

Gestion et valorisation sémantiques de données de biodiversité et d'études d'écosystèmes dans l'infrastructure ANAEE-France

Damien Maurice, INRA – AnaEE-France
Christian Pichot, INRA – AnaEE-France

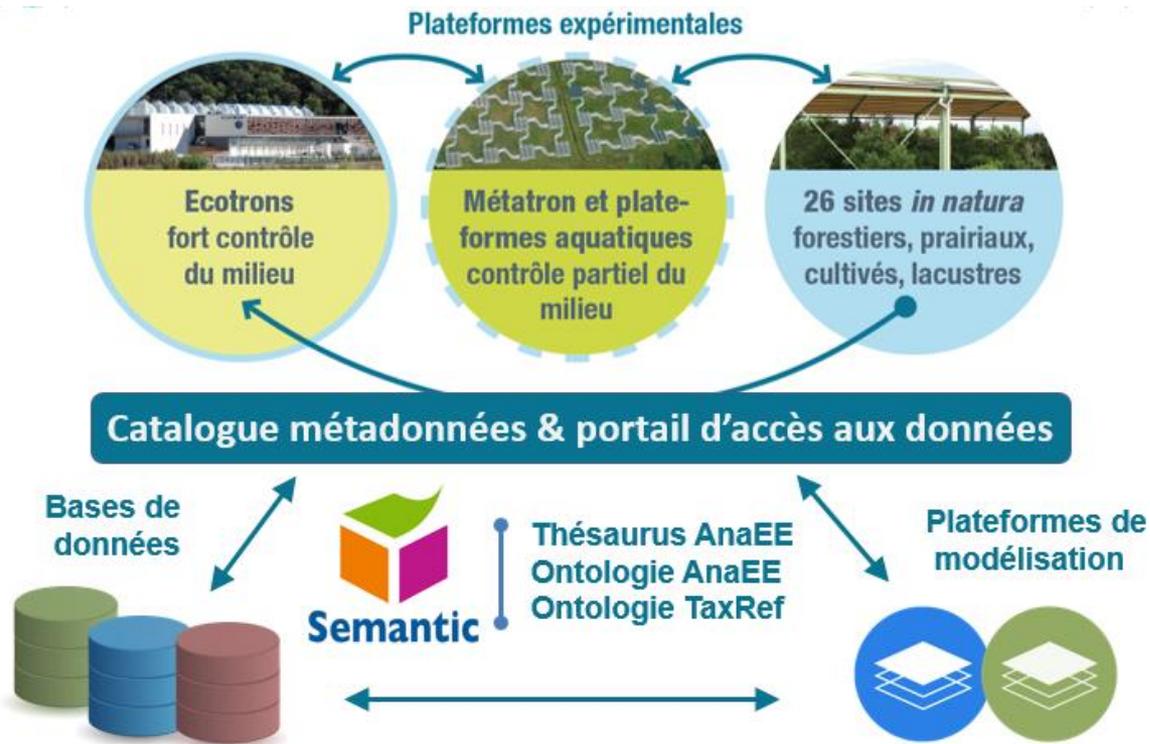
A. Chanzy, E. Aivayan, N. Beudez, C. Callou, P. Clastre, M. El-Hamadry,
L. Greiveldinger, B. Jaillet, F. Lafolie, A. Léturgie, A. Maire, C. Martin, D. Maurice,
N. Moitrier, G. Monet, H. Raynal, A. Schellenberger, R. Yahiaoui

École thématique e-Envir - Gif-sur-Yvette

28-31 octobre 2019



Contexte AnAEE



→ Mobilisation des technologies du web sémantique pour la gestion et l'exploitation de la connaissance sur les données, par les machines et un peu les humains ...

