

Séminaire rdatadev
5 mars 2024

Intégration sémantique de données de biodiversité et construction d'un graphe de connaissances en écologie

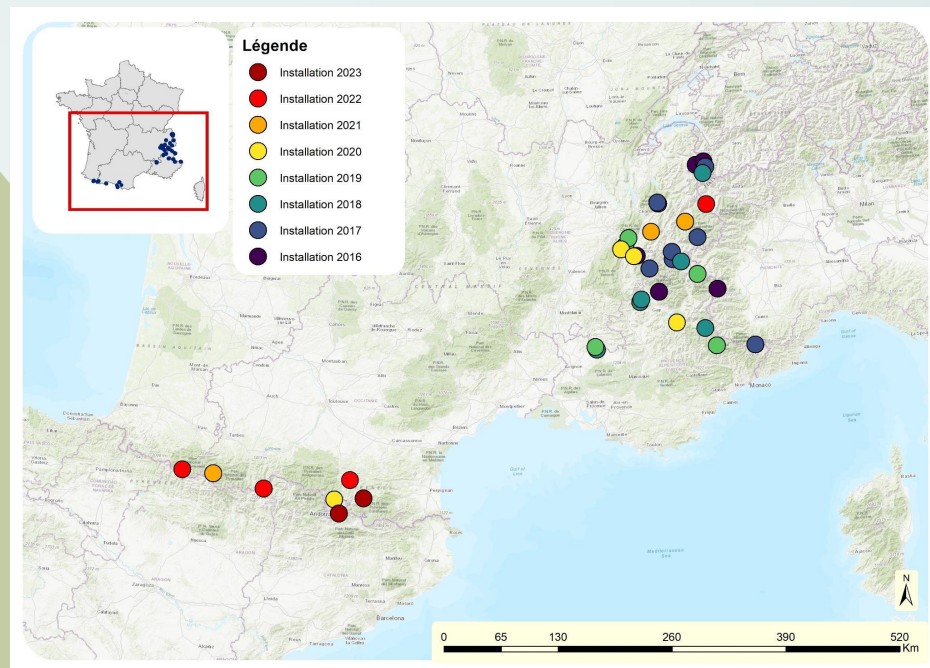
Nicolas Le Guillarme

Laboratoire d'Écologie Alpine, Grenoble, France

L'observatoire ORCHAMP

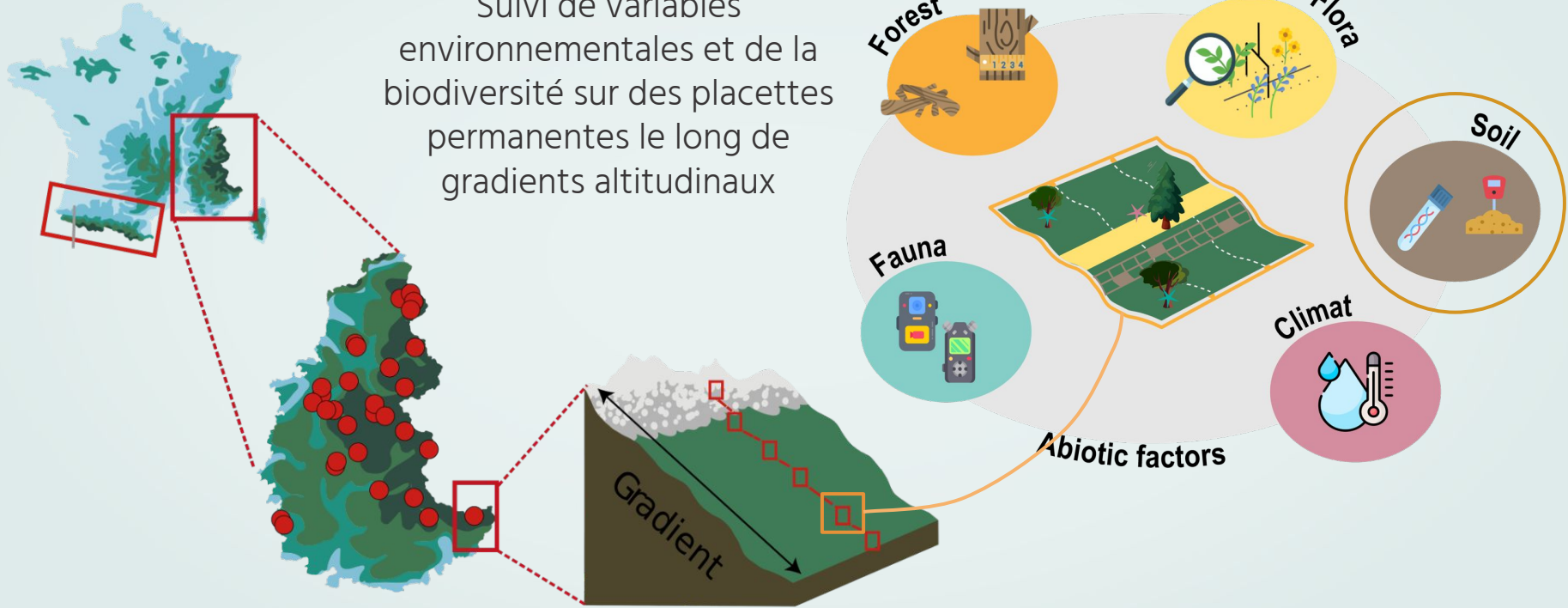


Observer, étudier et modéliser la biodiversité des territoires de montagne pour **comprendre et prédire** les répercussions des changements environnementaux sur la biodiversité et le fonctionnement des écosystèmes dans les Alpes françaises et les Pyrénées



L'observatoire ORCHAMP

Suivi de variables
environnementales et de la
biodiversité sur des placettes
permanentes le long de
gradients altitudinaux



La biodiversité du sol
représente près de 26 % des
espèces vivantes connues
de la planète

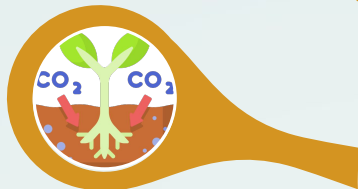
(Estimation en 2021)



Recyclage de la matière
organique, croissance
des plantes



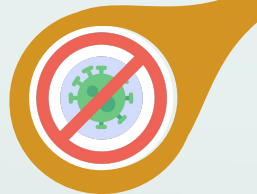
Stockage du carbone,
régulation du climat



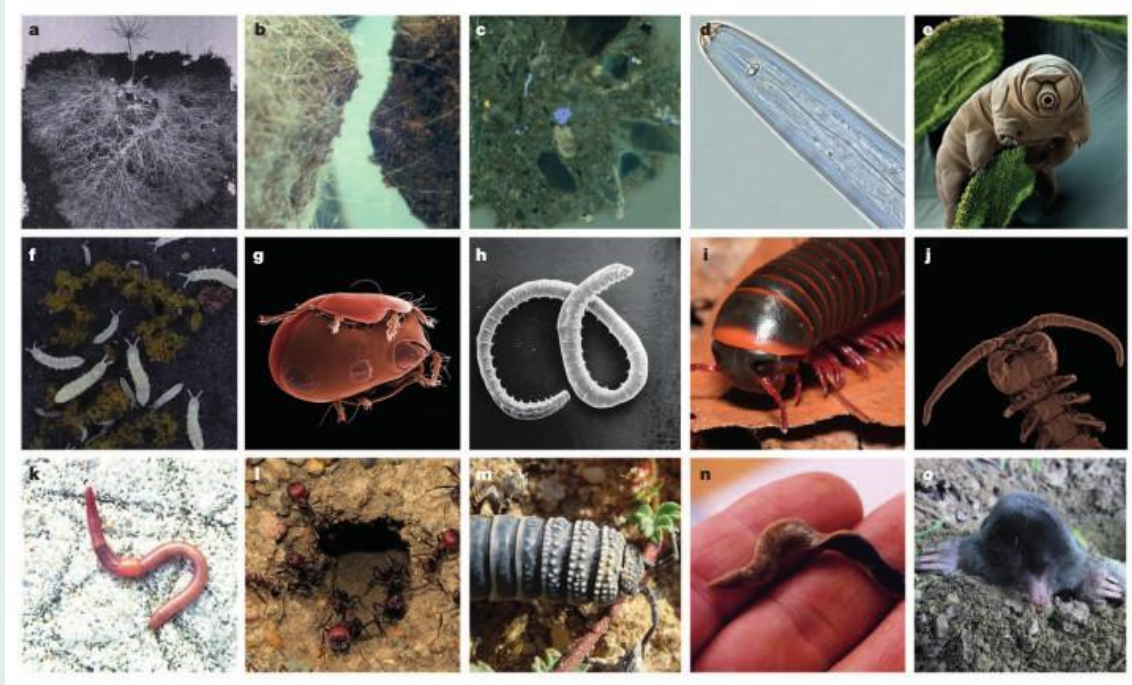
Régulation du cycle
de l'eau, purification



