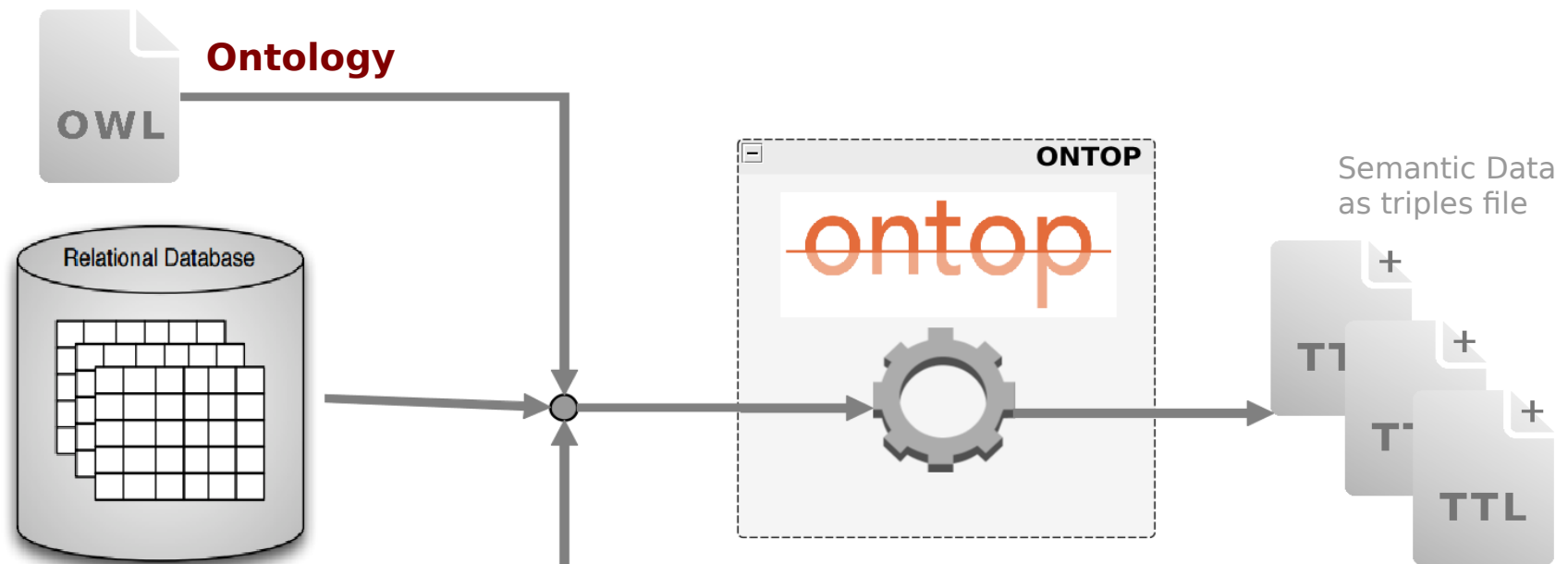
Something which says how relational data are transformed to semantic data graph

ONTOP is doing exactly what we want ....

But ??



# AUTOMATION APPROACH



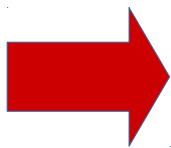
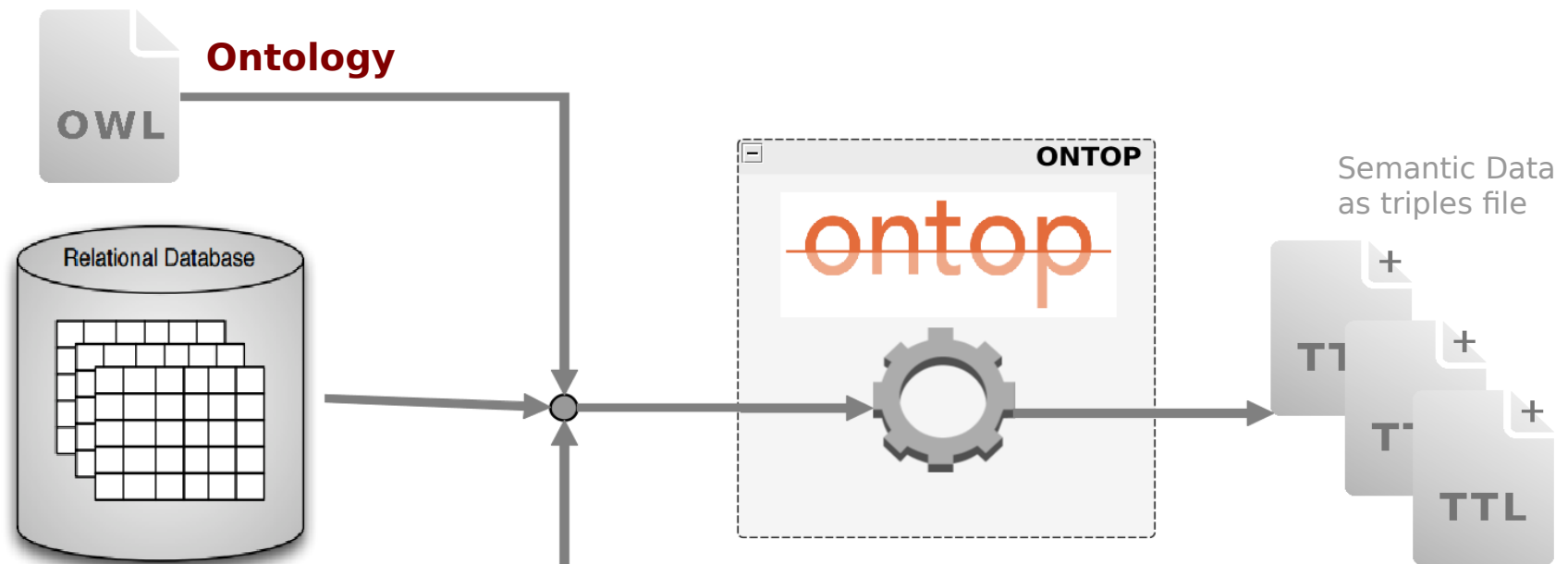
**Limitations**



- **Performance**  
( I'll explain right after )
- **Limited Inference ( due to QUEST )**
- **No Automation**



# AUTOMATION APPROACH



Limitations



**almost perfect**  
**~~semantic tool~~**

That we are  
looking for

In our case ( of course )

So, let's see how  
can these limitations be solved..