Un Système d'Information distribué ... et une approche sémantique

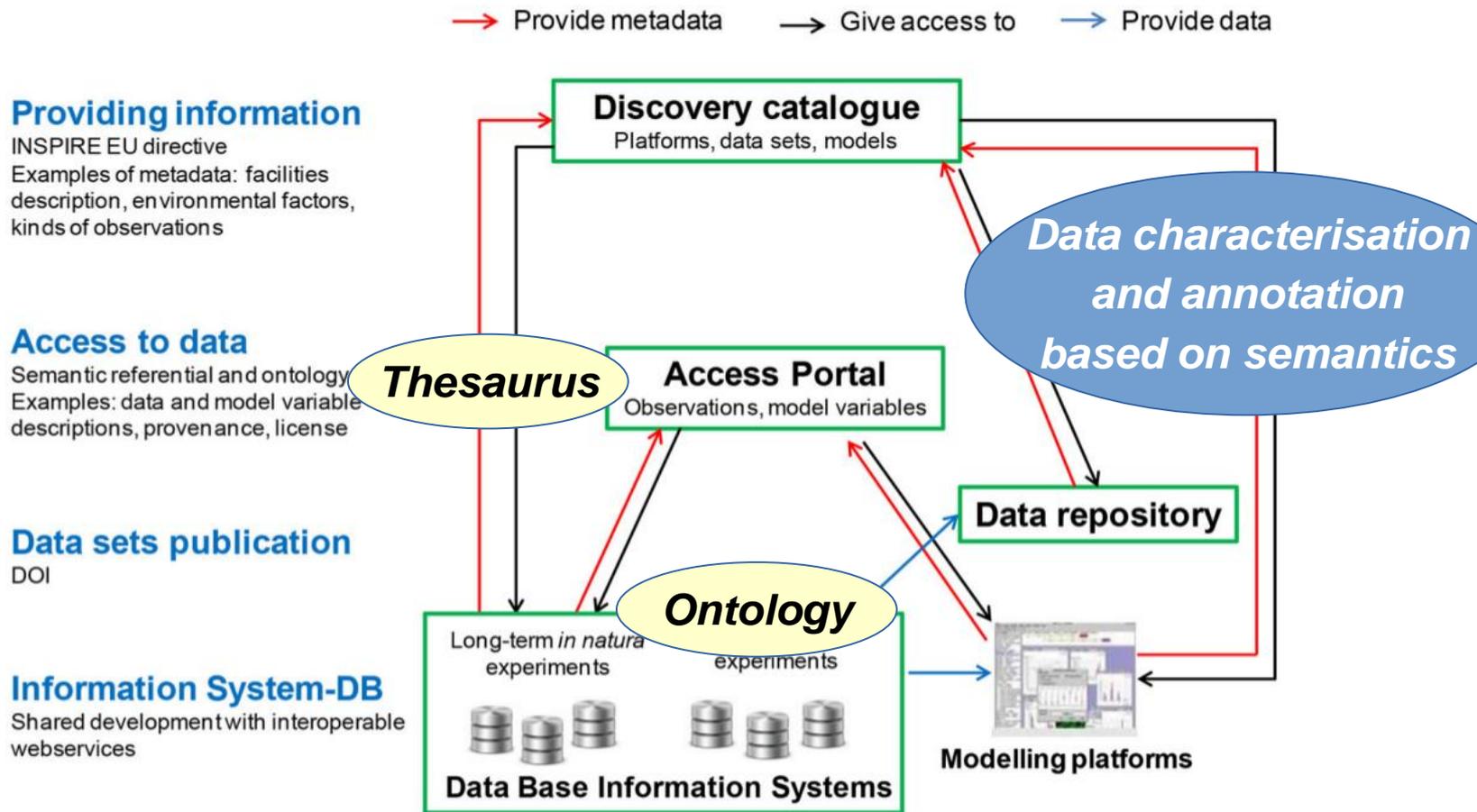
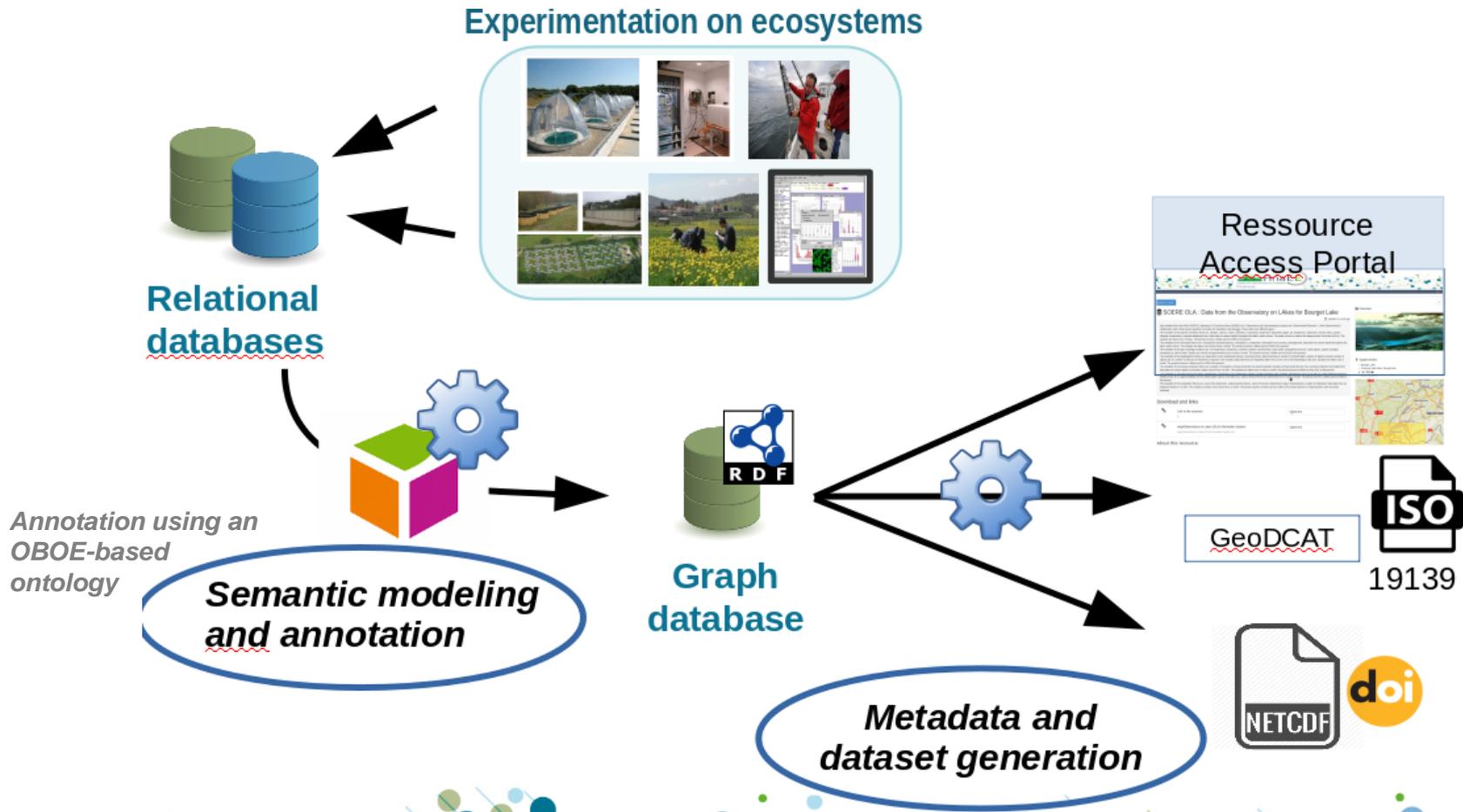


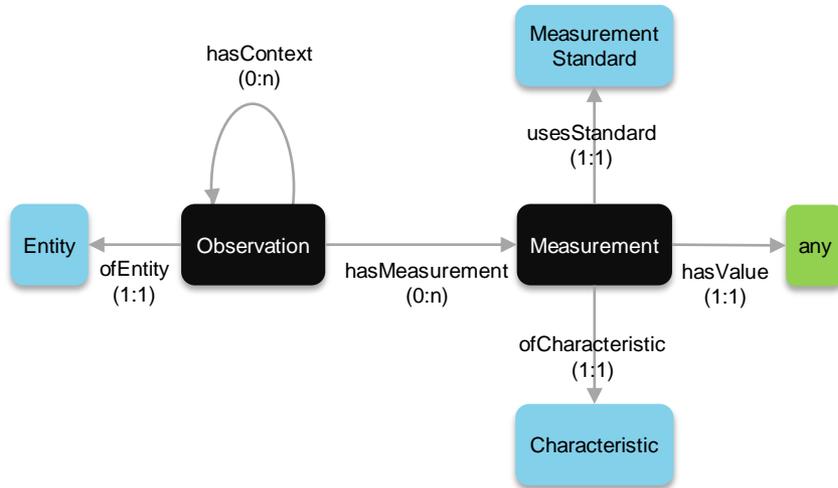
FIGURE 5 | The distributed architecture of the AnaEE-F information system includes a discovery catalog to access metadata information about platforms, datasets, or models, a portal to access metadata about observations or model variables including a semantic referential and an ontology, and a data repository to store digital object identifies (DOI) of data sets from information systems of *in natura* and mesocosm experiments. Data sets from experiments are linked with model factories to enable model parameterisation or data assimilation.