Production alimentaire,
santé humaine

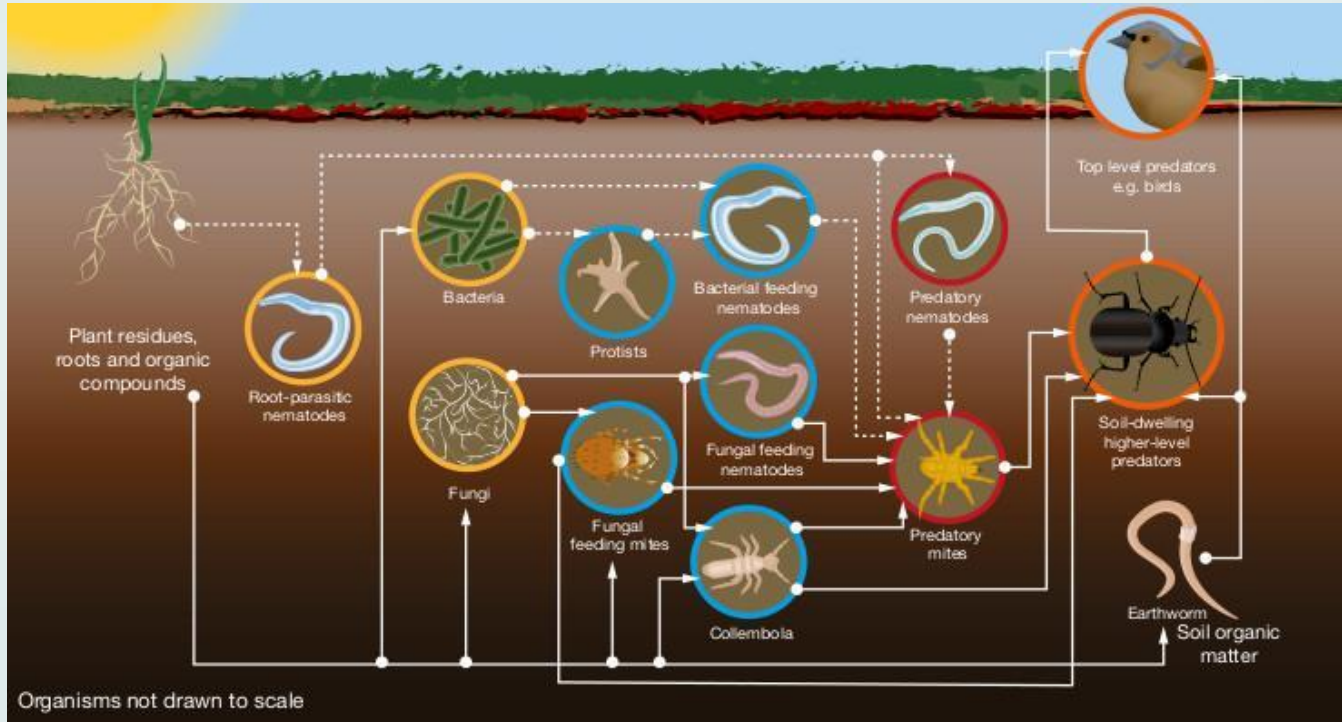


Un écosystème multitrophique complexe

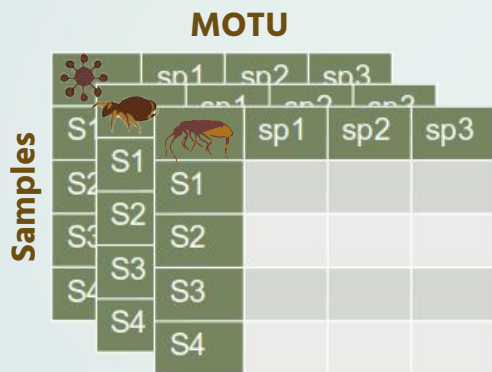


Bardgett & van der Putten, *Nature* (2014)

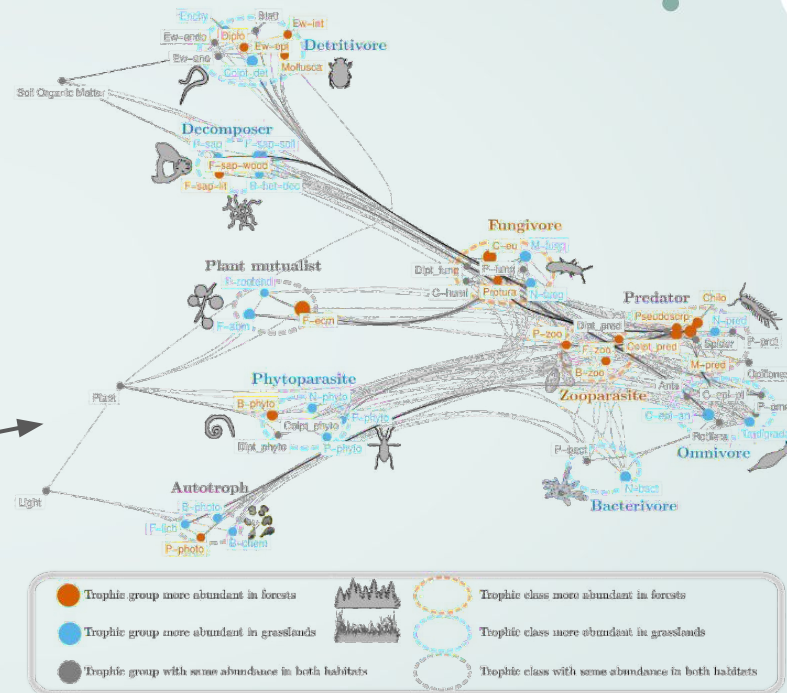
Un écosystème multitrophique complexe



Modélisation des réseaux trophiques du sol



Bases de données trophiques



Des sources d'informations distribuées et hétérogènes

Base de traits fonctionnels pour la fonge mondiale téléchargeable au format XLSX



Base de traits pour les invertébrés du sol téléchargeable au format CSV

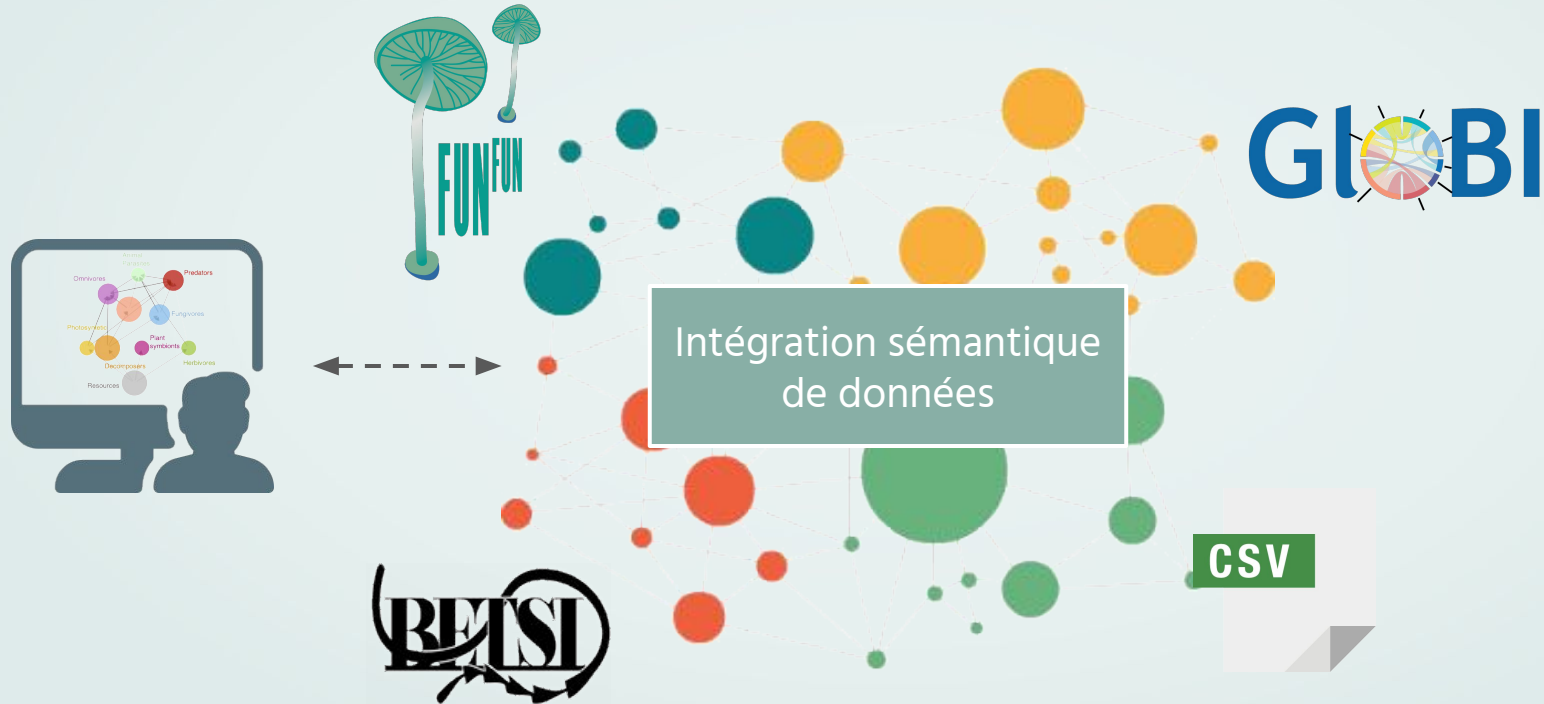


Base de données d'interactions écologiques en ligne accessible par une API

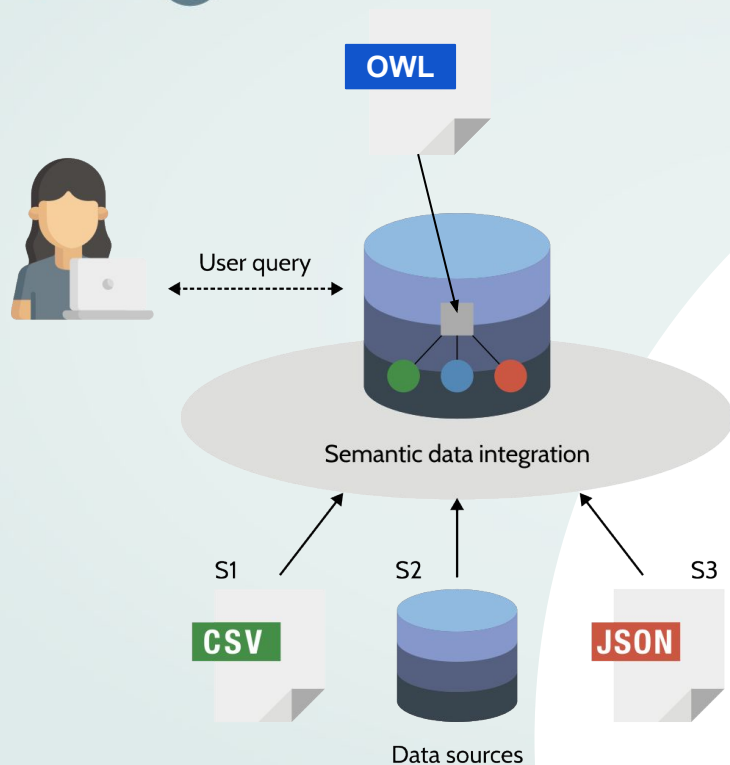
CSV

De nombreuses bases de données collectées et stockées en local par différentes équipes (*long-tail data*)