Mapping has to be  
created by hand

# AUTOMATION APPROACH



Level of  
Automation

## Ontop From Protegé \*

Ontop provides  
a nice  
integrated GUI  
into Protegé to  
Create mapping  
files.

Once Protegé  
opened, we  
distinguish  
**3 main parts**

Protegé :  
open-source  
Ontology editor

The screenshot displays the Protegé ontology editor interface. On the left, the 'Class hierarchy' pane shows a tree structure starting with 'Thing', followed by 'Characteristic', 'Name', 'O2', 'Physical Characteristic', 'Dimension', 'Relationship', 'Type', 'Characteristic Qualifier', 'Base Characteristic Qualifier', 'Average', 'Maximum', 'Minimum', 'Composite Characteristic Qualifier', 'Entity', 'Primitive Value', 'Boolean', 'Decimal', 'String', 'SpatialEntity', 'Lac', 'PlateForme', 'TemporalEntity', 'PickingDate', 'Water', 'Measurement', 'Context', 'MetaData', 'Observation', 'Observation Collection', 'Protocol', 'Standard', 'NominalStandard', 'DateTime', 'UniqueIdentifier', 'Unit', 'Base Unit', 'Composite Unit', 'Derived Unit', and 'Unit Conversion'. The 'Mapping editor' pane on the right shows a mapping for 'measurement-ph-water' with a 'Target (Triples Template)' and a 'Source (SQL Query)'. The 'Target' part is labeled '(2)' and the 'Source' part is labeled '(3)'. The 'Datasource editor' pane at the bottom right shows connection parameters for a PostgreSQL database, labeled 'DB Part (1)'. The 'Mapping count' is 7, and the 'Search' field contains 'pred: Measurement'.

**Target Part (2)**

How data are mapped

I'll explain this right after

**Source Part (3)**

Which data to map

**DB Part (1)**

Connection parameters

Connection URL: jdbc:postgresql://127.0.0.1/ola

Database User: ryahiaoui

Database Password: .....

Driver class: org.postgresql.Driver

Test Connection

Inform connection information to DB

Mapping count: 7 Search: pred: Measurement

# AUTOMATION APPROACH



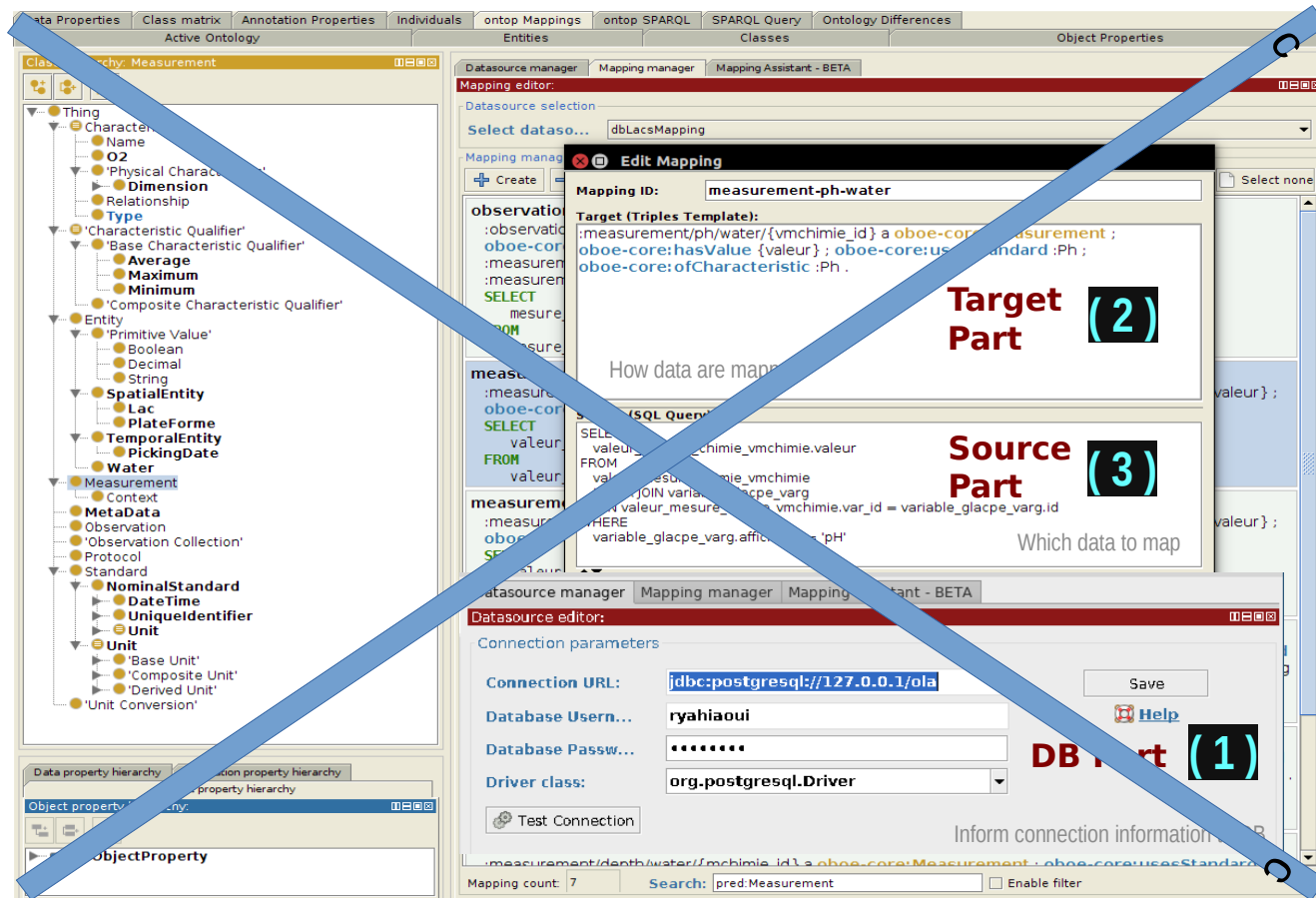
Level of  
Automation

Ontop From  
Protegé \*



This approach is

**Not an effective  
solution for  
automation**



11-13 Octobre 2017

PARIS

AnaEE - Envri+ f2f meeting

10

Take a look at what there ..