Un flux de gestion des données/métadonnées

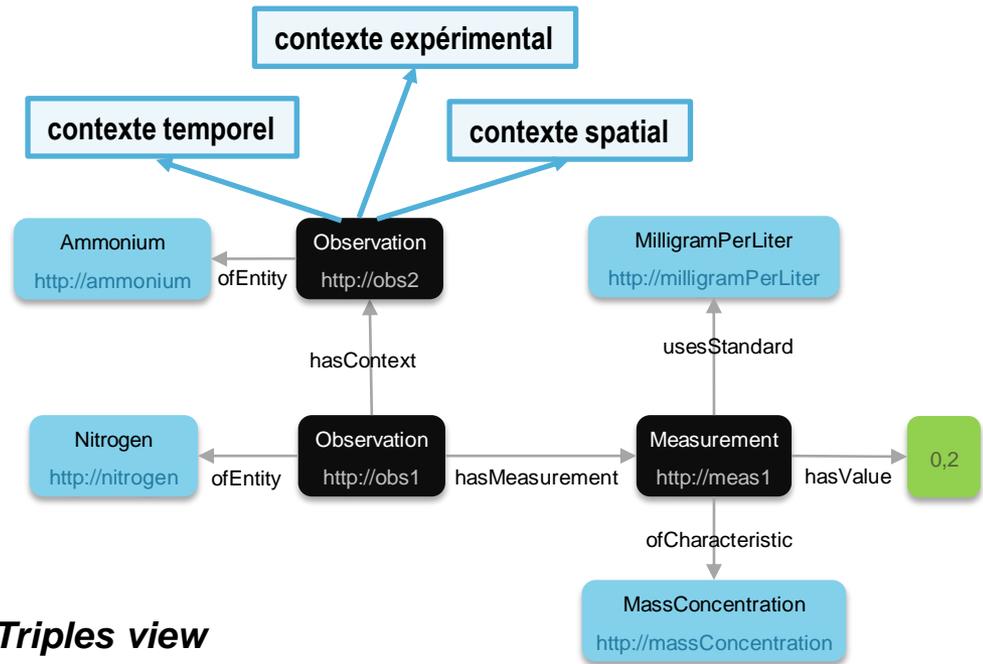


Modélisation sémantique et ontologie

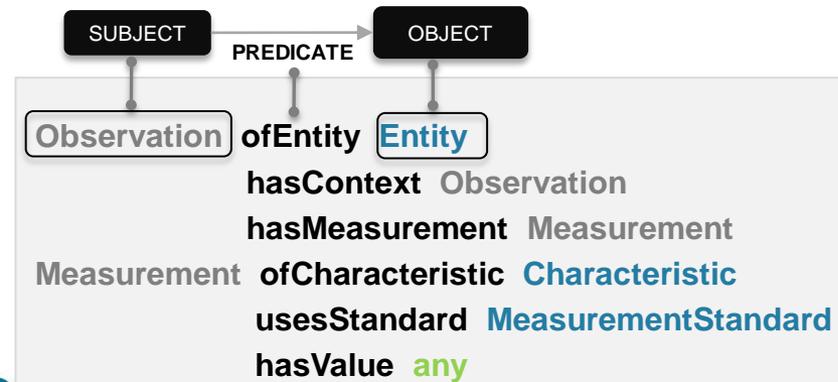
Grappe générique de l'ontologie OBOE



Exemple de graphe



Triples view

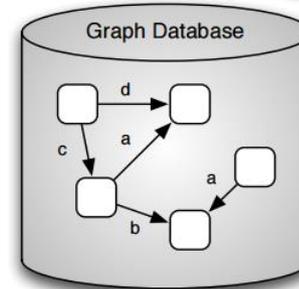
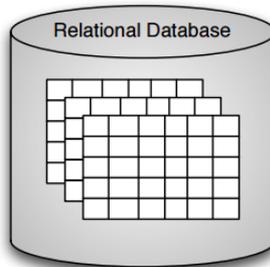
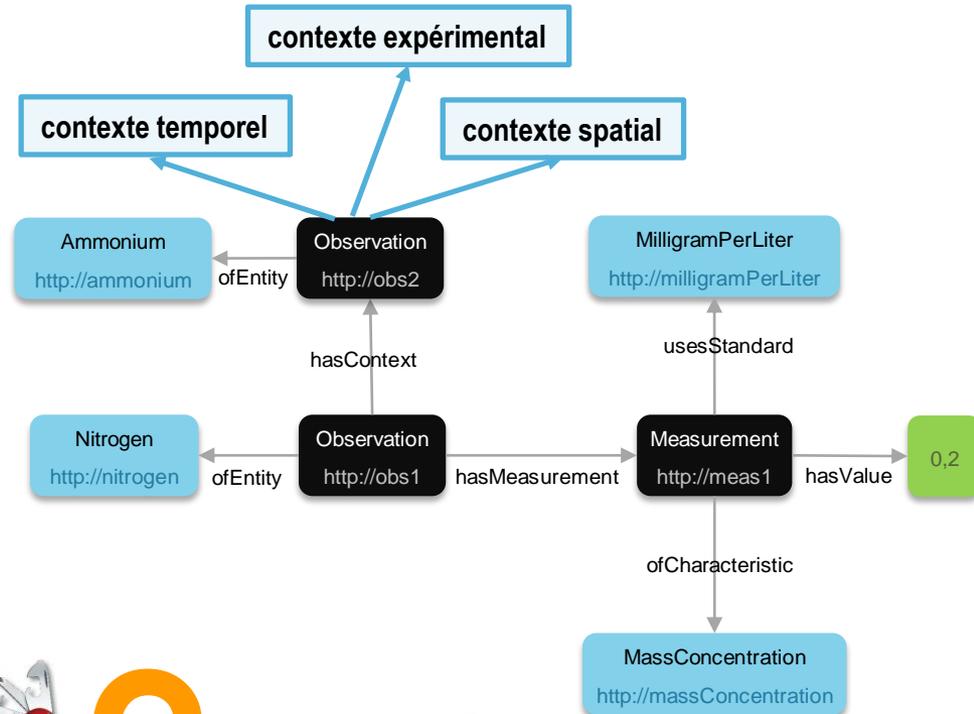
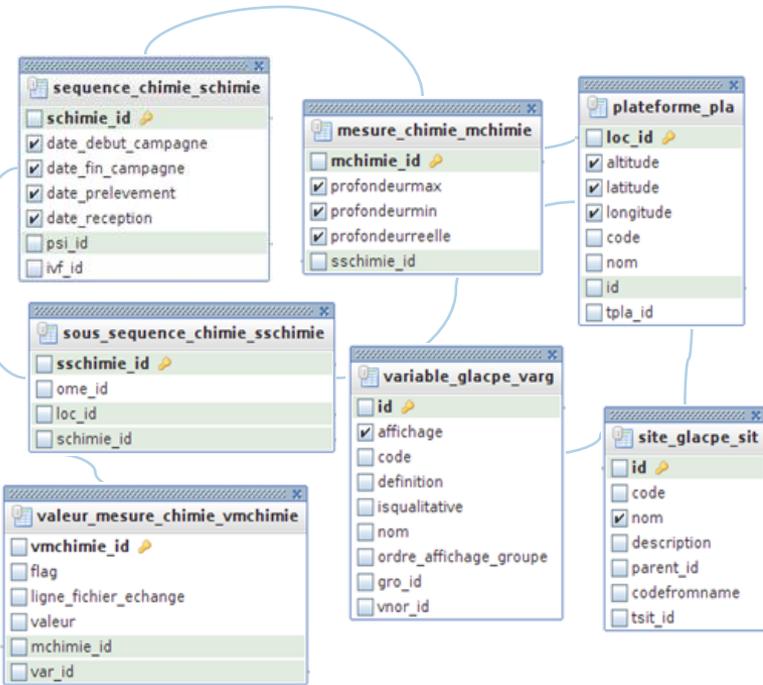