Des sources d'informations distribuées et hétérogènes



Intégration sémantique de données



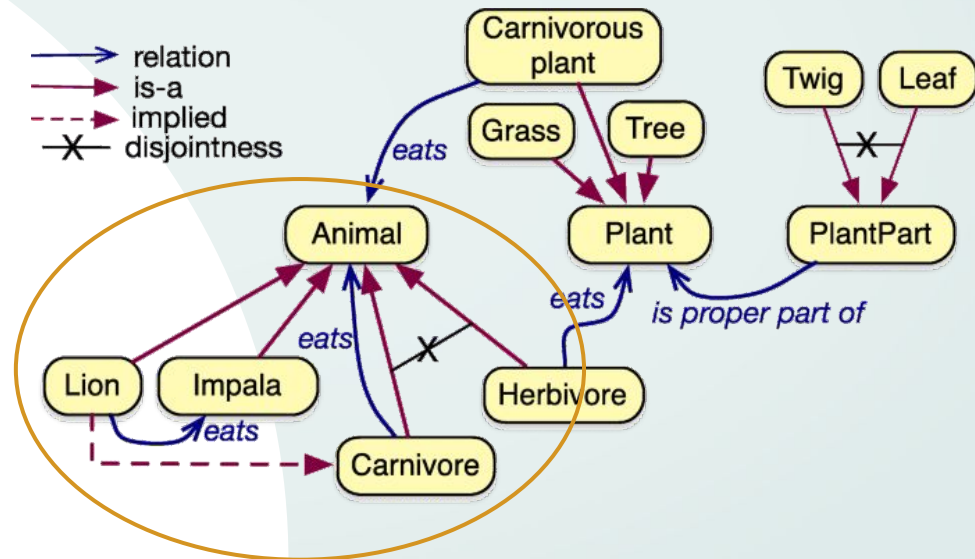
Intégration → combinaison de données issues de **sources hétérogènes** au sein d'une base de connaissance unique.

Sémantique → utilisation d'une **ontologie** pour réconcilier les différences dans l'interprétation de la "signification" des données.

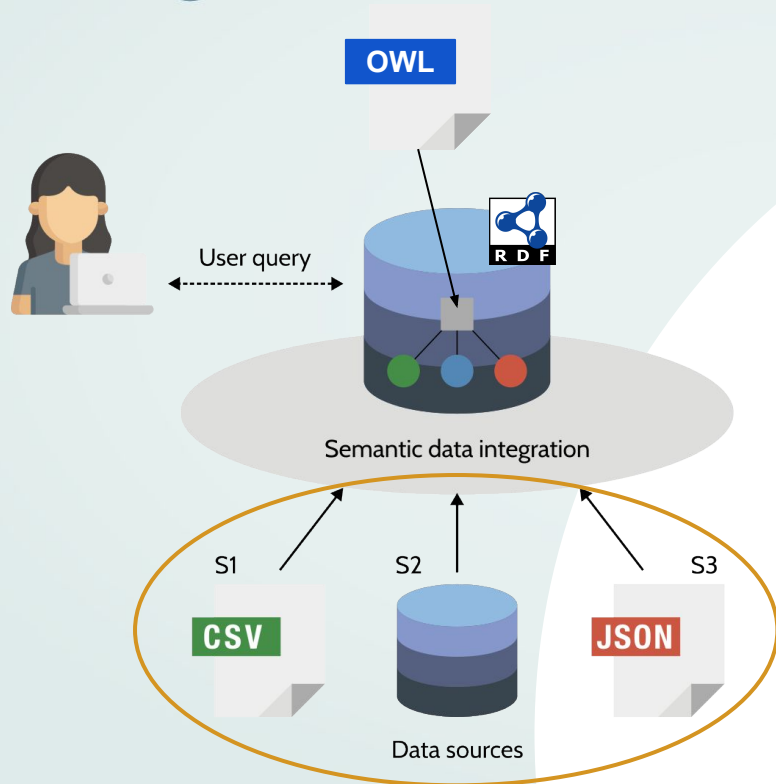
Ontologie pour la représentation des connaissances

Ontologie = un **modèle formel de la connaissance** dans un domaine qui décrit les concepts et leurs relations à l'aide de la logique mathématique.

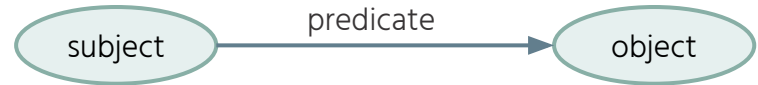
De nouvelles connaissances (implicites) peuvent être dérivées grâce au **raisonnement automatique**.



Le format RDF

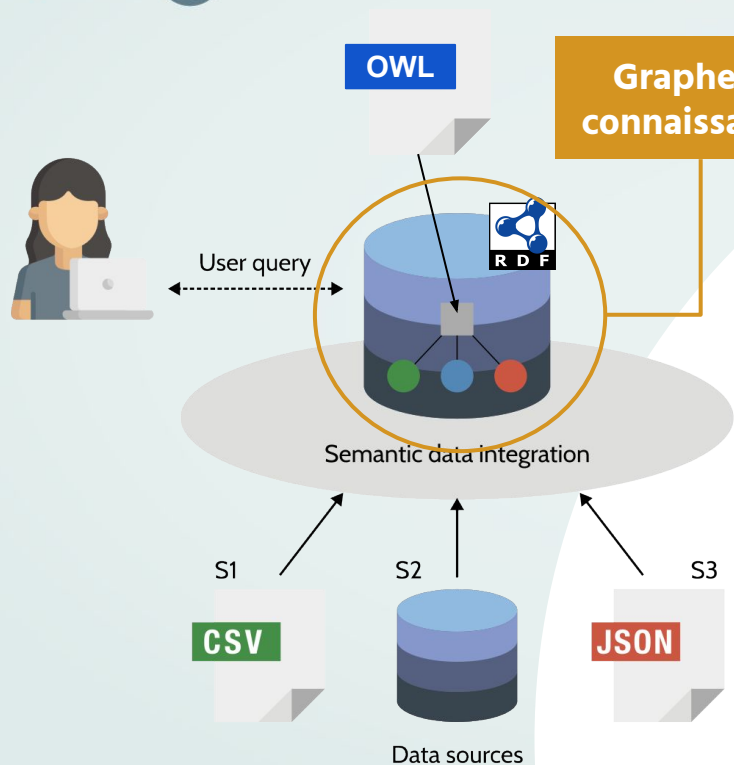


Les hétérogénéités schématiques (structurelles) sont résolues en transformant les données dans un format commun : **le format RDF** (Resource Description Framework).



Le langage de représentation des connaissances OWL est construit sur le modèle de données RDF.

Graphe de connaissance



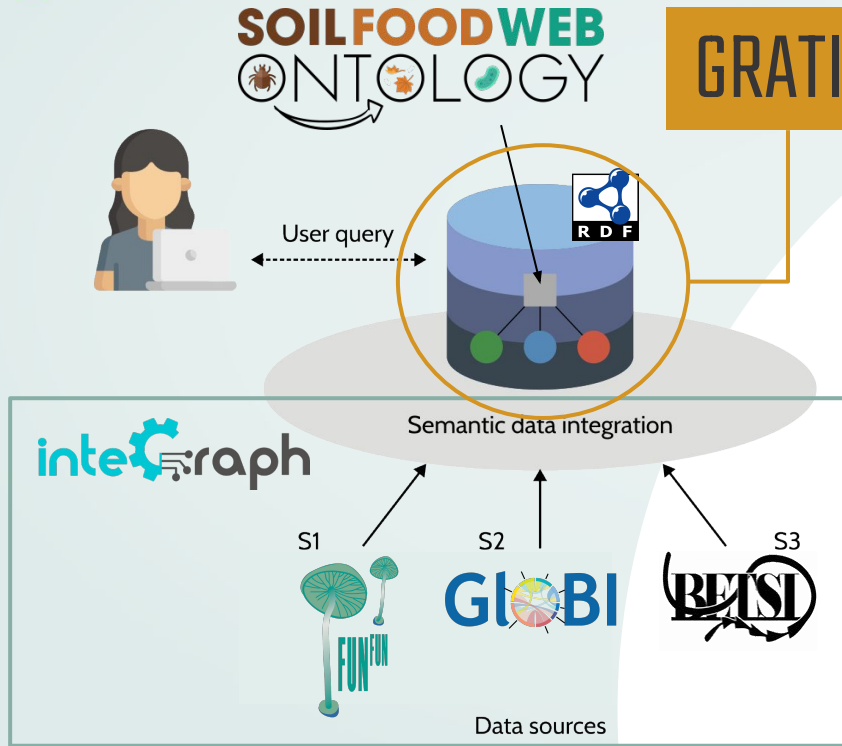
Graphe de connaissance = une base de données RDF dont la sémantique est contrôlée par une ontologie.