# AUTOMATION APPROACH



Level of Automation

## Behind the scene ... Ontop manipulates OBDA File ( based on [R2RML](#) language)

( [R2RML](#) is a W3C recommended RDB-to-RDF mapping language )

```
[PrefixDeclaration]
rdf: http://www.w3.org/1999/02/22-rdf-syntax-ns#
oboe-core: http://ecoinformatics.org/oboe/oboe.1.0/oboe-core.owl#
oboe-temporal: http://ecoinformatics.org/oboe/oboe.1.0/oboe-temporal.owl#
xsd: http://www.w3.org/2001/XMLSchema#
: http://www.anaee.fr/ontology/anaee-france_ontology#
oboe-standard: http://ecoinformatics.org/oboe/oboe.1.0/oboe-standards.owl#
oboe-characteristics: http://ecoinformatics.org/oboe/oboe.1.0/oboe-characteristics.owl#
oboe-spatial: http://ecoinformatics.org/oboe/oboe.1.0/oboe-spatial.owl#
oboe-standards: http://ecoinformatics.org/oboe/oboe.1.0/oboe-standards.owl#
rdfs: http://www.w3.org/2000/01/rdf-schema#

[SourceDeclaration]
sourceUri dbLacsMapping
connectionUrl jdbc:postgresql://127.0.0.1/ola?sendBufferSize=5000
username ryahiaoui
password yahiaoui
driverClass org.postgresql.Driver

[MappingDeclaration] @collection []
mappingId (52) ola_characteristic_depthRelativeToSurface_min
target :ola/characteristic/depthRelativeToSurface/min a :DepthRelativeToSurface
oboe-core:hasQualifier :Minimum
source SELECT id from (values ('1')) s(id) ;

]]
```

Interesting thing .. The 3 parts that was discussed previously

how can we automate the generation of this kind of the mapping files

(1) Inform connection information to DB

(2) How data are mapped

(3) Which data to map

## 3 Important parts

### \* DB Part

#### (1) DB Access Informations

### \* Target Part

#### (2) How data are mapped

**Rule :** Graphs are composed of nodes, each non terminal node is identified by an URI.

**Target Part = Graph + URIs**

Based on **turtle** syntax

### \* Source Part

#### (3) Which data to map

Using SQL Queries

## AUTOMATION APPROACH



Level of  
Automation

**Piece of Target Part**

According to the modeling graphs presented earlier, and according to the rule : **Target Part = Graph + URI**

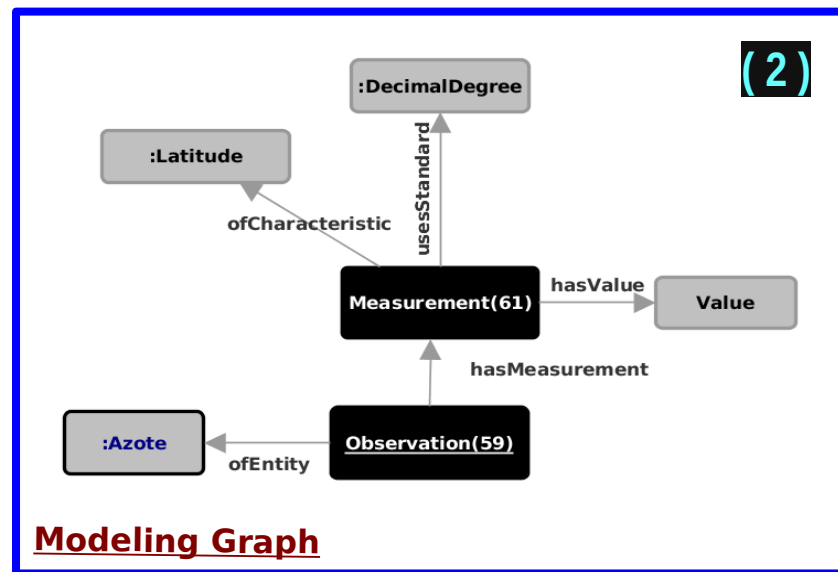
**Modeling graphs can be considered as a piece of the target Part.**

What is missing here is something which allows us to uniquely identify each node of the graph, these unique things are URI

In semantic approach, the simplest way to modelize is to use Graphs



**Why not transform  
this Modeling  
Graphs to ODBA  
mapping file ???**



# AUTOMATION APPROACH



Level of  
Automation

Target Part

Connection information to DB (1)

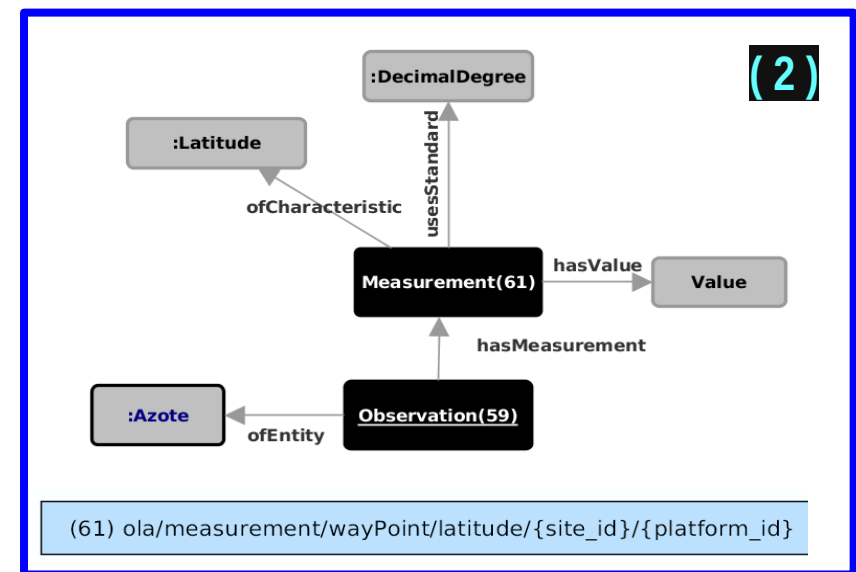
obda-sourceUri : dbLacsMapping

obda-connectionUri : jdbc:postgresql://127.0.0.1/ola?sendBufferSize=5000

obda-username : ryahiaoui

obda-password : yahiaoui

obda-driverClass : org.postgresql.Driver



yedGen was specifically written to generate obda files from modeling graphs



Assign SQL Query for each non terminal node (3)