Triples view

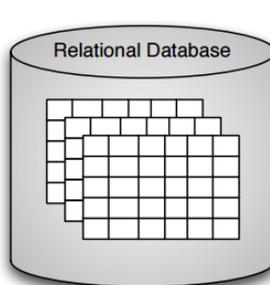
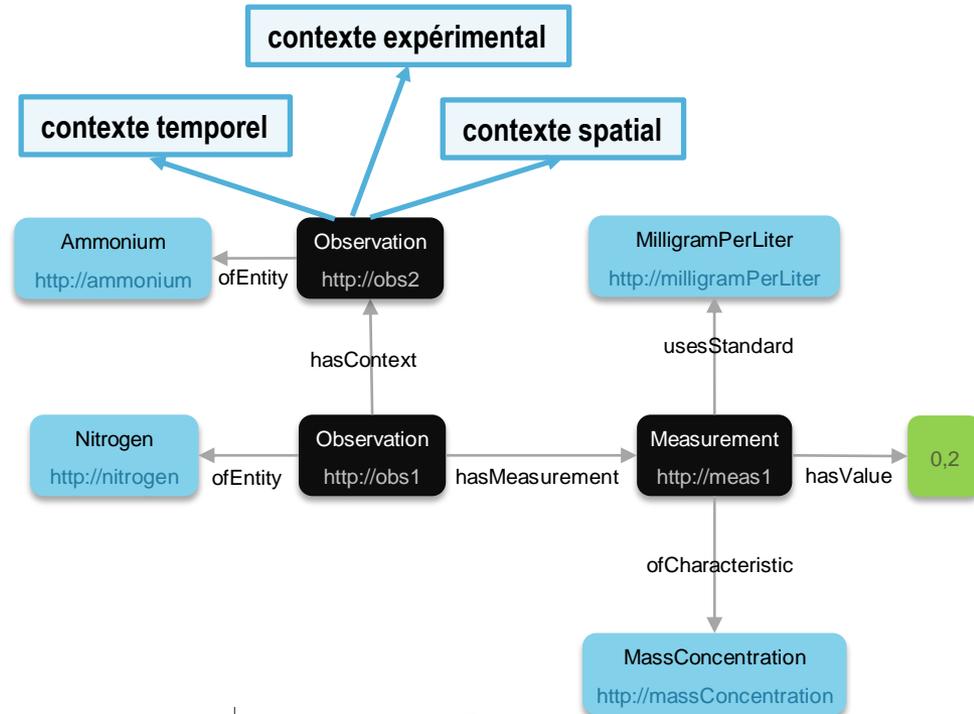
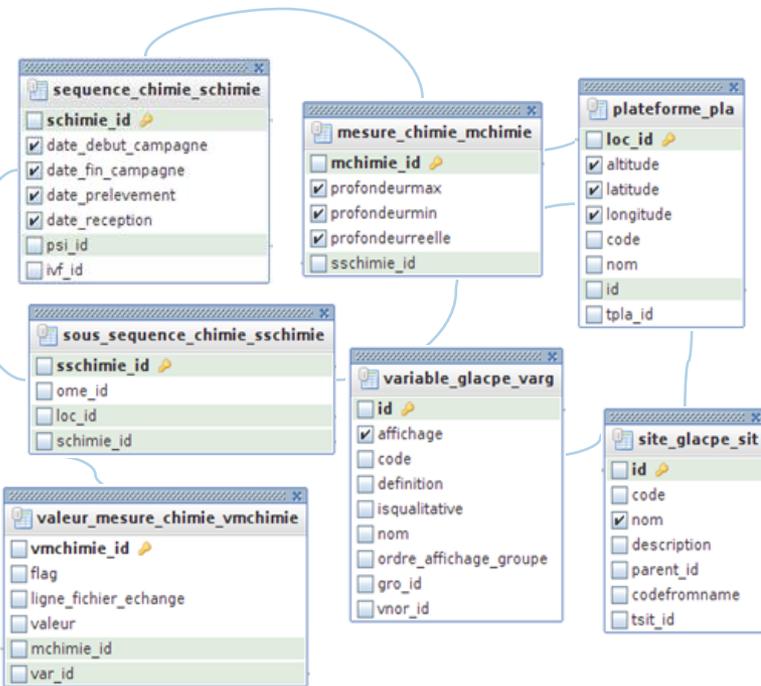
```

http://obs1 ofEntity http://nitrogen
            hasContext http://obs2
            hasMeasurement http://meas1
http://obs2 ofEntity http://ammonium
http://meas1 ofCharacteristic http://massConcentration
            usesStandard http://milliGramPerLiter
            hasValue 0,2
  
```

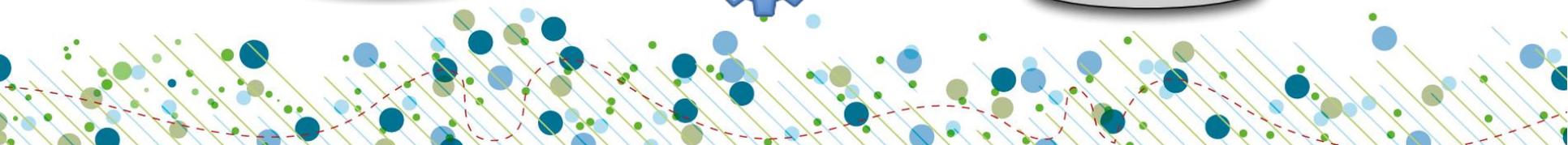
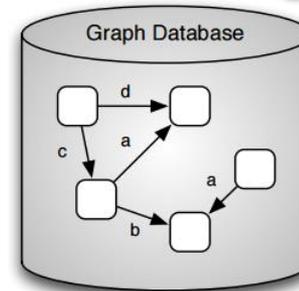
Comment passer des SI initiaux (ici BDD) au(x) graphe(s) ?



Comment passer des SI initiaux (ici BDD) au(x) graphe(s) ?

