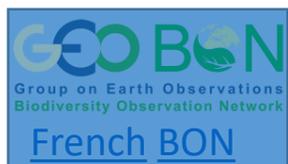# Pôle national de données de Biodiversité
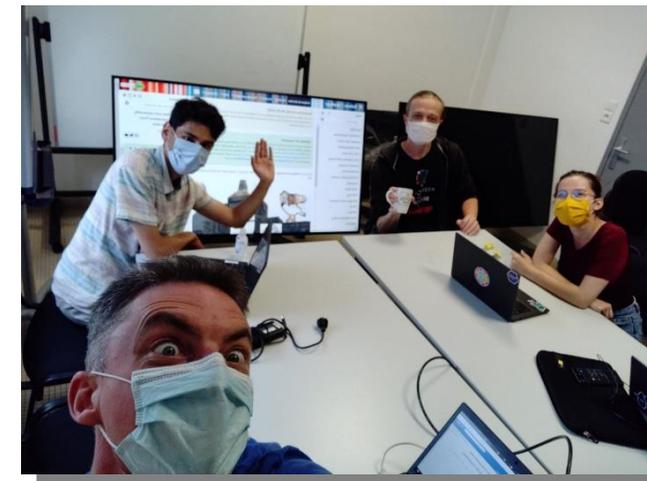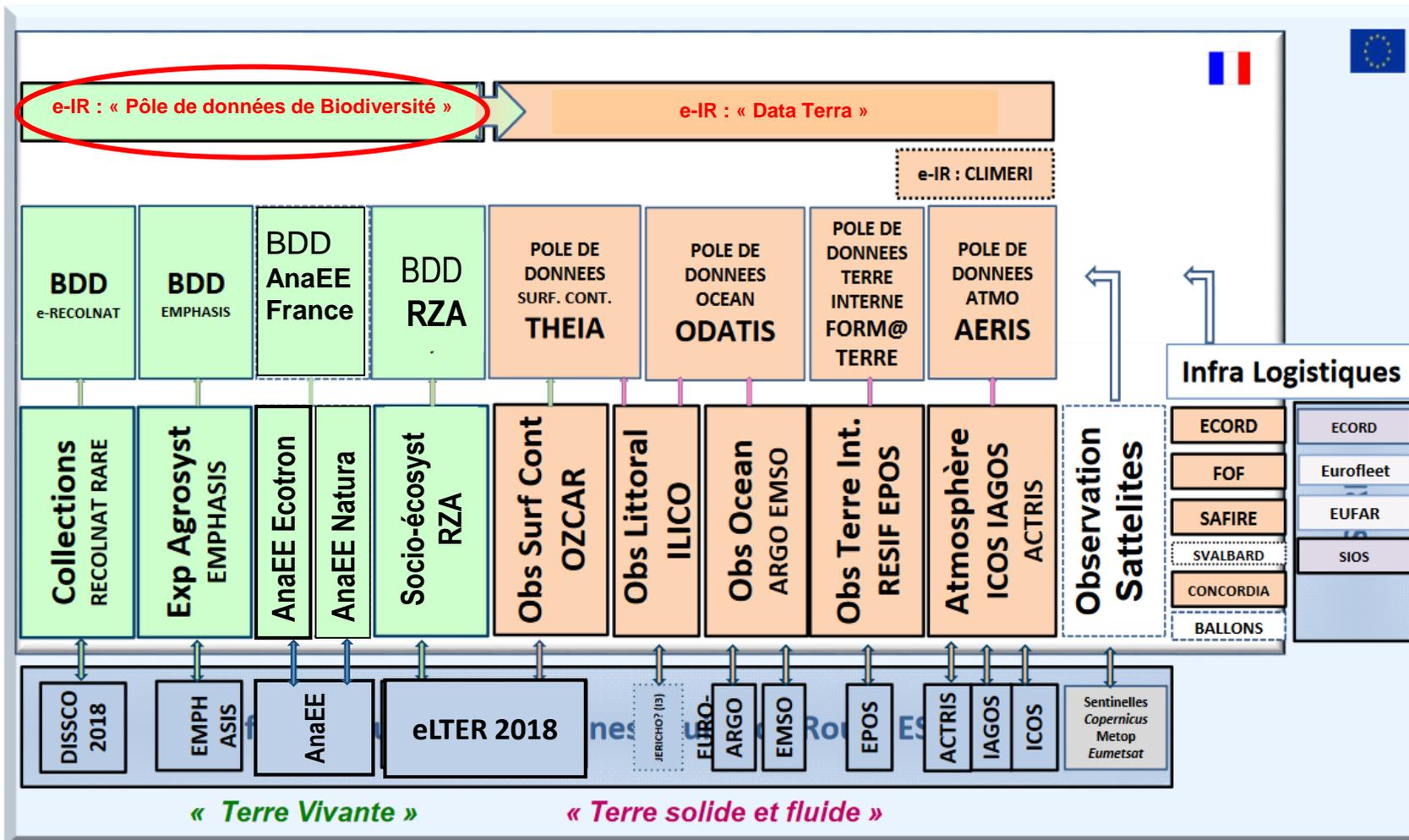## de la donnée de biodiversité au calcul scientifique via la métadonnée