Query\_(61) : SELECT pla.loc\_id AS platform\_id, site.id AS site\_id, pla.latitude AS latitude  
FROM  
public.site\_glacpe\_sit site INNER JOIN public.plateforme\_pla pla ON site.id = pla.id

Source Part

This is how was approached the Automation problem

11-13 Octobre 2017

PARIS

AnaEE - Envri+ f2f meeting

13



# AUTOMATION APPROACH



## Genericity

The genericity we talk about here concerns the functioning of yedGen

## About Genericity ?

The idea is to generate multiple instances of the same Graph according to different variables described in file ( CSV ).

**Why ? Because these variables has the same Structuration in the data Bases.**

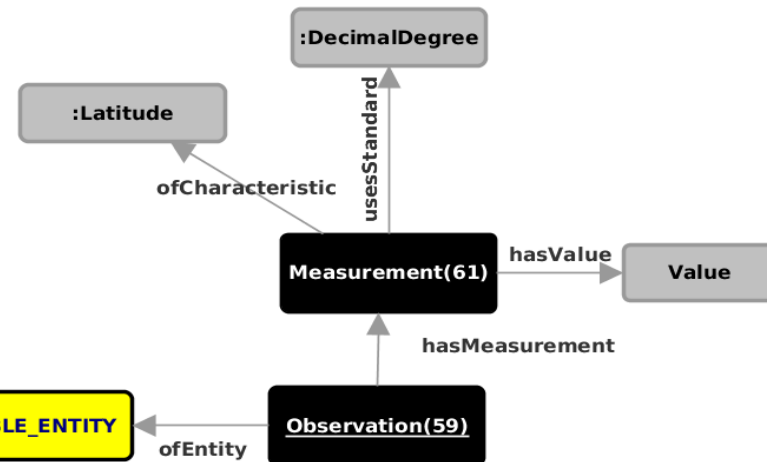
Instead to create one graph per variable, we use a graph type ( graph designed for several variables ) in order to create instances of this graph

## csv file of variable semantic description

|   | AnaEE Standar       | Entity          | Context     | .. |
|---|---------------------|-----------------|-------------|----|
| 1 | cumulative rainfall | cumulative rain |             | .. |
| 2 | air carbon dioxide  | carbon dioxide  | atmosphere, | .. |
| 3 | atmospheric air sta | air             | atmosphere  | .. |

Apply on VARIABLE\_ENTITY  
each value of the column  
Enty from CSV ,  
which gives us ...

?VARIABLE\_ENTITY



11-13 Octobre 2017

PARIS

AnaEE - Envri+ f2f meeting

14



# AUTOMATION APPROACH

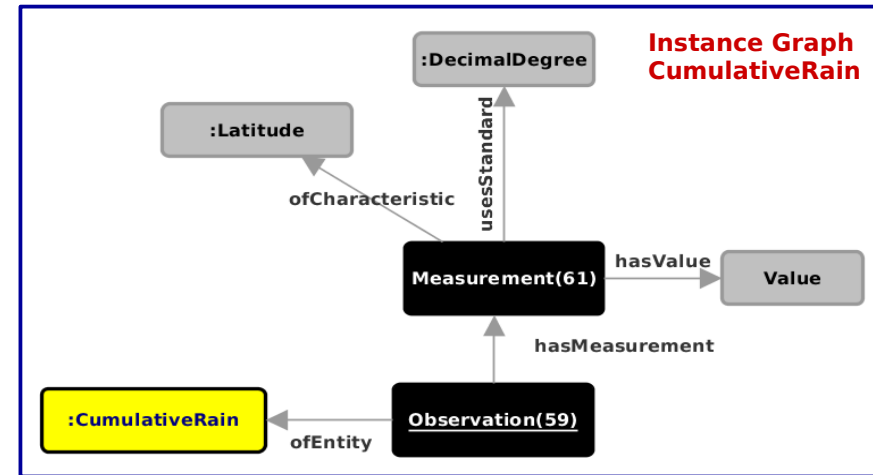


Genericity

## About Genericity ?

### csv file of variable semantic description

|   | AnaEE Standar       | Entity          | Context     | .. |
|---|---------------------|-----------------|-------------|----|
| 1 | cumulative rainfall | cumulative rain |             | .. |
| 2 | air carbon dioxide  | carbon dioxyde  | atmosphere, | .. |
| 3 | atmospheric air sta | air             | atmosphere  | .. |



# AUTOMATION APPROACH

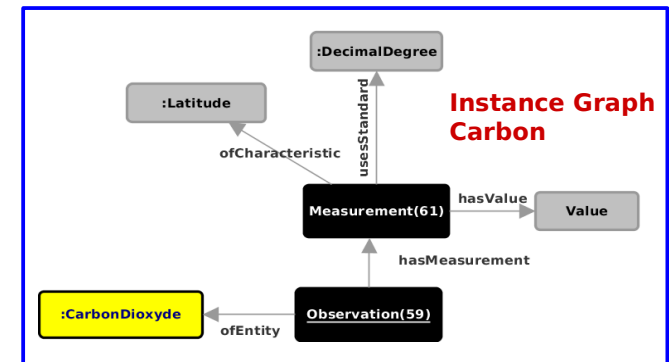
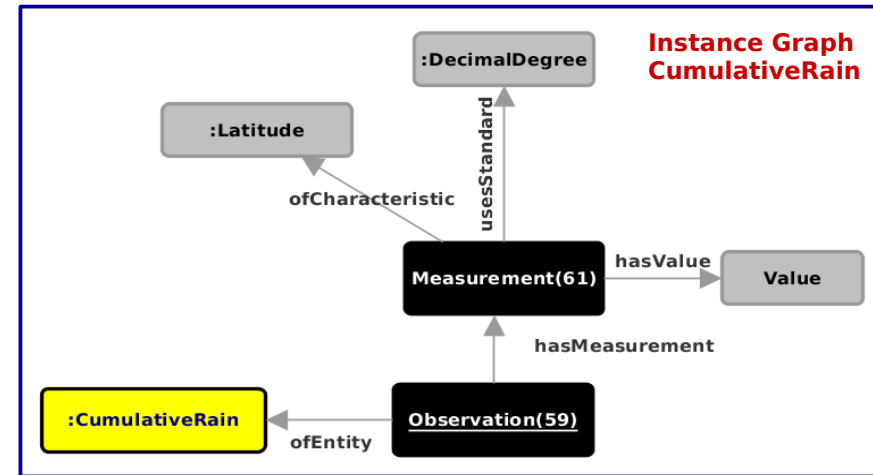