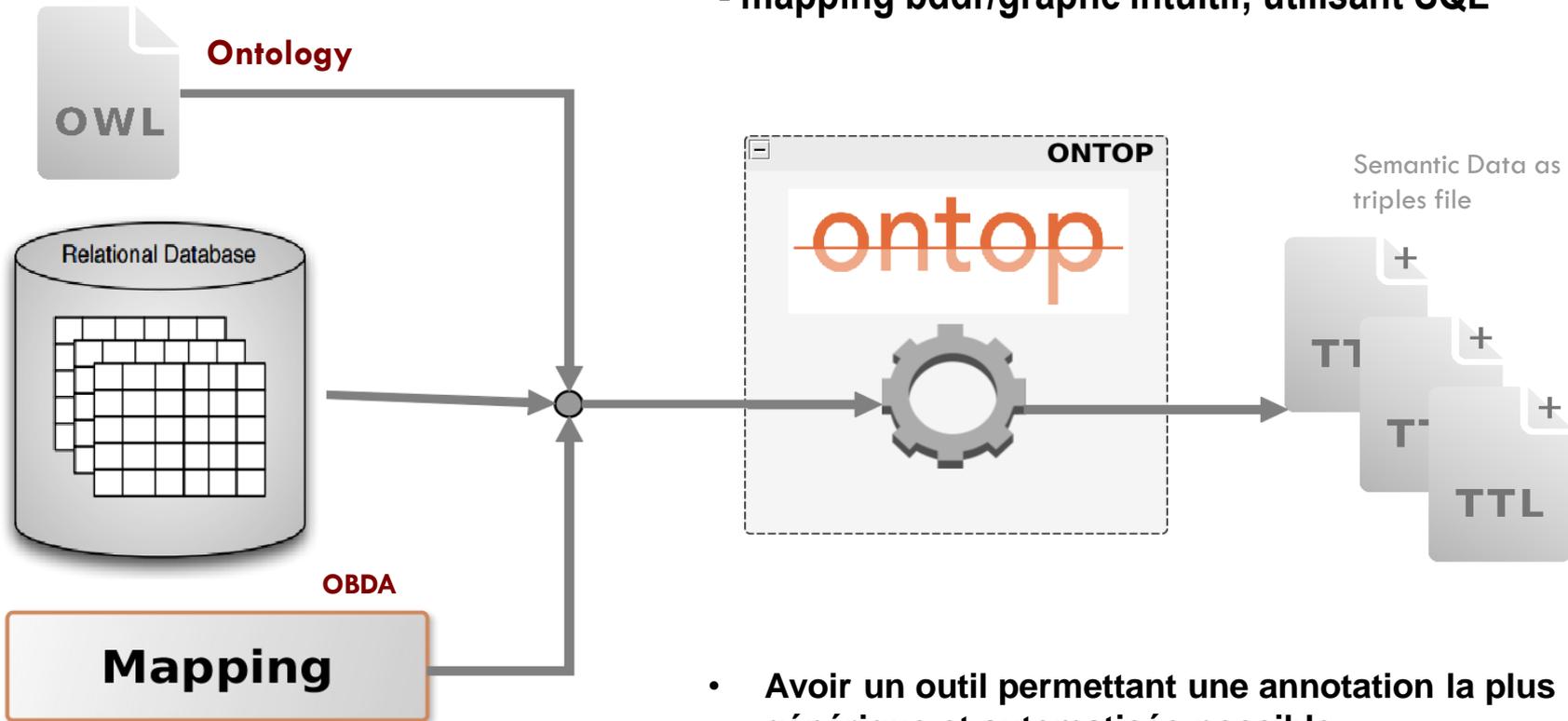


ontop



Comment passer des SI initiaux (ici BDD) au(x) graphe(s) ?

- transformation à la volée utilisant une ontologie
- mapping bddr/graphe intuitif, utilisant SQL



Fichier spécifique indiquant comment transformer les données relationnelles en graphes sémantiques

- Avoir un outil permettant une annotation la plus générique et automatisée possible
- En utilisant un outi open source et des développements spécifiques

Comment effectuer le mapping requis?

1. Modéliser les graphes (selon l'ontologie)

2. A chaque nœud d'un graphe doivent être associés :

- un URI :

→ URI fixe (ex: classe de l'ontologie)

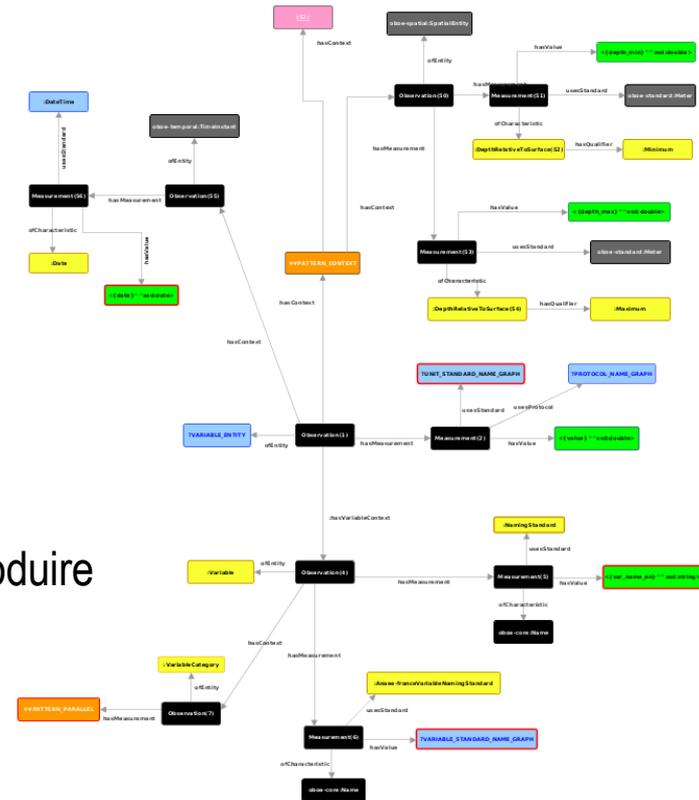
Ex: <http://anaee/massConcentration>

→ URI dynamique contenant des valeurs issues des bdd value

Ex : http://anaee/ola/observation/water{measure_id}

- une requête SQL pour renseigner les URI dynamiques et produire les triplets à la volée

`SELECT measure_id, value FROM table`



modèle d'annotation pour le pipeline de production des triplets

automatiser et générer le plus possible

1 modèle d'annotation pour n variables [350 variables déclarées dans l'ontologie]

Variable	Category(ies)	Context(s)	Entity	Characteristic	Unity
DissolvedAmmoniumNitrogenMassConcentration	PhysicalChemistry	Water, Solutes, Ammonium	Nitrogen	MassConcentration	MilligramPerLiter
CalciumMassConcentration	PhysicalChemistry	Water	Calcium	MassConcentration	MilligramPerLiter
WaterPH	PhysicalChemistry		Water	pH	pHUnit
...

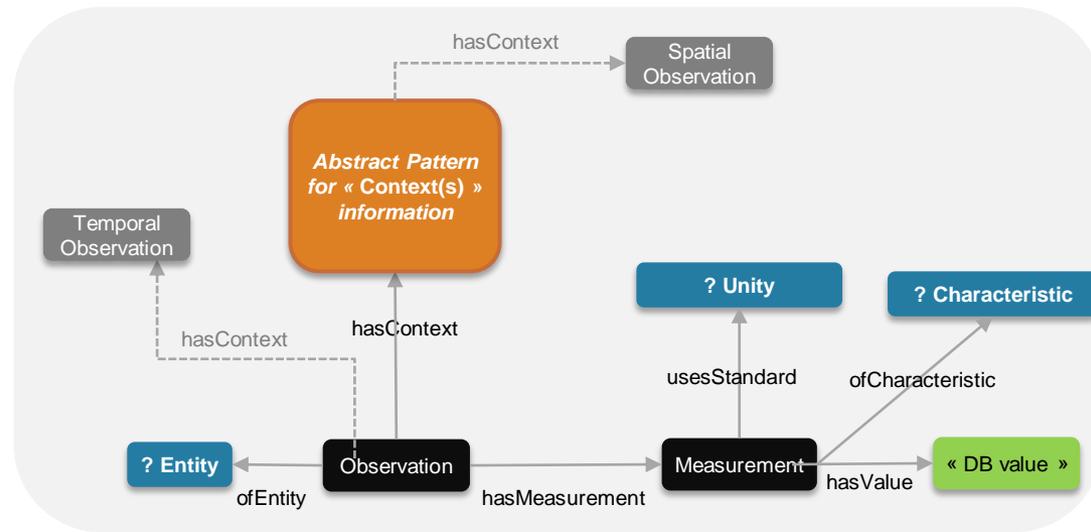
? Information

Unique information

Abstract Pattern

Multiple informations

Principe : générer automatiquement des fichiers de mapping ontop de plusieurs variables à partir d'un même modèle d'annotation (dépendant de la modélisation de la bddr).