GEO BON
Group on Earth Observations
Biodiversity Observation Network
French BON

GEOSS
EuroGEOSS
A EUROPEAN APPROACH TO GEOSS
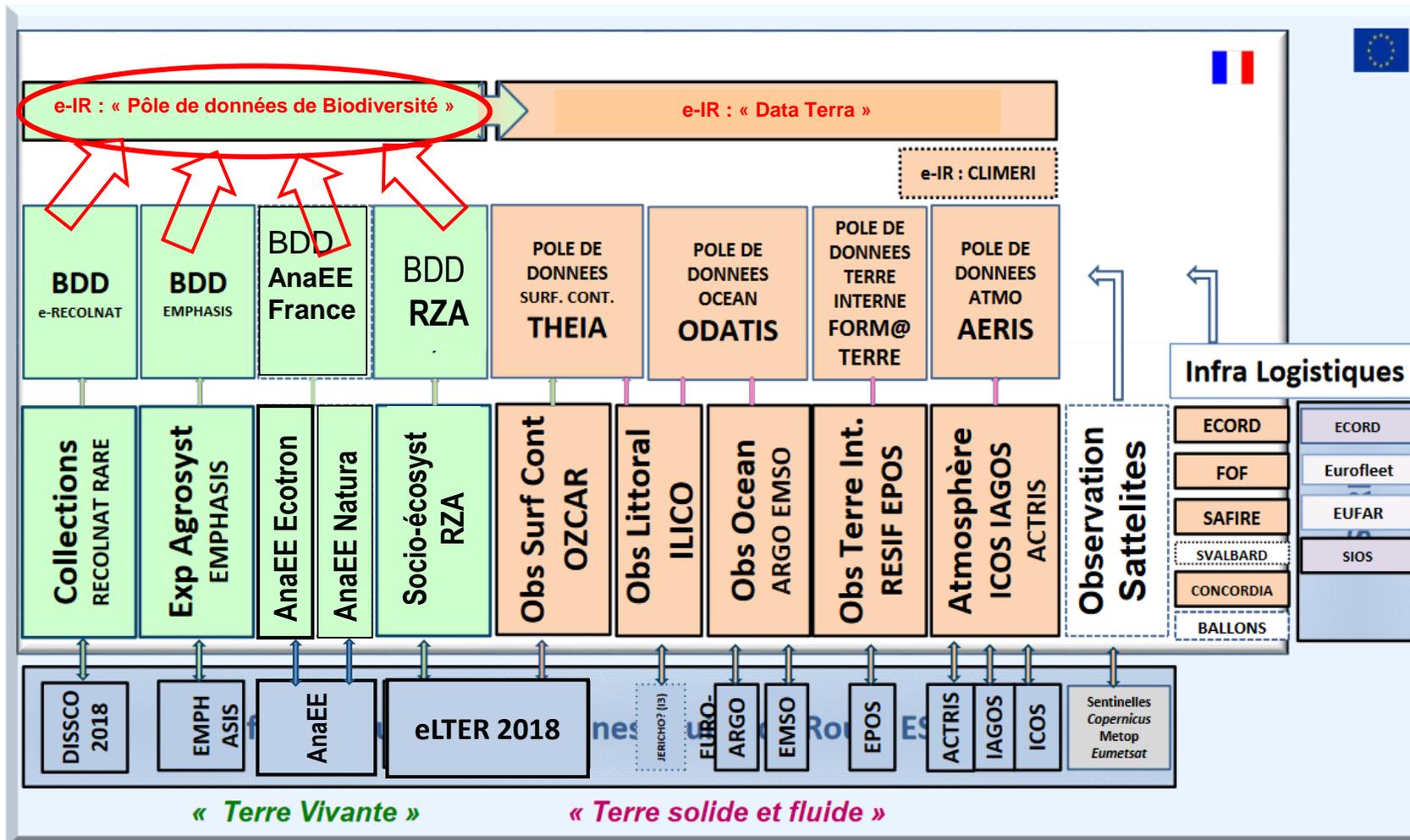
Biodiversity &
Ecosystems Action
Group