“A knowledge graph acquires and integrates information into an ontology and applies a reasoner to derive new knowledge.”

[Ehrlinger and Wöß, 2016]

Intégration sémantique de données = la construction d'un graphe de connaissances à partir de sources de données hétérogènes.

Construction d'un graphe de connaissance en écologie trophique du sol



La **Soil Food Web Ontology** est une ontologie OWL qui décrit formellement les concepts liés aux réseaux trophiques du sol.

inteGraph est une boîte à outil permettant de générer et de contrôler l'exécution de pipelines de sémantification de données de biodiversité.

GRATIN est un graphe de connaissance fournissant un accès unifié à des informations multisources, multitaxons, multitrophiques.

SOILFOODWEB ONTOLOGY

<https://soilfoodwebontology.github.io/>

20+ contributeurs

500+ concepts

1000+ termes

160+ groupes trophiques
pour la classification de la
faune du sol et de la fonge

Un développement
continu ouvert à tous

AgroPortal Parcourir Alignements Recommandeur Annotateur Panorama Rechercher dans AgroP Connexion FR Assistance

ontologies > SFWO

Soil Food Web Ontology (SFWO) OWL [View license](#)

Date de la dernière soumission 22 septembre 2023

Summary **Classes** Properties Instances Notes Mappings Widgets Sparql All languages

Details Instances (0) Visualization Notes (0) Classes Mappings (0)

Identifiant http://purl.org/sfwo/SFWO_0000015

Nom préféré **malacophage**

Définitions A zoophage that primarily eats mollusks such as gastropods, bivalves, brachiopods and cephalopods.

Raw data

Jump to

- food resource
- material anatomical entity
- object
 - chemical entity
 - organism
 - organism or virus or viroid
 - autotroph
 - chemotroph
 - heterotroph
 - browser
 - carnivore
 - predator
 - acariphage
 - insectivore
 - malacophage**
 - nematophage

+ axiomatisation des concepts :
malacophage \equiv *organism* and
(eats some 'member of' value *Mollusca*)

Le Guillarme, N., Hedde, M., Potapov, A. M., Martínez-Muñoz, C. A., Berg, M. P., Briones, M. J., ... & Thuiller, W. (2023). The soil food web ontology: Aligning trophic groups, processes, resources, and dietary traits to support food-web research. *Ecological Informatics*, 102360.

SOILFOODWEB ONTOLOGY

<https://soilfoodwebontology.github.io/>

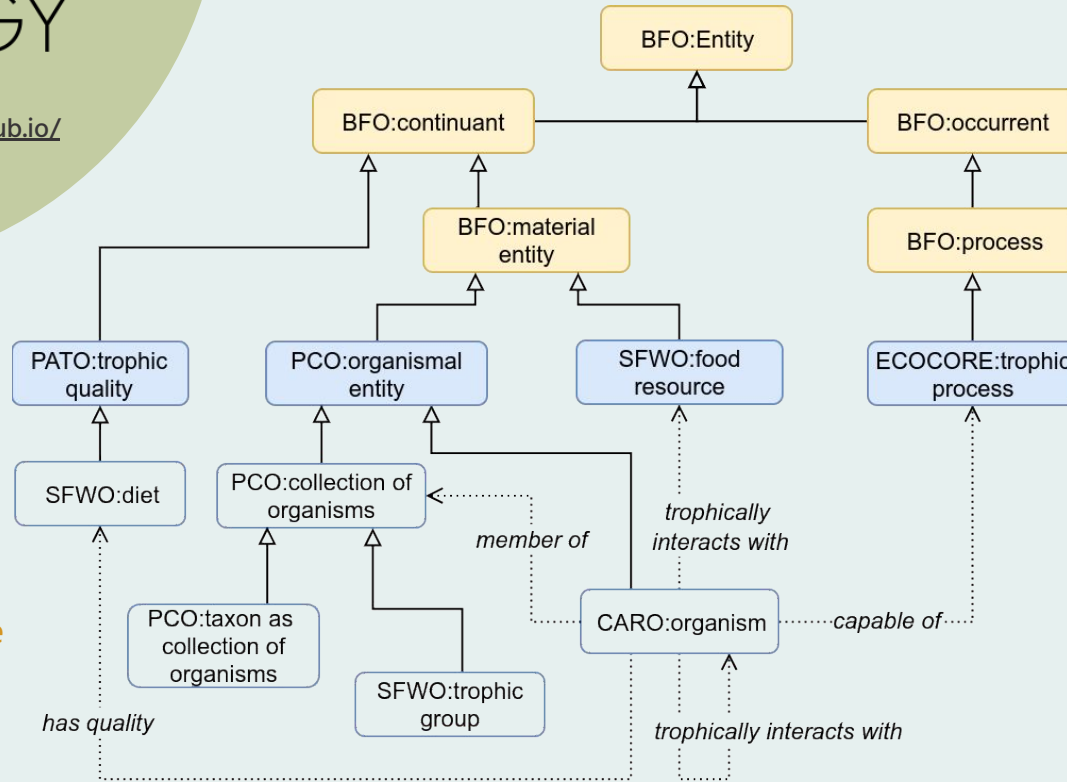
20+ contributeurs

500+ concepts

1000+ termes

160+ groupes trophiques pour la classification de la faune du sol et de la fonge

Un développement continu ouvert à tous



top-level concepts

domain-specific concepts

rdfs:subClassOf →

property - - - - ->

External ontologies

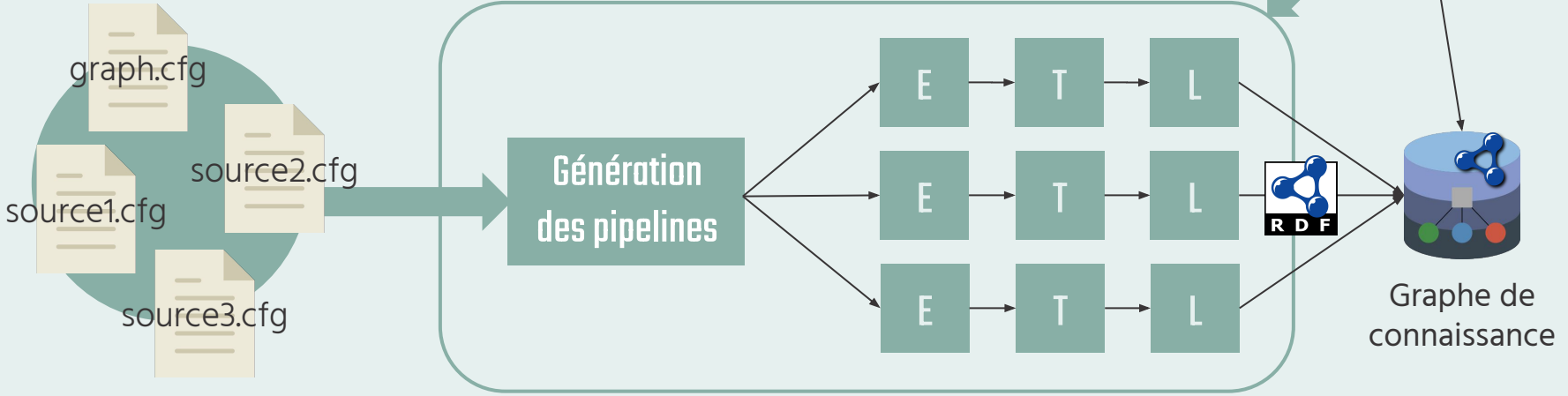
- BFO : Basic Formal Ontology
- CARO : Common Anatomy Reference Ontology
- ECOCORE : ontology of core ecological entities
- PATO : Phenotype and Trait Ontology
- PCO : Population and Community Ontology

Le Guillarme, N., Hedde, M., Potapov, A. M., Martínez-Muñoz, C. A., Berg, M. P., Briones, M. J., ... & Thuiller, W. (2023). The soil food web ontology: Aligning trophic groups, processes, resources, and dietary traits to support food-web research. *Ecological Informatics*, 102360.



<https://nlequillarme.github.io/inteGraph/>

Une boîte à outils open-source pour faciliter l'intégration et la publication de données sur la biodiversité sous la forme de graphes RDF.