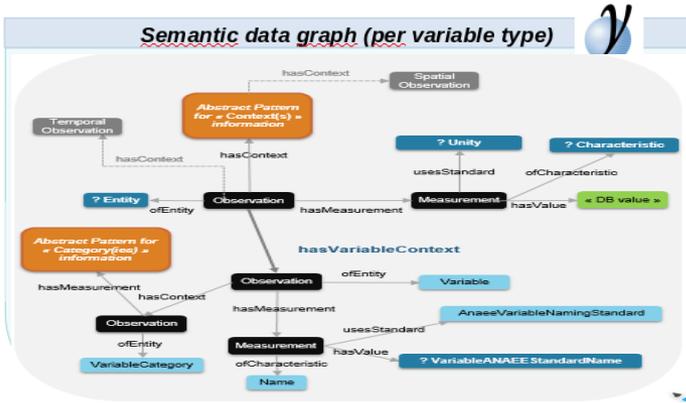


Vue générale du pipeline d'annotation sémantique

Variable semantic description



AnaEE standard	Category	Context	Entity	Characteristic	Protocol	Unit	variable DB name	DB category
Phytoplankton	Biodiversity	Water	Phytoplankton	Volume Per Volume		MicroMeter Cubed Per Millimeter	phytoplankton	biodiversité
WaterPH	Physical Chemistry		Water	pH		pHUnit	pH	physicochimie



YedGen



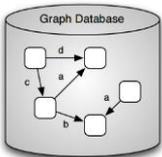
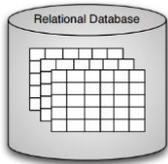
Mapping files for Ontop