Genericity

## About Genericity ?

### csv file of variable semantic description

|   | AnaEE Standar       | Entity          | Context     | .. |
|---|---------------------|-----------------|-------------|----|
| 1 | cumulative rainfall | cumulative rain |             | .. |
| 2 | air carbon dioxide  | carbon dioxyde  | atmosphere, | .. |
| 3 | atmospheric air sta | air             | atmosphere  | .. |

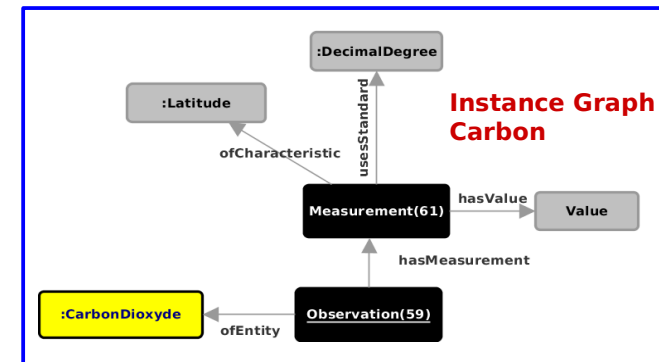
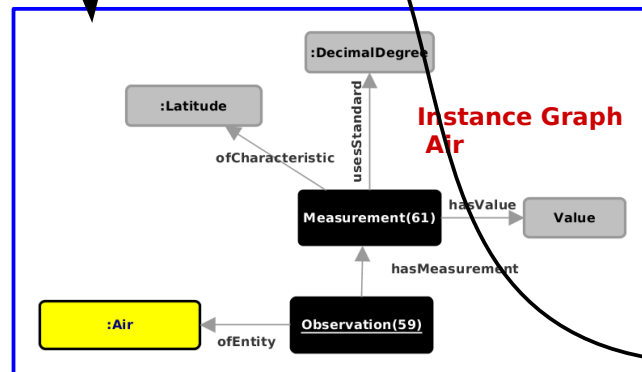
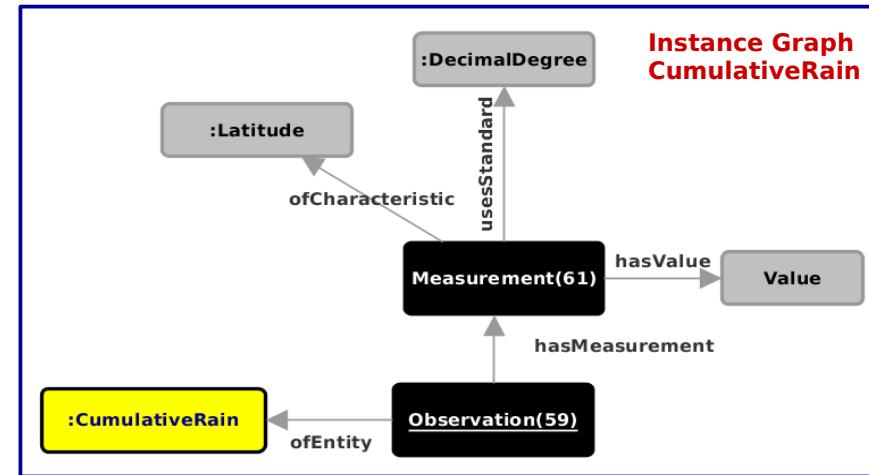




## About Genericity ?

### csv file of variable semantic description

| AnaEE Standar         | Entity          | Context     | .. |
|-----------------------|-----------------|-------------|----|
| 1 cumulative rainfall | cumulative rain | ..          | .. |
| 2 air carbon dioxide  | carbon dioxide  | atmosphere, | .. |
| 3 atmospheric air sta | air             | atmosphere  | .. |