Le Guillarme, N., & Thuiller, W. (2023). A practical approach to constructing a knowledge graph for soil ecological research. *European Journal of Soil Biology*, 117, 103497.

integrateGraph

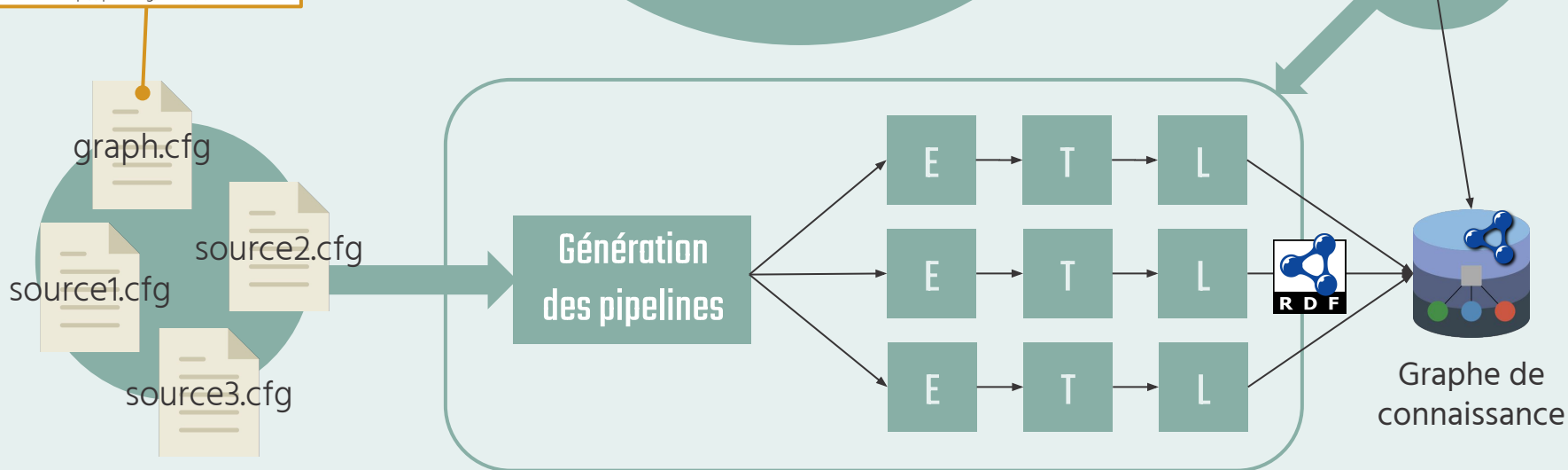
<https://nlequillarme.github.io/inteGraph/>

```
[graph]
id=http://leca.osug.fr/gratin

[sources]
dir=sources

[load]
id=graphdb
conn_type=http
host=129.88.204.79
port=7200
repository=gratin

[ontologies]
sfwo="http://purl.org/sfwo/sfwo.owl"
```

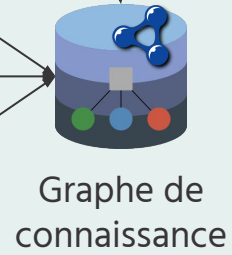
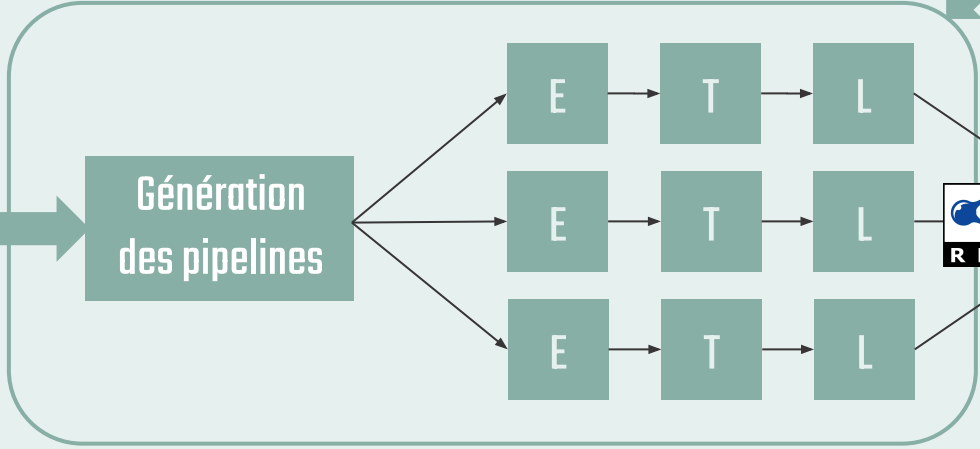
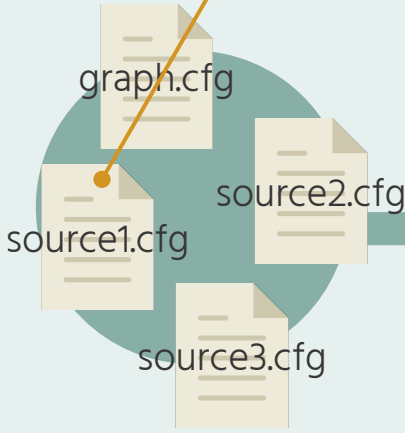


Le Guillarme, N., & Thuiller, W. (2023). A practical approach to constructing a knowledge graph for soil ecological research. *European Journal of Soil Biology*, 117, 103497.

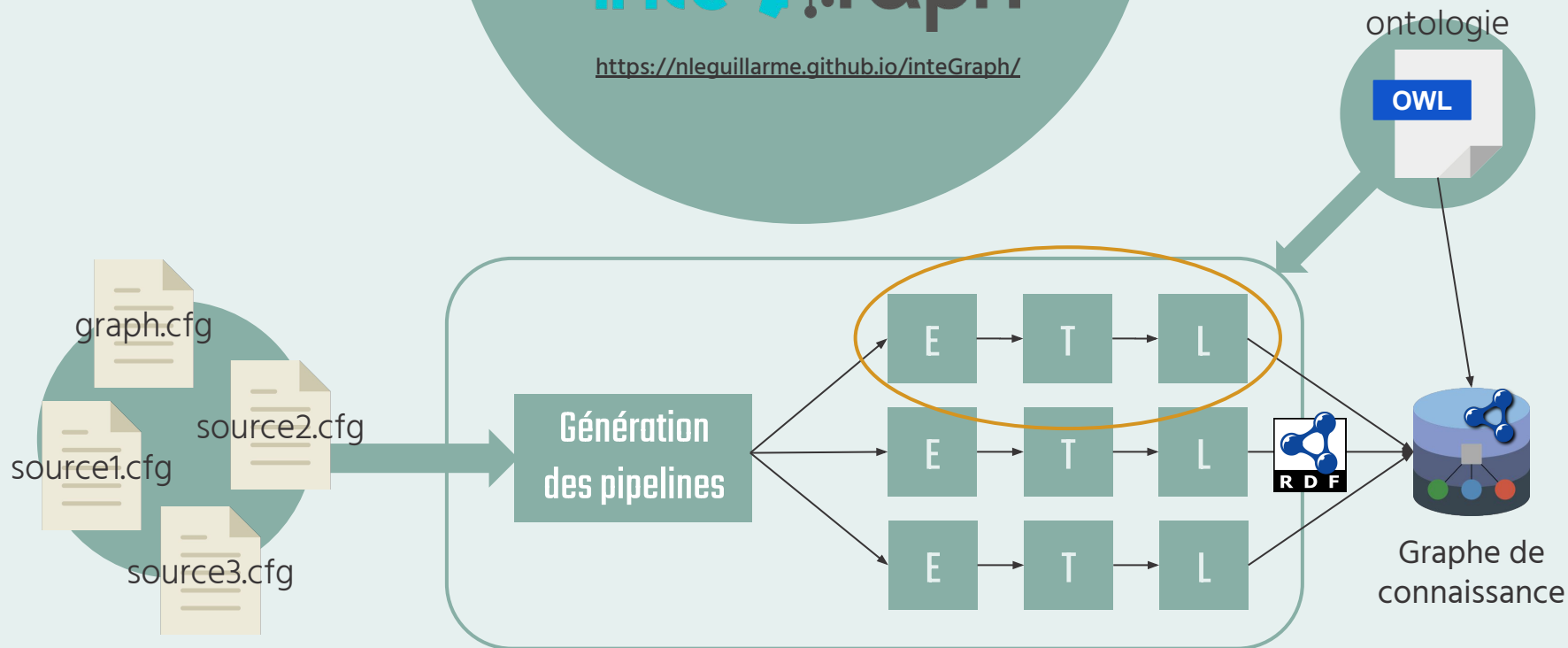


<https://nlequillarme.github.io/inteGraph/>

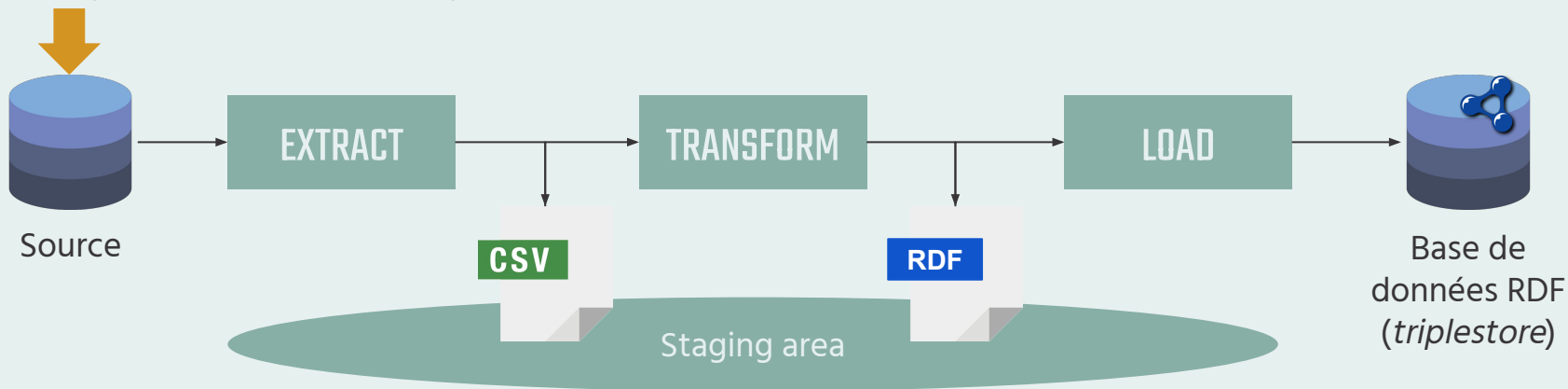
```
[source]
id=lavigne_asilidae
[extract.file]
file_path=http://www.geller-grimm.de/catalog/prey.xls
[transform.cleanser]
script="clean.py"
[transform.annotate.taxon]
label=scientificName
annotators=["TaxonAnnotator"]
[transform.annotate.trait]
label=traitValue
annotators=["SFOW", "YAMLMap"]
[transform.triplify]
mapping=mapping.xlsx
```