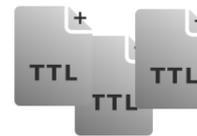
ontop



Ontology



raw data with inferred triples

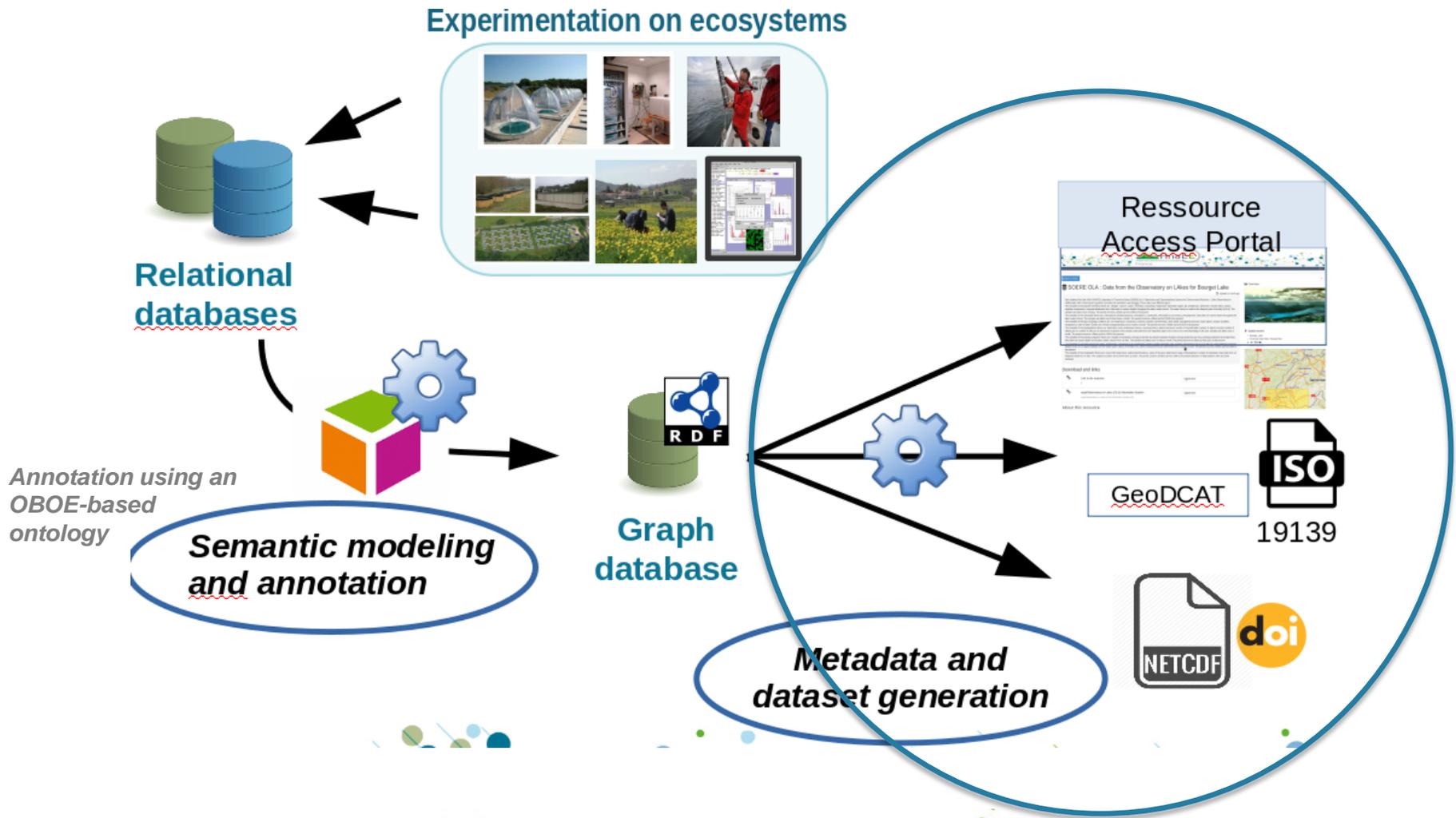


raw data

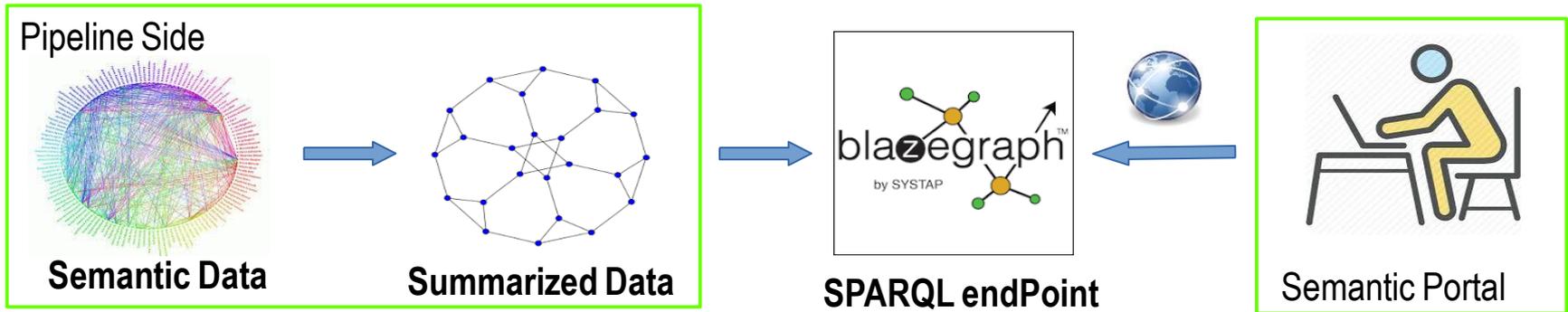


End point

Un flux de gestion des données/métadonnées



Ressource access portal



Accueil [Qui sommes-nous ?](#) [Services de l'infrastructure](#) [Espace enseignement](#) [Ressources](#)

Vous êtes ici : [Accueil](#) > [Ressources](#) > [Bases de données](#)

Ce sont toutes les données issues des plateformes d'expérimentations d'AnaEE.

Les bases de données à long terme...

Recherche

Tri par :

Type de ressources

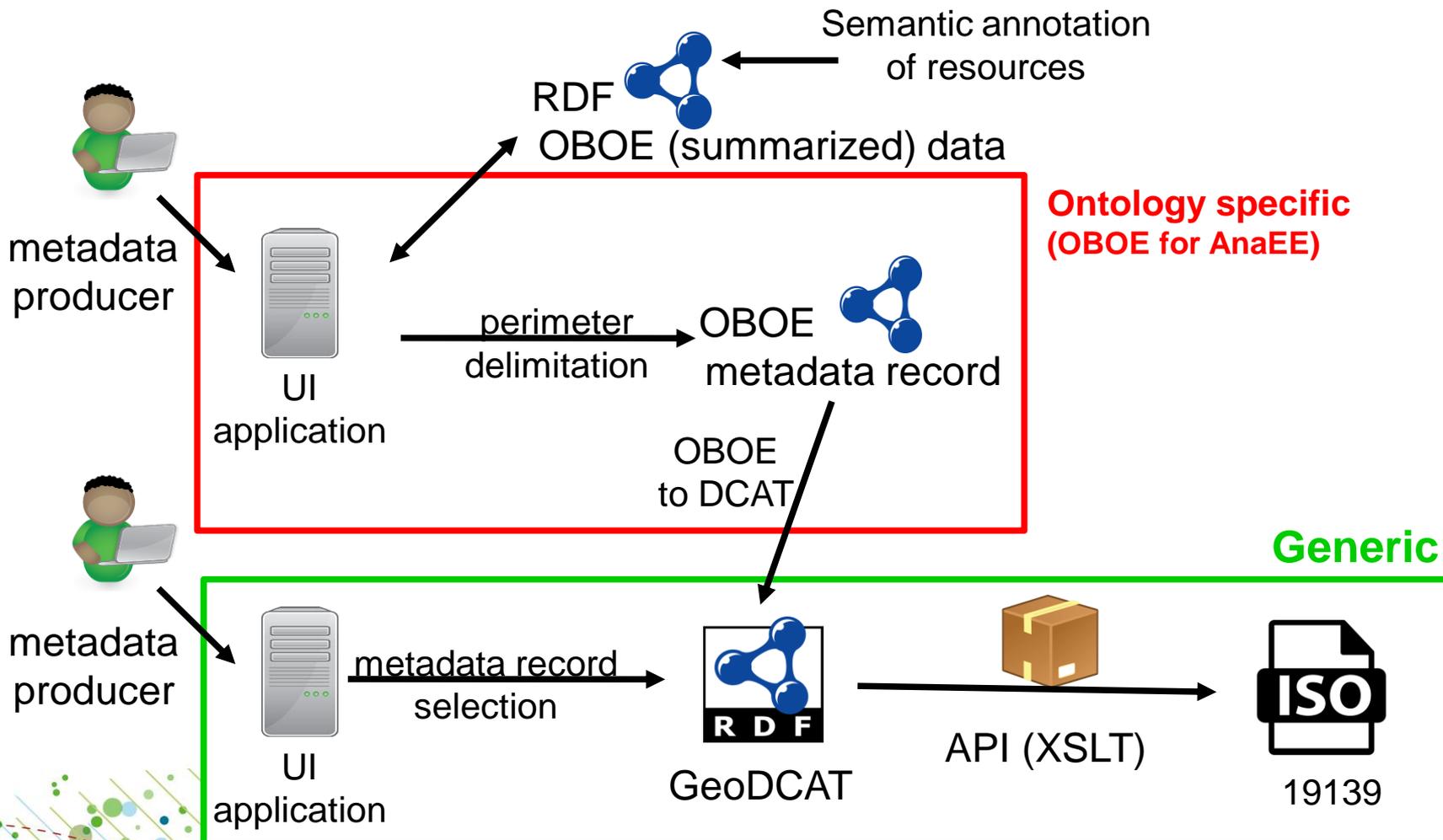
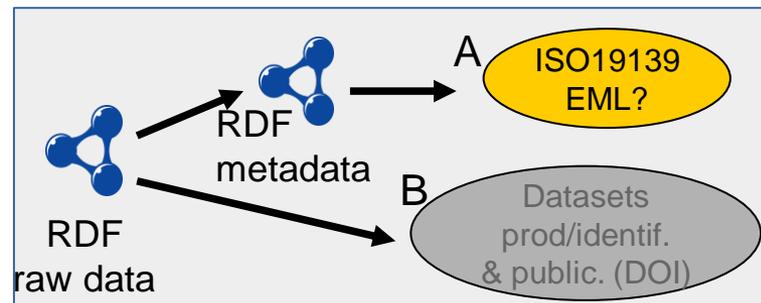
ecosystemes	especies	Année	lieu	variables	auteur
<input type="text" value="Sélectionnez une..."/>	<input type="text" value="Insectes"/>	<input type="text" value="1988"/>	<input type="text" value="Theix"/>	<input type="text" value="0"/>	<input type="text" value="0"/>