Yvan LE BRAS, Chef de projet (UMS PatriNat)
Sandrine Pavoine, MNHN (UMR CESCO)
Anne-Sophie ARCHAMBEAU, GBIF-France (UMS PatriNat)
Cécile CALLOU, Dir UMS BBEES (CNRS-MNHN)
Aurélie DELAVAUD (FRB)
Dominique JOLY, DAS CNRS (INEE)
Thomas Milon, Chef de projet "Système d'information sur la biodiversité«  (UMS PatriNat)
Laurent PONCET, Dir. UMS PatriNat, en charge du Centre de données (MNHN)
Jean-Denis VIGNE, DGD-Recherche, expertise, valorisation, enseignement MNHN

#PNDB @Yvan2935 @earnaud @ColineRoyaux @jusana_photos
*Yvan.le-bras@mnhn.fr, elie.arnaud@mnhn.fr,*
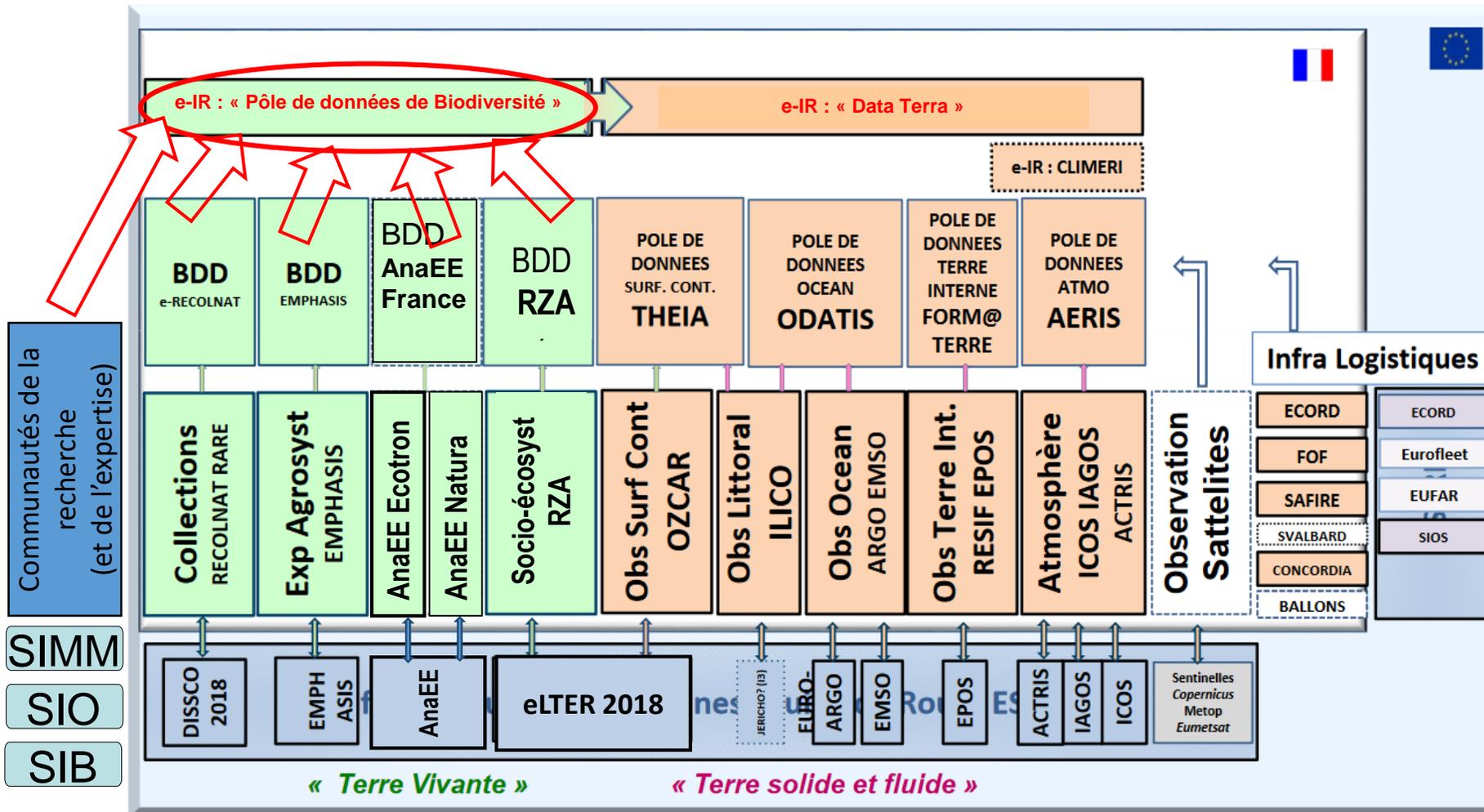*coline.royaux@mnhn.fr, julien.sananikone@mnhn.fr*