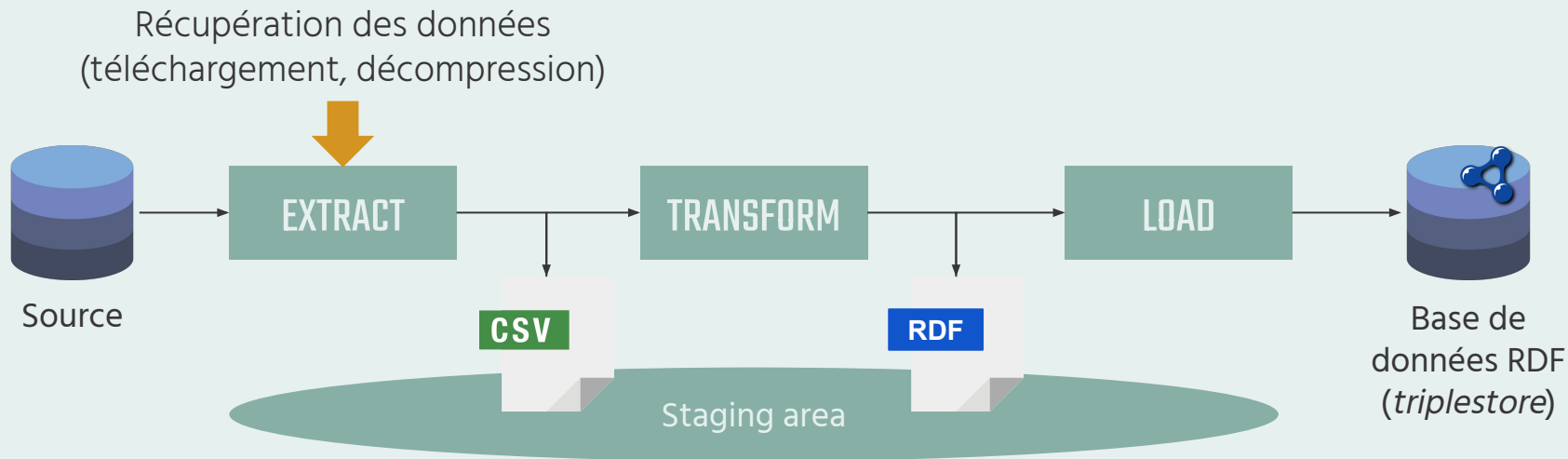


Le Guillarme, N., & Thuiller, W. (2023). A practical approach to constructing a knowledge graph for soil ecological research. *European Journal of Soil Biology*, 117, 103497.

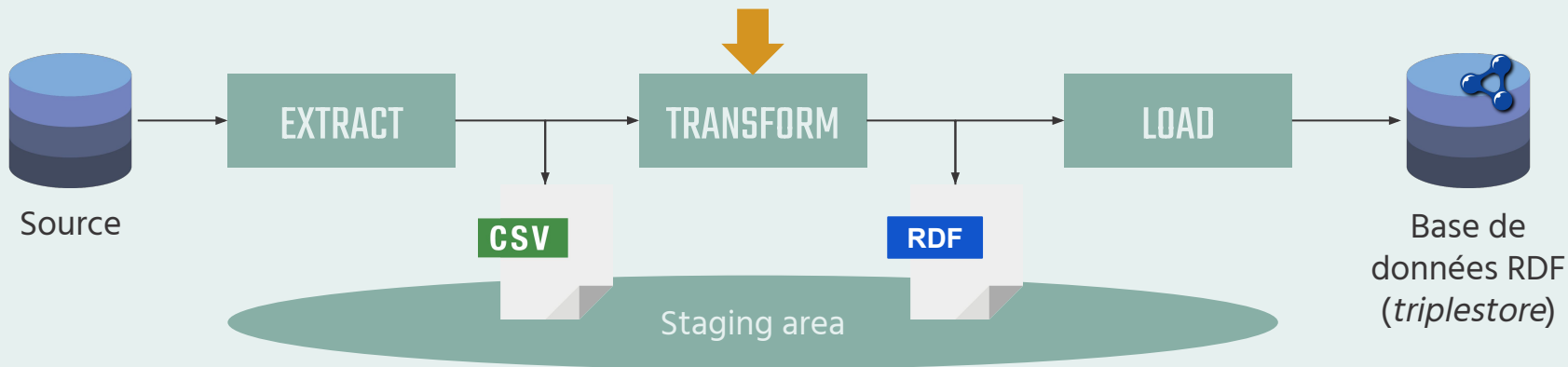


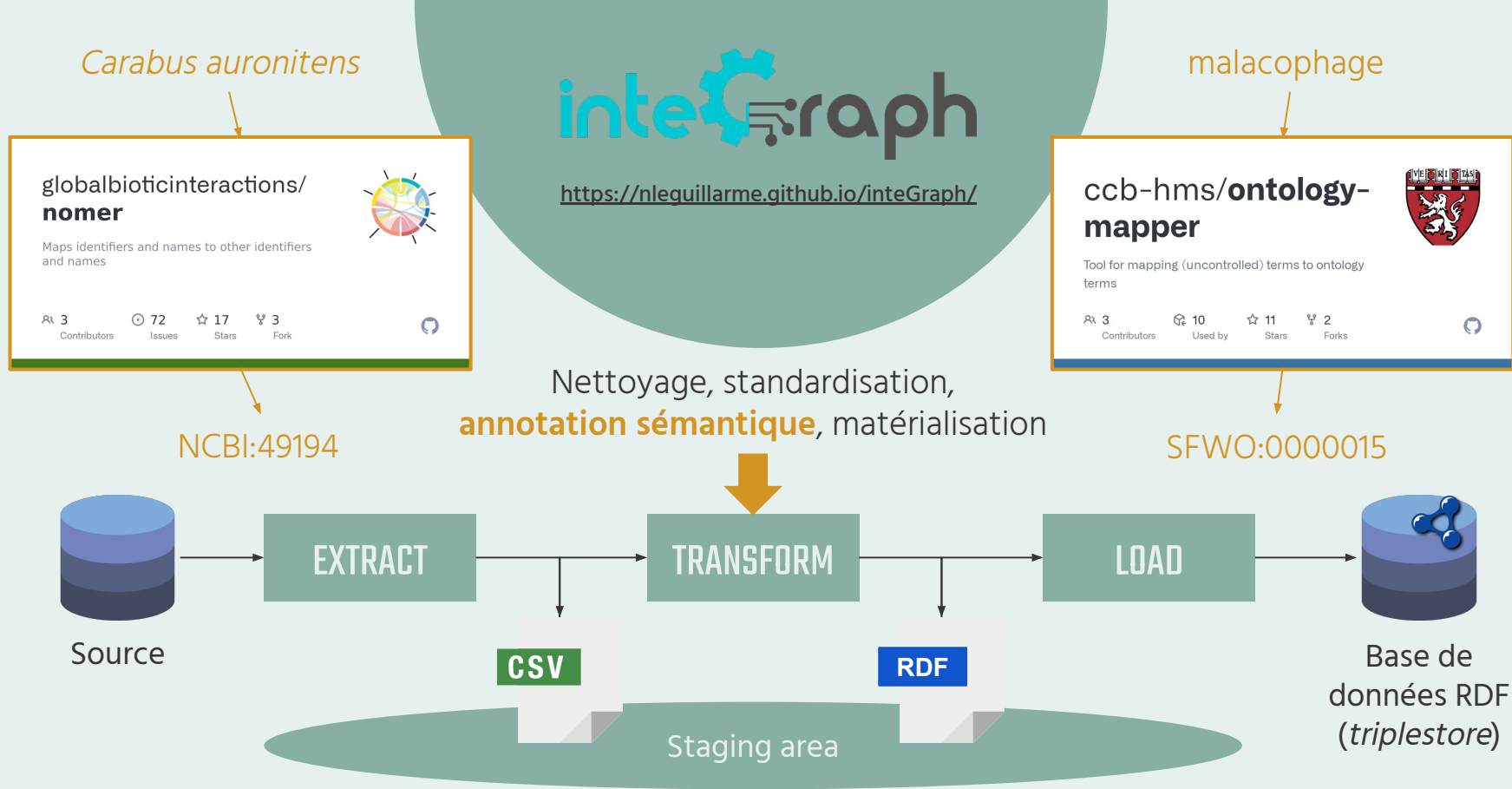
- Données structurées
- Fichiers, archives, API
- Stockage local ou source en ligne