1 2 3 4 5

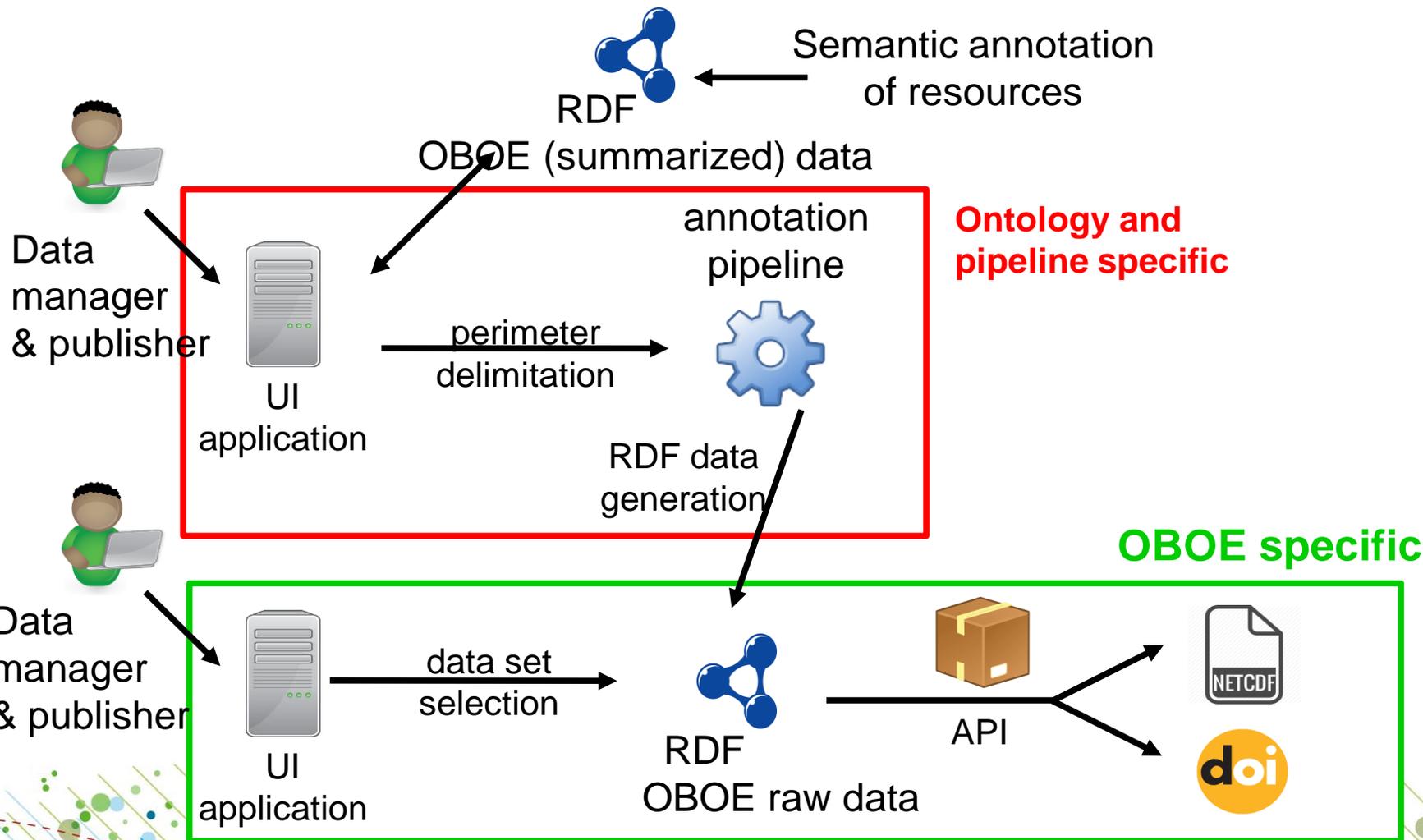
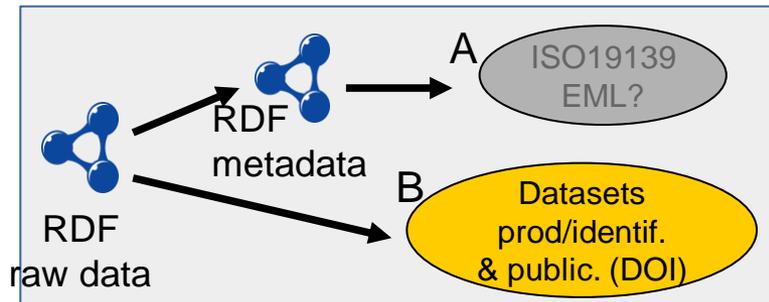
Affichage #

Objectif : Produire un nouveau graphe de données élaborées en utilisant le pipeline d'annotation et le publier dans un SPARQL Endpoint spécifique qui sera interrogé par le portail AnaEE-France

Pipeline for semantic generation of metadata (& data) sets



Pipeline for semantic generation of (metadata &) data sets



interface des services de génération de (méta)données

Homepage

Scopes creation

Scopes deletion

Data and metadata
production

Executed scopes
management

DOI and metadata
publication

Executions in
progress

Help

Select a scope

waterTempPormenazOla

Select scope

Variables	Variable categories	Sites	Infrastructures	Ecosystems	Year	Values number	Add to the scope
WaterTemperature	*	Pormenaz	SoereOla	*	*	2548	+

Create a new selection

Variables

Variables
DissolvedMagnesiumMassConcentration
DissolvedNitrateMassConcentration
DissolvedNitriteMassConcentration
DissolvedOrganicCarbonMassConcentrati..
DissolvedOrthoPhosphorusMassConcent..
DissolvedOxygenMassConcentration
DissolvedPotassiumMassConcentration

Apply filters

Variable categories

Classes de variables

Infrastructures

Infrastructures

Ecosystems

Ecosystèmes

Variables

Variable categories

Sites

Infrastructures

Ecosystems

Year

Values nb

0

Validate selection

Reset selection

After applying the filters

Create a new scope

Name

Save scope

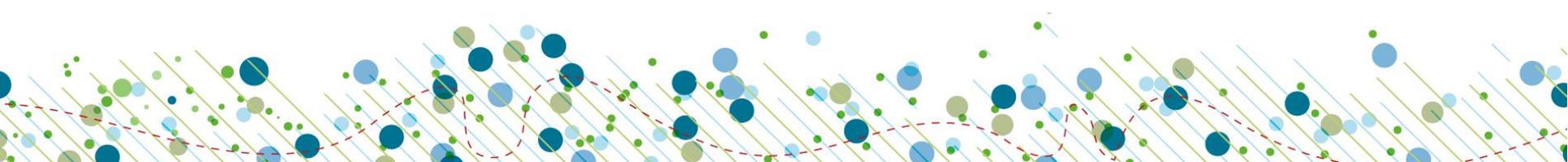
Bilan

LES PLUS

- technologies standards
- interopérabilité native
- FAIR compatible
- approche données et/ou métadonnées
- réutilisation de référentiels existants
- rapproche scientifiques et informaticiens
- généricité des pipelines (=> portefeuille de services d'ENVRI-plus)

LES MOINS

- beaucoup de nouvelles compétences à acquérir + outils
- gestion du volume des triplets



ETAT D'AVANCEMENT / PERSPECTIVES

Référentiels

- publication ontologie AnaEE et v2 du thésaurus AnaEE
- alignements avec d'autres référentiels

Pipeline d'annotation sémantique

- consolidation des développements
- tests de déploiement et performances
- déploiement sur d'autres SI d'AnaEE...autres infra de recherche

Pipelines de génération de données et métadonnées

- poursuite des développements
- prise en charge d'autres formats de sortie que NetCDF ?
- prise en charge du format EML en + de l'ISO19115 ?