GOFAIR
BiodiFAIRse

FRB
FONDATION
POUR LA RECHERCHE
SUR LA BIODIVERSITÉ

brgm
Géosciences pour une Terre durable

cirad
LA RECHERCHE AGRONOMIQUE
POUR LE DÉVELOPPEMENT

cnrs

Ifremer

INERIS
maîtriser le risque |
pour un développement durable

INRAe
la science pour la vie, l'humain, la terre

iRD
Institut de Recherche
pour le Développement
FRANCE

MUSÉUM
NAL HIST
NATURELLE

M

OFB
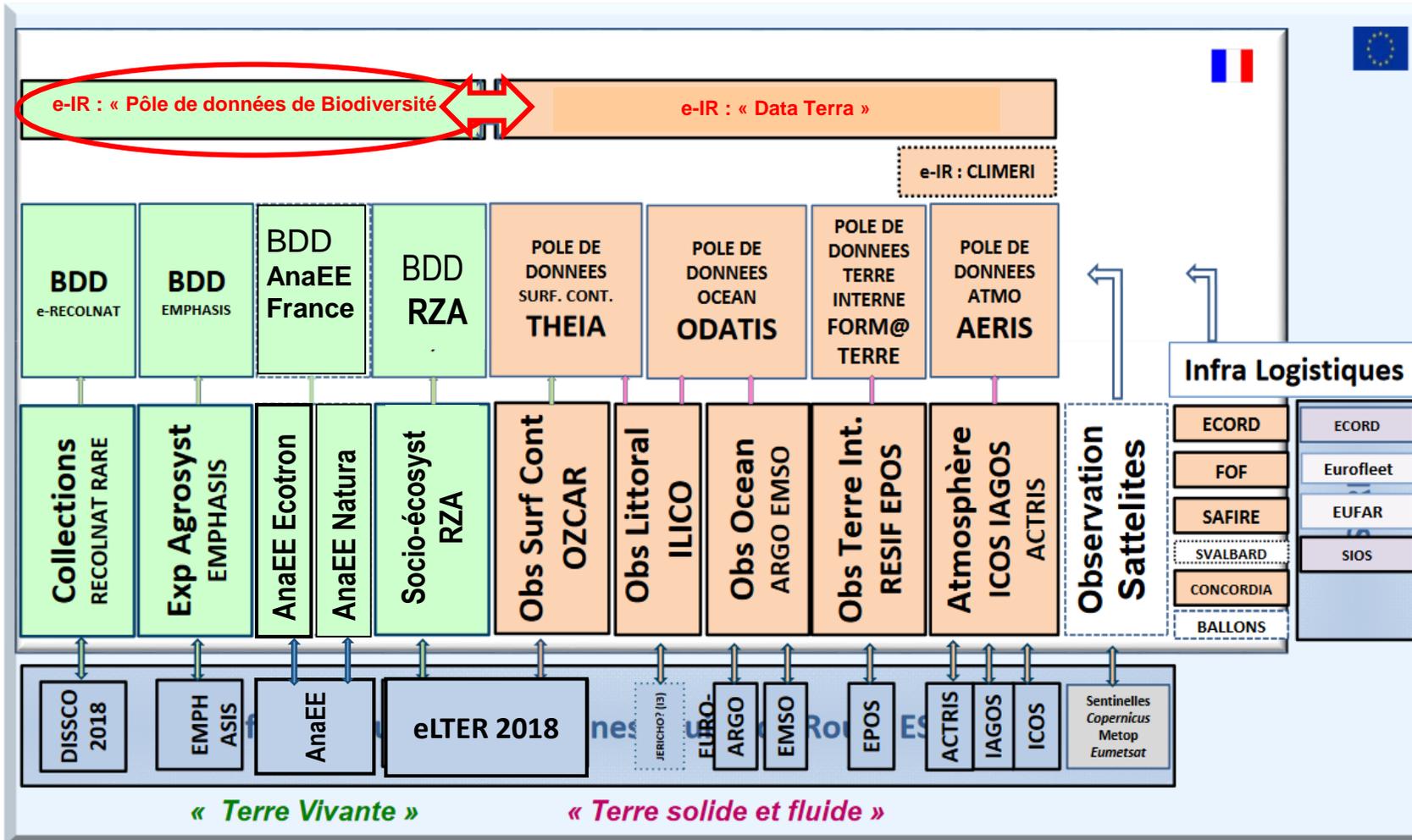OFFICE FRANCAIS
DE LA BIODIVERSITÉ

**Interface « terre vivante » IR**

**Interface « terre vivante » IR + chercheurs surtout via organismes + SI fédérateurs OFB**

Interface « terre solide » IR
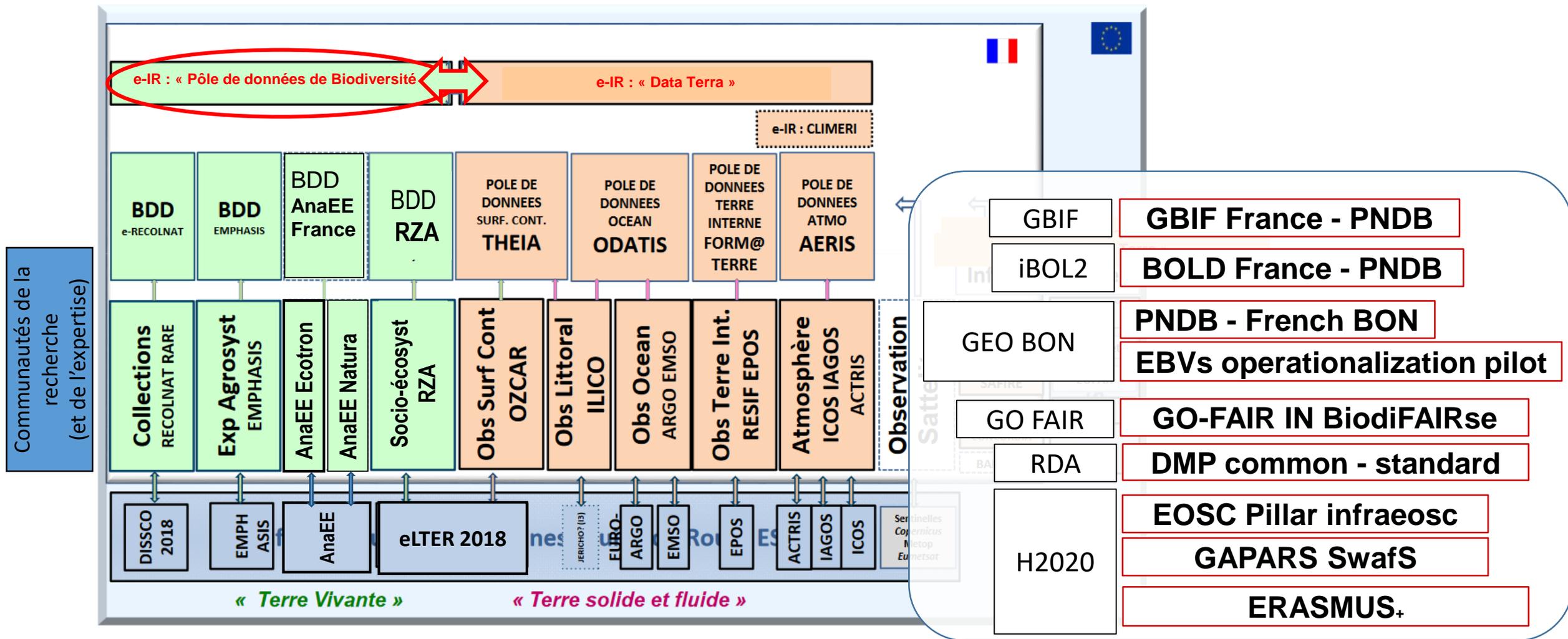
## Cohérence avec les dispositifs internationaux

# Le PNDB, e-Infrastructure nationale de recherche

- **<u>Infrastructure numérique</u>** inscrite sur la feuille de route du MESRI depuis mars 2018