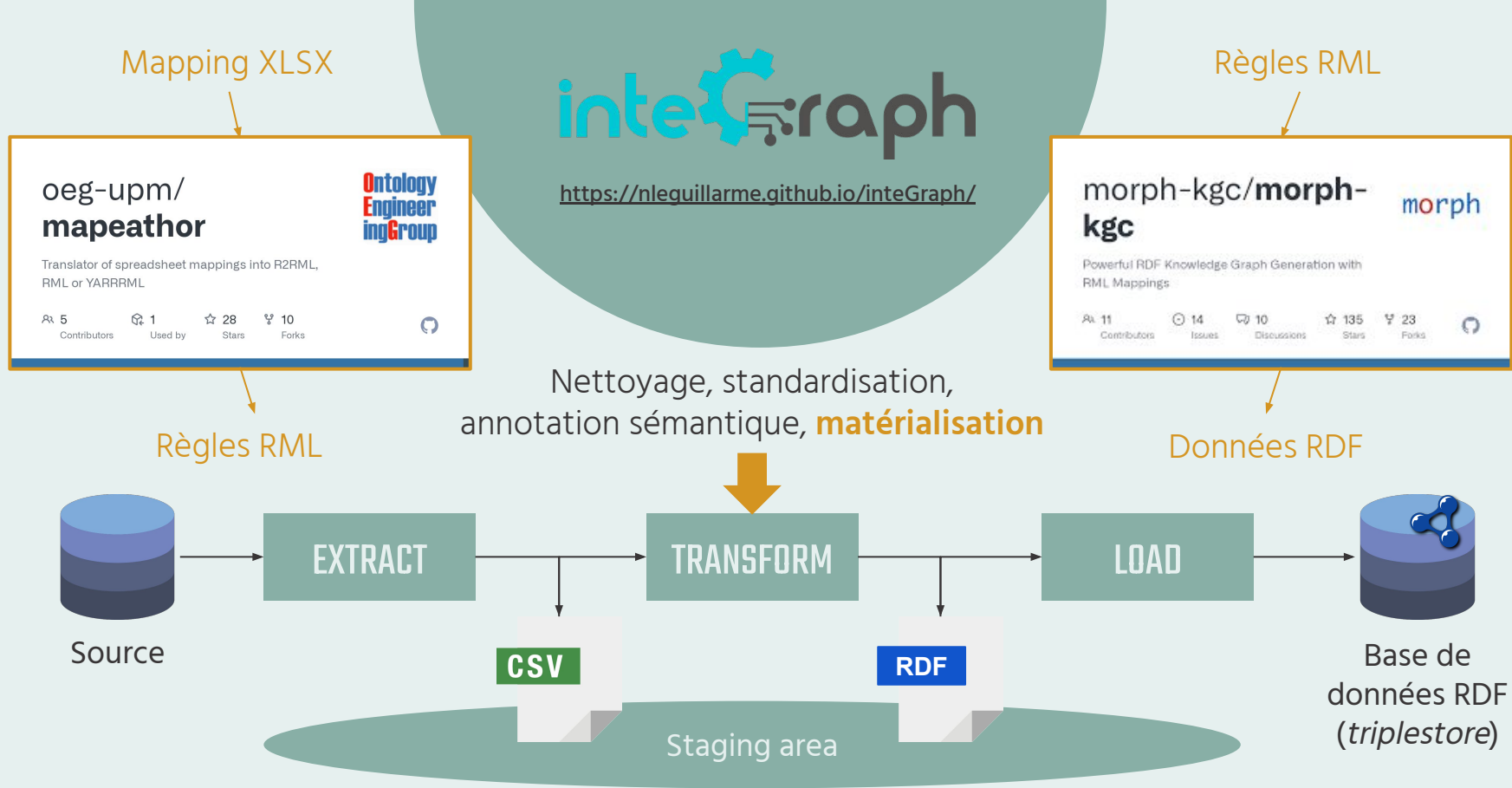


Nettoyage, standardisation,
annotation sémantique, matérialisation

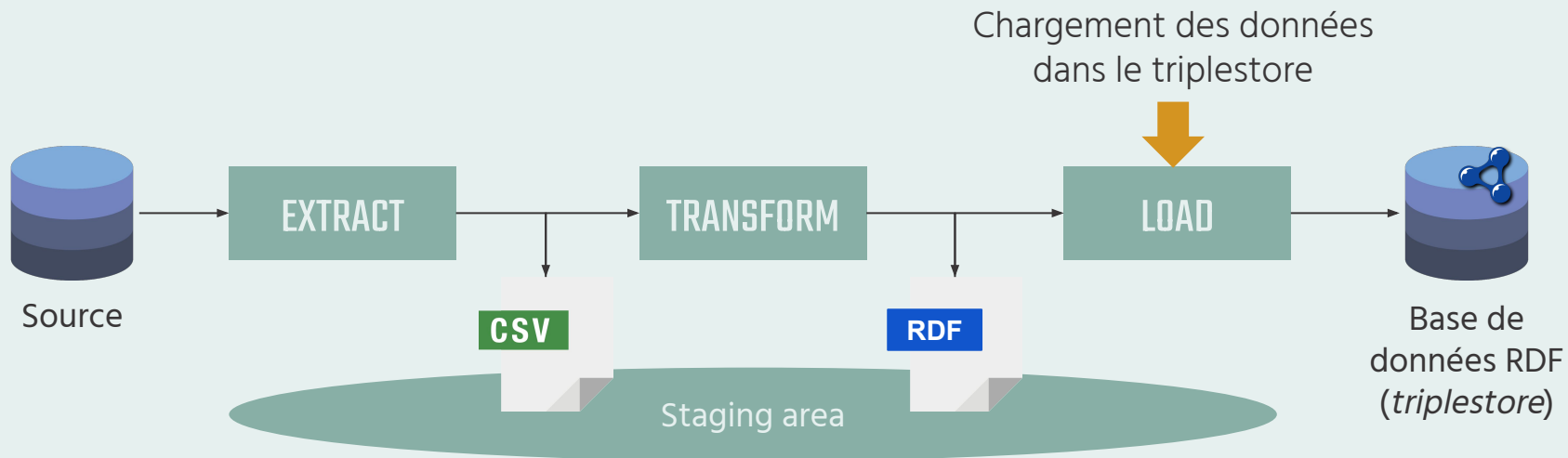




Le Guillarme, N., & Thuiller, W. (2023). A practical approach to constructing a knowledge graph for soil ecological research. *European Journal of Soil Biology*, 117, 103497.

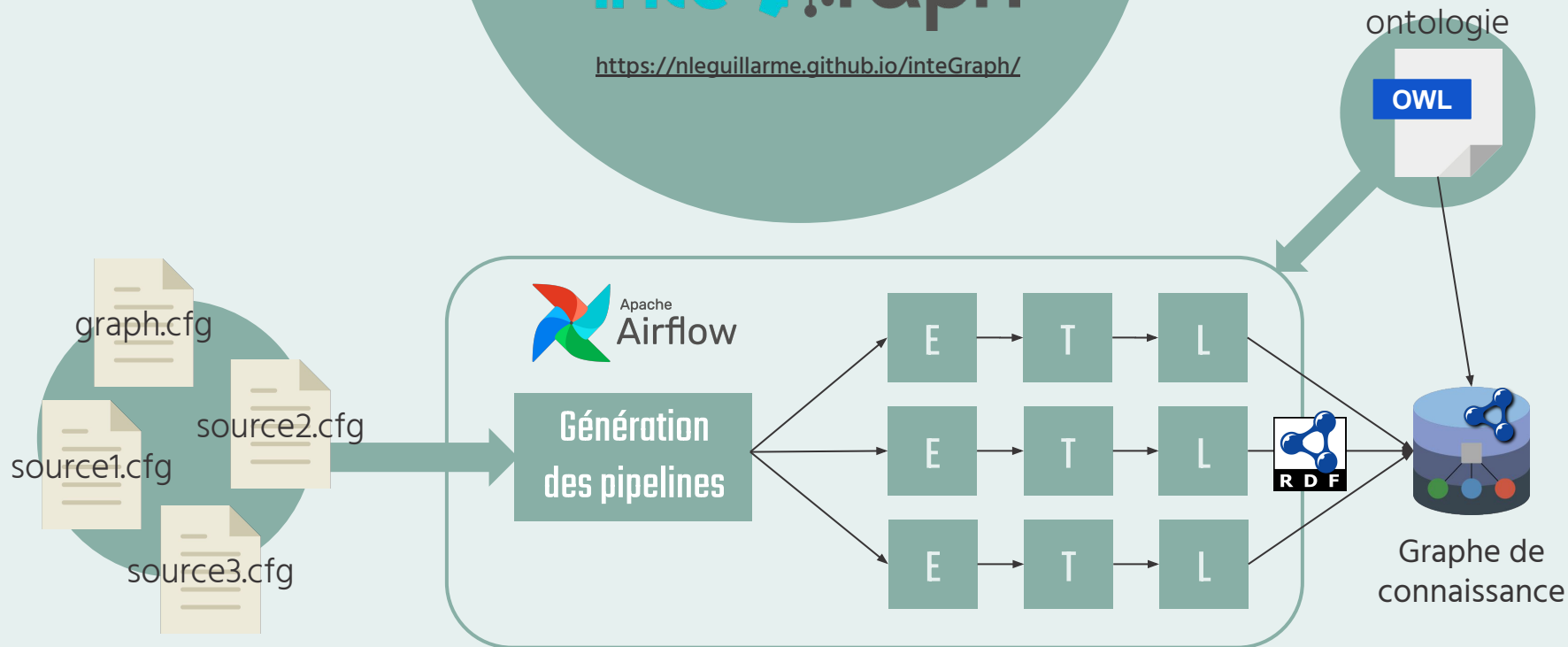


Le Guillarme, N., & Thuiller, W. (2023). A practical approach to constructing a knowledge graph for soil ecological research. *European Journal of Soil Biology*, 117, 103497.



integrate

<https://nlequillarme.github.io/inteGraph/>



Le Guillarme, N., & Thuiller, W. (2023). A practical approach to constructing a knowledge graph for soil ecological research. *European Journal of Soil Biology*, 117, 103497.