# AUTOMATION APPROACH



Genericity

## About Genericity ?

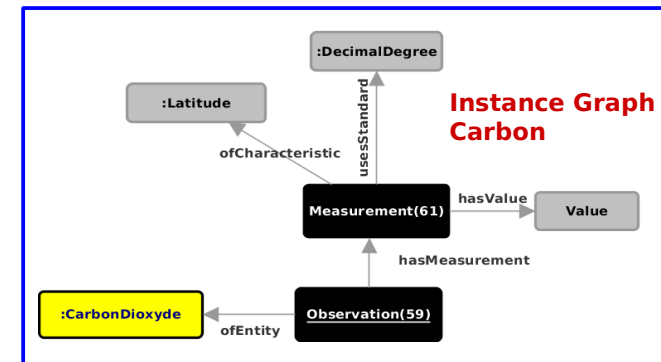
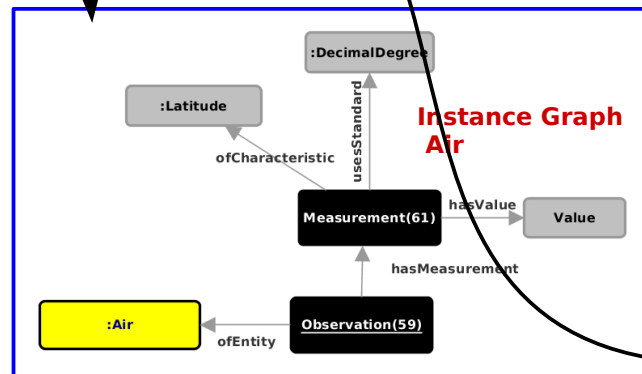
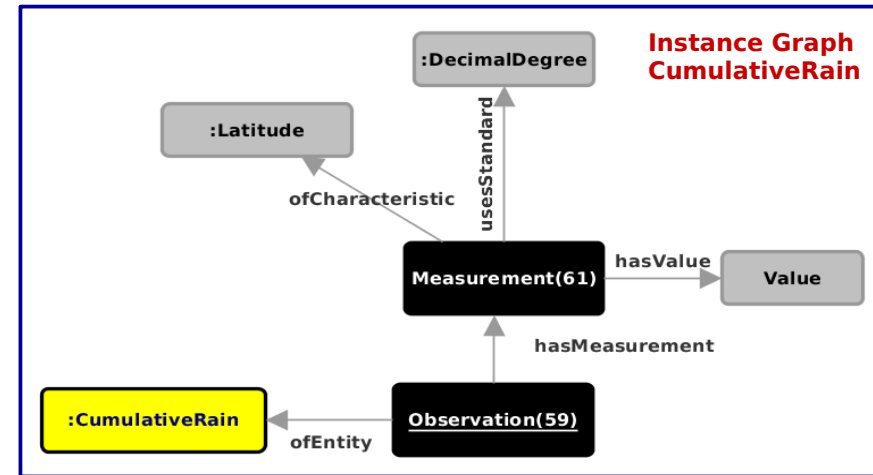
### csv file of variable semantic description

| AnaEE Standar         | Entity          | Context     | .. |
|-----------------------|-----------------|-------------|----|
| 1 cumulative rainfall | cumulative rain |             | .. |
| 2 air carbon dioxide  | carbon dioxyde  | atmosphere, | .. |
| 3 atmospheric air sta | air             | atmosphere  | .. |

The same process is repeated for each line..



This is how was approached the Genericity problem



# AUTOMATION APPROACH



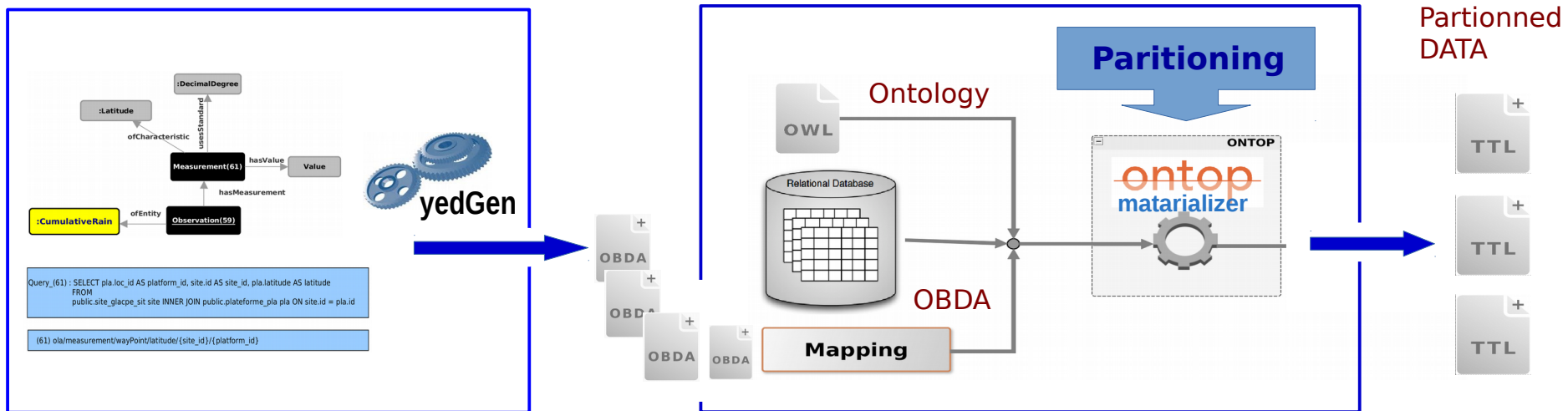
Performance(1)

Partitioning Data

## About Performance ( on large amounts of data ) ?

It sometimes happens that the amount of data to be processed by ONTOP and BlazeGraph exceeds the memory capacity, in this case Outofmemoryerrors are raised.

**Solution : Volume data Partionning** → Processing data by chunk



→ Process « infinite » volume of data

Filtering : a way to  
increase performance

## AUTOMATION APPROACH



Performance(2)

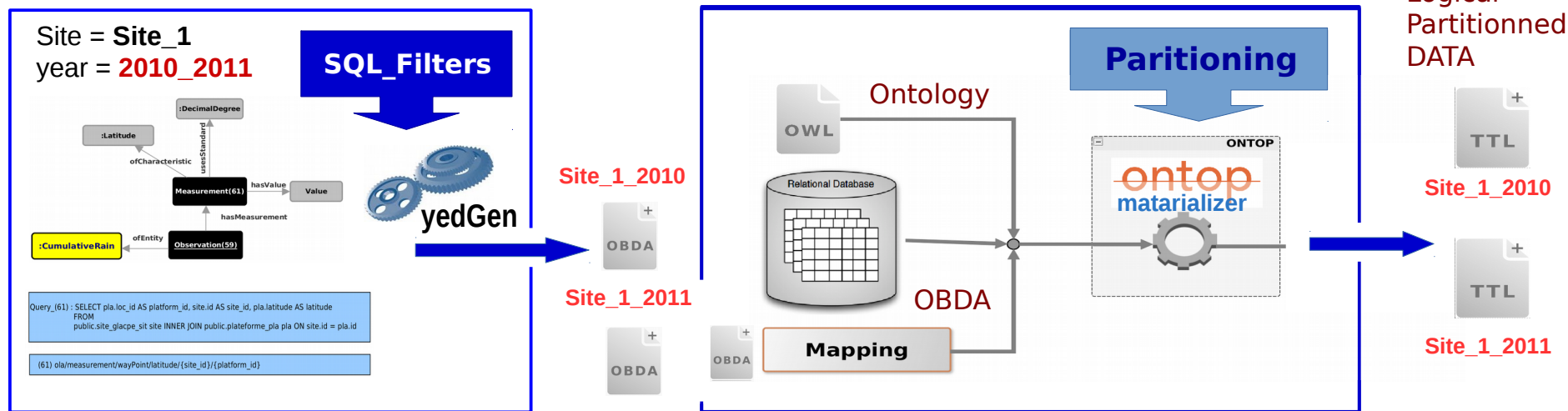
### About Performance ( Filtered Data )