- **<u>Consortium maître d'ouvrage</u> :**

  (22 partenaires institutionnels, avec le <u>soutien financier</u> du MESRI



- **<u>Maître d'œuvre</u> :**

  MNHN depuis mars 2018 (UMS PatriNat)

- **Gouvernance :**

  Comité de pilotage
  Conseil scientifique
  Comité exécutif

  Comités d'experts scientifique :
  **1. Variabilité génétique des populations domestiques ou sauvages**
  **2. Espèces, traits fonctionnels, communautés**
  **3. Ecosystèmes et socio-écosystèmes, variables SHS : structures et fonctions**
  **4. Données, scénarios, modélisation ; télédétection**

# Cahier des charges du PNDB

## 10 lignes de force

1. Orientation FAIR (aussi FAIR et *Open* que possible)
2. Relier/mutualiser avec les composantes existantes du Syst. Terre-Env.
3. Articuler/mutualiser avec le SIB-SIMM
4. Offrir des services à valeur ajoutée
5. Privilégier la qualité (*quality data*), au-delà de la quantité (*big data*),
6. Promouvoir la flexibilité des services (« à la carte »)
7. Développer, faciliter et favoriser la description fine des données
8. Viser une portée internationale (couverture & utilisation)
9. Articuler le PNDB avec les initiatives européennes et internationales
10. Pour commencer, s'appuyer sur un petit nombre de cas d'étude

(preuve de concept)

## 3 volets

1. **Accès aux métadonnées et données**
2. **Animation et accompagnement**
3. **Accès aux outils de traitement, de couplage, d'analyse (calcul)**