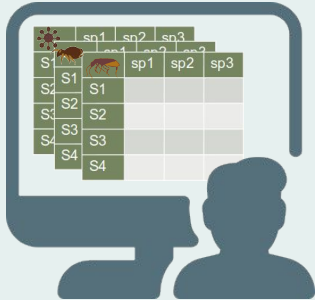
GRATIN

A graph of trophic information

GRATIN est un graphe de connaissances sur l'écologie trophique du sol qui intègre une douzaine de bases de données (bactéries, champignons, protistes, invertébrés, nématodes...)

GratinNavigatoR :
un package R pour
faciliter l'accès à GRATIN

```
get.trophic.groups(sciName="Achipteria")
```



GRATIN est un graphe de connaissances sur l'écologie trophique du sol qui intègre une douzaine de bases de données (bactéries, champignons, protistes, invertébrés, nématodes...)

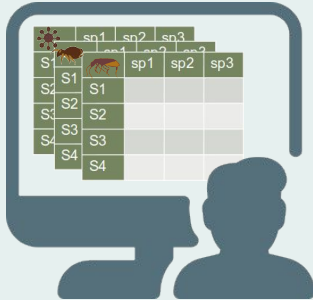
GRATIN

A graph of trophic information



GratinNavigatoR : un package R pour faciliter l'accès à GRATIN

get.trophic.groups(sciName="Achipteria")



```
SELECT ?sciName ?groupName WHERE  
{  
  ?query rdfs:label "Achipteria"  
  ?taxon rdfs:subClassOf ?query  
  ?organism ro:member_of ?taxon  
  ?organism ro:member_of ?groupIRI  
  ?taxon rdfs:label ?sciName  
  ?groupIRI rdfs:label ?groupName  
}
```



GRATIN

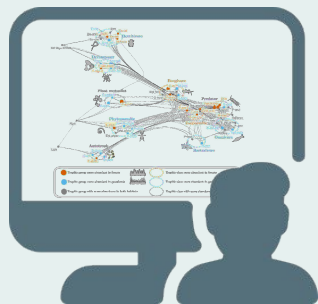
A graph of trophic information



GRATIN est un graphe de connaissances sur l'écologie trophique du sol qui intègre une douzaine de bases de données (bactéries, champignons, protistes, invertébrés, nématodes...)

GRATIN

A graph of trophic information

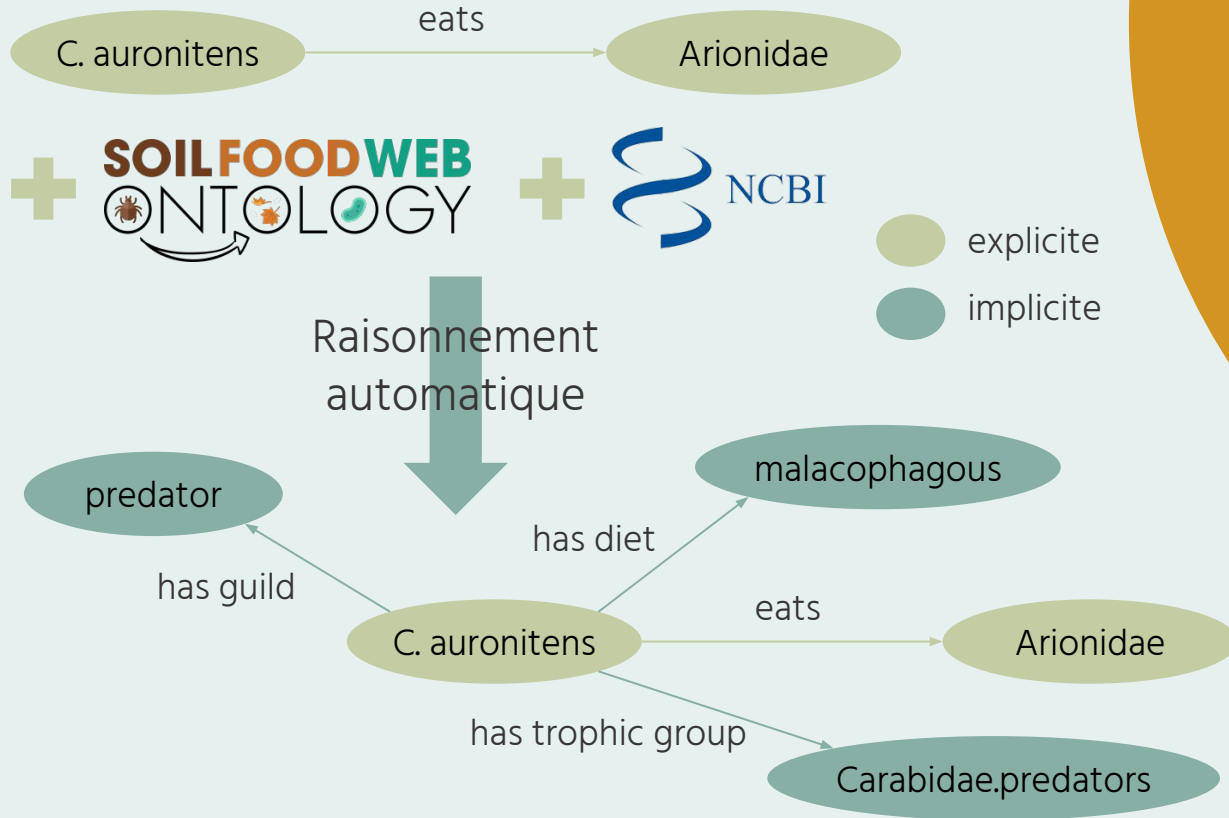


sciName	groupName
Achipteria coleoptrata	Acari.all
Achipteria coleoptrata	Oribatida.all
Achipteria coleoptrata	Oribatida.detritivores



GRATIN est un graphe de connaissances sur l'écologie trophique du sol qui intègre une douzaine de bases de données (bactéries, champignons, protistes, invertébrés, nématodes...)

Inférence dans GRATIN



GRATIN

A graph of trophic information



Inférence dans GRATIN

Il existe plusieurs variantes du langage OWL permettant de trouver un compromis entre l'**expressivité** et l'**efficacité** du raisonnement :

- OWL 2 Full
- OWL 2 EL
- OWL 2 QL
- OWL 2 RL

Critères de choix du triplestore :

- Prise en charge du profil OWL 2 RL
- Version gratuite
- Scalabilité



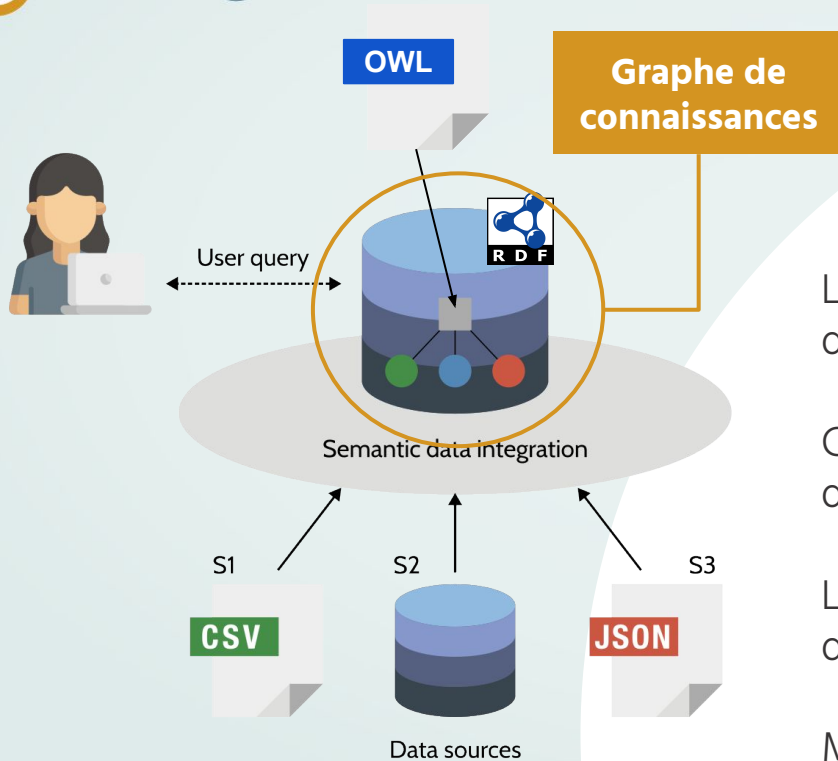
RDFOX
(in-memory)

GRATIN

A graph of trophic information



Pour conclure...



Les graphes de connaissances ont un fort potentiel de facilitation de l'accès à l'information en écologie.

Ce potentiel dépend fortement de la disponibilité d'ontologies spécialisées suffisamment expressives.

La complexité (réelle ou perçue) des technologies du web sémantique rendent leur adoption difficile.

Manque une solution de triplestore libre/gratuite.

MERCI !

Des remarques ? Des questions ?
N'hésitez pas à me contacter par mail :

nicolas.leguillarme@univ-grenoble-alpes.fr



Mickaël Hedde
(Eco&Sols, INRAE)



Wilfried Thuiller
(LECA, CNRS)



Nicolas Le Guillarme
(LECA, UGA)



Irene Calderon Sanou
(LECA, CNRS)