For some use cases, we have to **extract only** the data that users **needs**

**Solution : Logical data partitionning**

Generate Data for a specific variable with specific Site and specific range of years for example.

the more you filter, the  
less you process data,  
the more you're  
performant



Volume data partitioning & Logical  
data partitioning can be **combined**

→ Process « infinite » volume of data

11-13 Octobre 2017

PARIS

AnaEE - Envri+ f2f meeting

20

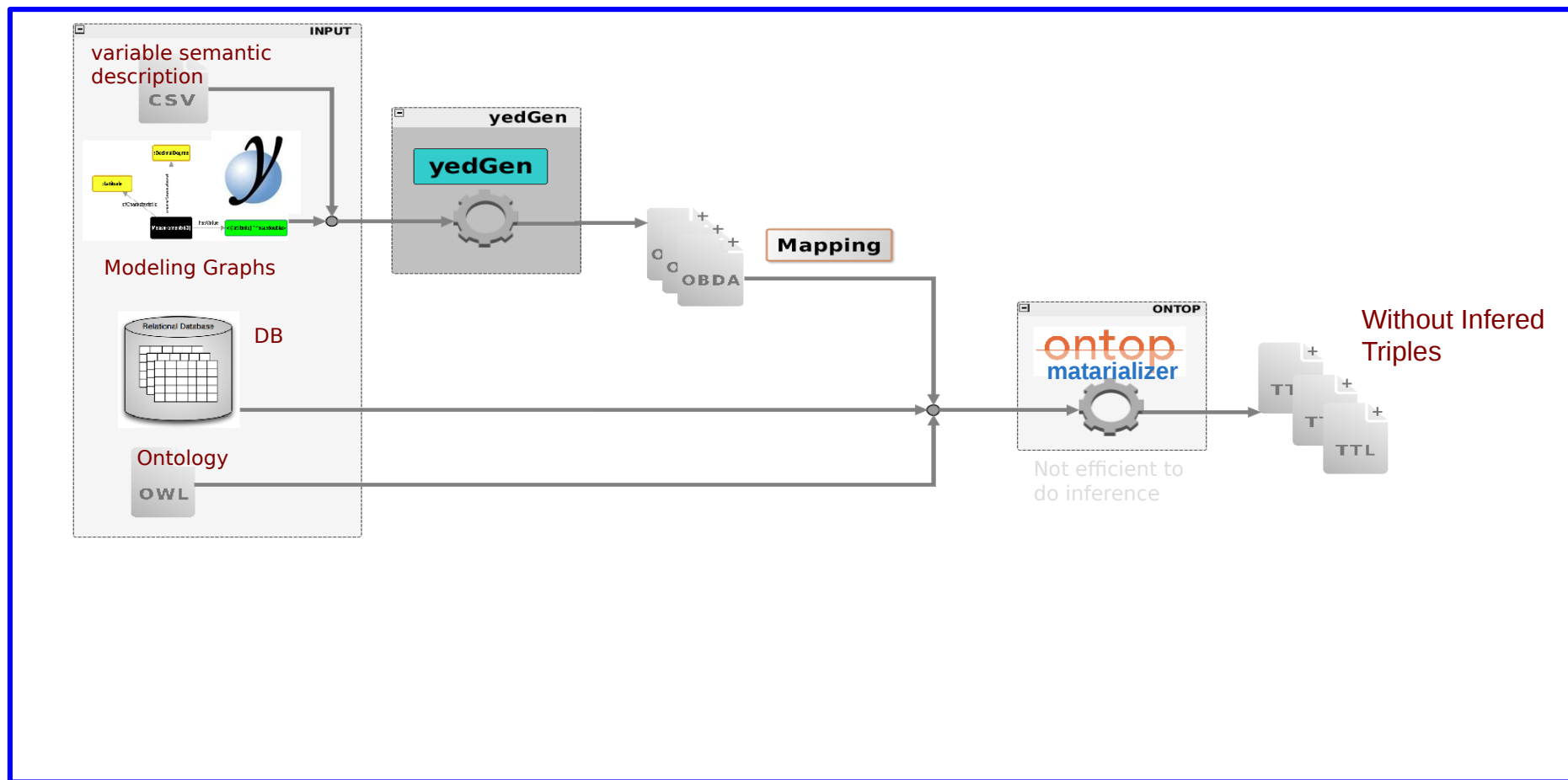


According to what  
said previously

# AUTOMATION APPROACH



Recap

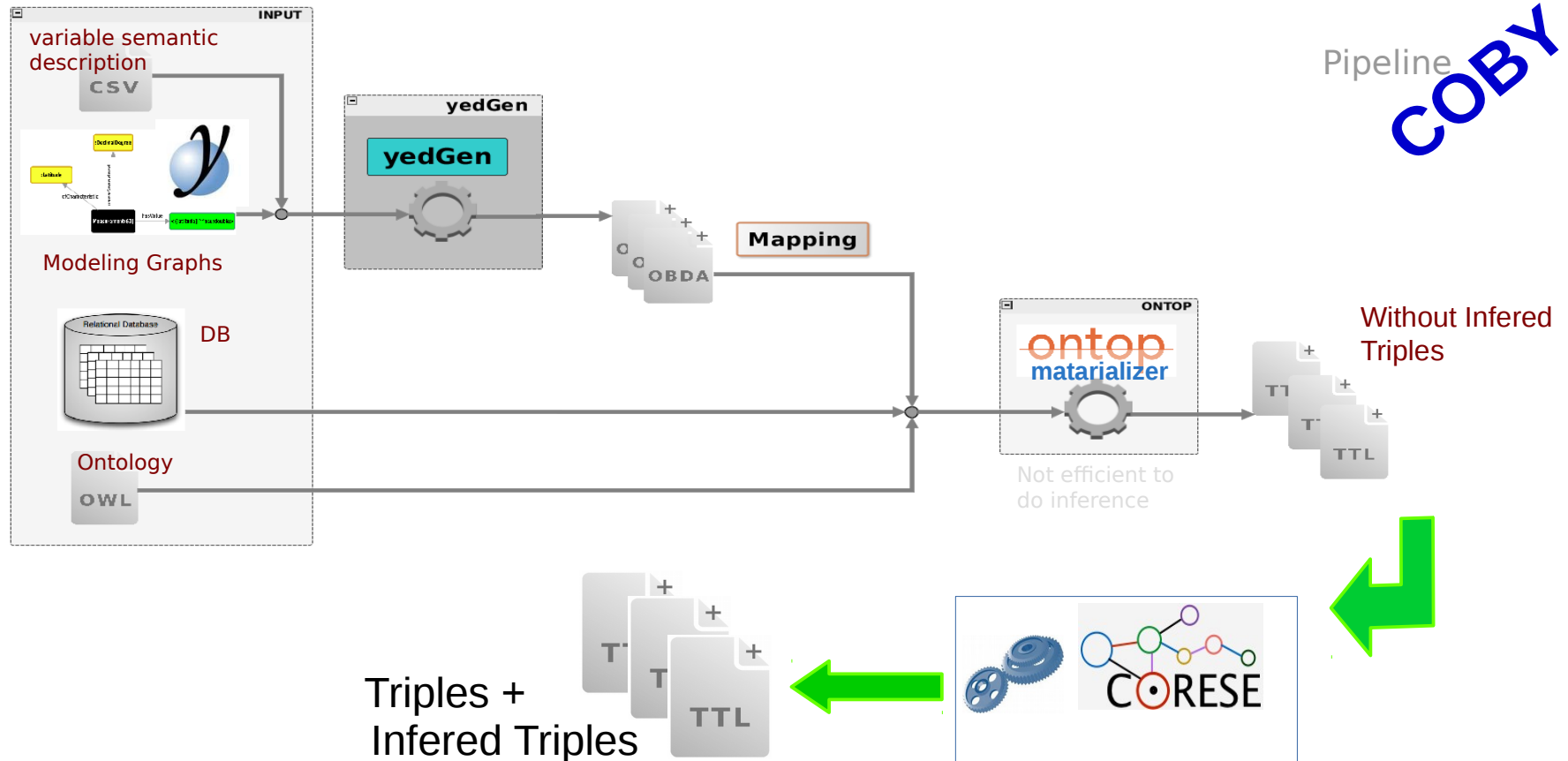


# AUTOMATION APPROACH



Recap

Pipeline  
**COBY**

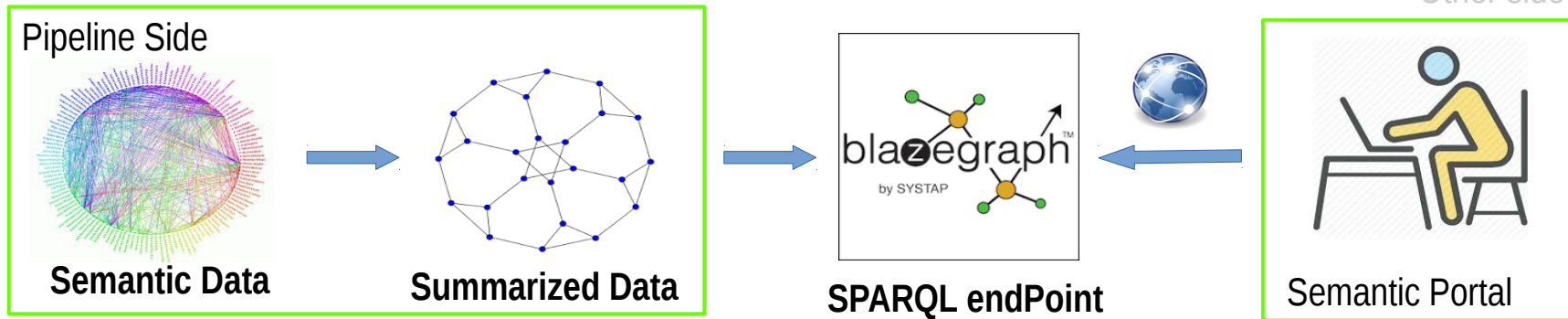


# USE CASES



Concrete use cases  
of the Pipeline ?