FRB FONDATION POUR LA RECHERCHE SUR LA BIODIVERSITÉ • brgm Géosciences pour une Terre durable • cirad LA RECHERCHE AGRONOMIQUE POUR LE DÉVELOPPEMENT • CNRS dépasser les frontières • Ifremer • INERIS maîtriser le risque pour un développement durable • INRAE la science pour la vie, l'humain, la terre • iRD Institut de Recherche pour le Développement FRANCE • MUSÉUM NAT'L HIST NATURELLE • UNIVERSITÉ MONTPELLIER • OFB OFFICE FRANÇAIS DE LA BIODIVERSITÉ

Journées Calcul et Données JCAD 2020 (03 décembre 2020)

Le paysage **(méta)données** via l' *Ecological Metadata Language*

# Le paysage **analyse** via *Github, Conda, Containers, Cloud* et *Galaxy*

**Codes** + **dépendances** → **Containers**

RStudio

**Autres plateformes**

D4science

Jupyter Notebook

**Machines virtuelles (local ou cloud)**

**Galaxy**

Codes sur Github

Dépend de

script R 1

script R n

script R y

Rproj Rmd

RStudio

ypnb

Jupyter

Ui.R Server.R

Shiny

Galaxy
PAMPA
SEEE
RapidMiner
Taverna
Nextflow ...

scripts Python
...

package R 1

package R n

package R x

Codes sur Github

Recette conda 1

SOFTWARE  CONTAINER  WORKFLOW

BioContainer

QUAY

docker hub

Codes sur Github

The Comprehensive R Archive Network

Recette conda 1

Recette conda n

Recette conda x

**Workflows**  **Histories**

# Les projets

**FNSO - OpenMetaPaper**

**ANR Challenge IA & Biodiversité**

| GBIF | **GBIF France - PNDB** |
| iBOL2 | **BOLD France - PNDB** |
| GEO BON | **PNDB - French BON** |
| | **EBVs operationalization pilot** |
| GO FAIR | **GO-FAIR IN BiodiFAIRse** |
| RDA | **DMP common - standard** |
| H2020 | **EOSC Pillar infraeosc** |
| | **GAPARS SwafS** |
| | **ERASMUS+** |

# Merci de votre attention

## PNDB team

**Coline Royaux** – Ingénieure développement scientifique R / Galaxy (workflows Galaxy pour calcul indicateurs espèces / communautés)

**Elie Arnaud** – Ingénieur développement scientifique via R Shiny / connaissance métadonnées

**Julien Sananikone** – Ingénieur DevOps / administration système & réseau / développeur web

**Yvan Le Bras** – Béta testeur    *Yvan.le-bras@mnhn.fr*

https://www.pndb.fr/

Démo saisie et publication données / métadonnées MetaShARK
https://youtu.be/OVViSMzRGtw
Démo portail de données et métadonnées
https://youtu.be/STwsYDHEt2A
Démo Galaxy Europe
- https://youtu.be/HelAHggX6D4
- Essential biodiversity variables on Galaxy: implementing the PAMPA application
- Producing biodiversity indicators from citizen science projects: update of birds and bats monitoring schemes on Galaxy-E

Visitez la version préliminaire du site web du PNDB pour suivre le projet
https://www.pndb.fr/

PNDB    Outils ▾    Animation ▾    Projets ▾    FAQ

# Pôle National de Données de Biodiversité!

Un pôle de données au service des scientifiques produisant, gérant et analysant des données de biodiversité

**Accéder aux données »**    **Tester MetaShaRK »**

En 2018, le Ministère de l'Enseignement supérieur, de la recherche et de l'innovation a inscrit sur sa feuille de route la création d'une nouvelle infrastructure intitulée Pôle National de données de biodiversité (PNDB). Les missions du PNDB s'inscrivent dans une approche FAIR (Facile à trouver, Accessible, Intéropérable, Réutilisable), et consistent à :

1. fournir un accès aux jeux de données et de métadonnées, à des services associés et à des produits dérivés des analyses
2. promouvoir l'animation scientifique pour identifier les lacunes et favoriser l'émergence de dispositifs portés par des communautés d'utilisateurs et producteurs
3. faciliter le partage des pratiques avec les autres communautés de recherche, favoriser le partage des données et leur réutilisation, s'insérer dans la réflexion de la future infrastructure Système Terre.
4. favoriser la cohérence avec les efforts nationaux, européens et internationaux relatifs à l'accès et à l'exploitation des données de recherche sur la biodiversité, à la promotion de produits et services.