## 1 - Semantic Summarized Data for Portal Use Case 1



**Purpose :** Produce a semantic summarized data using the Pipeline and publish it on a specific SPARQL Endpoint ( blazegraph ) which be consumed by an external entity ( AnaEE-F Portal )

## 2 - Production of netCDF file Use Case 2



**Purpose :** Produce **Filtered** Semantic data ( in TTL Format ) that will be used to produce **netCDF** Files

Session for each of these use cases

11-13 Octobre 2017

PARIS

AnaEE - Envri+ f2f meeting

23

## EVOLUTION



### **\*\* SOERE Level :**

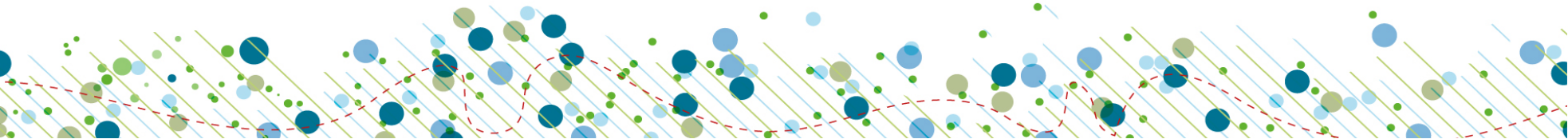
Modeling new Data types

Consists on Creation of new annotation models for other variables stored in relationnal data bases with Yed Graph Editor

### **\*\* Pipeline Level :**

**Improve performance by introducing distributed processing [ With Docker\*\* for example ]** We will have  
**One OBDA File ( Mapping ) per Docker Container ?? )**

**Docker** : software container platform



# Thank You!