Le PNDB est porté par le Muséum national d'Histoire naturelle, plus particulièrement par l'UMS 2006 PatriNat, unité MNHN CNRS et AFB. Le projet est en lien étroit avec la FRB et plusieurs de ses institutions fondatrices (AFB, BRGM, CIRAD, CNRS, Ifremer, INERIS, INRA, IRD, IRSTEA, MNHN, Univ. Montpellier).

## Portail de données        ## MetaShark        ## Galaxy

Visitez le portail de données du PNDB et jouez avec les premiers jeux de données-métadonnées
https://data.pndb.fr/

Testez l'interface de saisie de données et métadonnées du PNDB
https://metashark.test.pndb.fr/

Testez la plateforme d'analyse / couplage de données du PNDB https://ecology.usegalaxy.eu/



Tutoriels : https://training.galaxyproject.org/

Codes sources : https://github.com/65MO/Galaxy-E
https://github.com/galaxyecology/tools-ecology

# Pourquoi ce fameux *Ecological Metadata Language* ?

## Les LTER -> THE exemple US & Australie

# Pourquoi ce fameux *Ecological Metadata Language* ?

## Les LTER -> THE exemple US & Australie



EML (MTD std) https://ropensci.github.io/EML/articles/creating-EML.html

EML Assembly Line https://ediorg.github.io/EMLassemblyline/index.html

Metacat (data portal) https://github.com/NCEAS/metacatui

OAI-PMH (harvesting)

NCEAS / EDI / KNB

DataOne

# Pourquoi ce fameux *Ecological Metadata Language* ?

## Les LTER -> THE exemple Europe

# Pourquoi ce fameux *Ecological Metadata Language* ?

## Les LTER -> THE exemple Europe

~2010

**Drupal Ecological Information Management System**



2018

**DEIMS-SDR Dynamic Ecological Information Management System Site and Dataset Registry**

# Pourquoi ce fameux *Ecological Metadata Language* ?

## Amazing EML

Language + Modules -> flexibility!

### Semantic Module functionnality

**Principle**

Write annotations: kind of "sentence" composed with:

- a subject
- a property
- a property's value

(similar to RDF statement)

An annotation is added inside an *Attribute* as its direct child (cf. example).

### Example

*Natural language:*

attribute 'weight' uses standard unit 'gram'

*EML specification:*

'attribute7' 'usesStandard' 'Gram'

Semantics!

```
<attribute id="att.12">
  <attributeName>biomass</attributeName>
  ...
  <annotation>
    <propertyURI label="of characteristic">http://ecoinformatics.org/oboe/oboe.1.2/oboe-core.owl#ofCharacteristic</propertyURI>
    <valueURI label="Mass">http://ecoinformatics.org/oboe/oboe.1.2/oboe-characteristics.owl#Mass</valueURI>
  </annotation>
  <annotation>
    <propertyURI label="of entity">http://ecoinformatics.org/oboe/oboe.1.2/oboe-core.owl#ofEntity</propertyURI>
    <valueURI label="Plant Sample">http://example.com/example-vocab-1.owl#PlantSample</valueURI>
  </annotation>
</attribute>
```

Provenance tracking!

Constraints definition!

Deploy Analytics environment on the cloud!

= uniq ID

Proposition de rôle via BED

- Liens communautés ZA -> PNDB

- Support technique PNDB

- Retours / conseils vis-à-vis initiatives internationales (GEO BON dont EBV / DataOne ..)

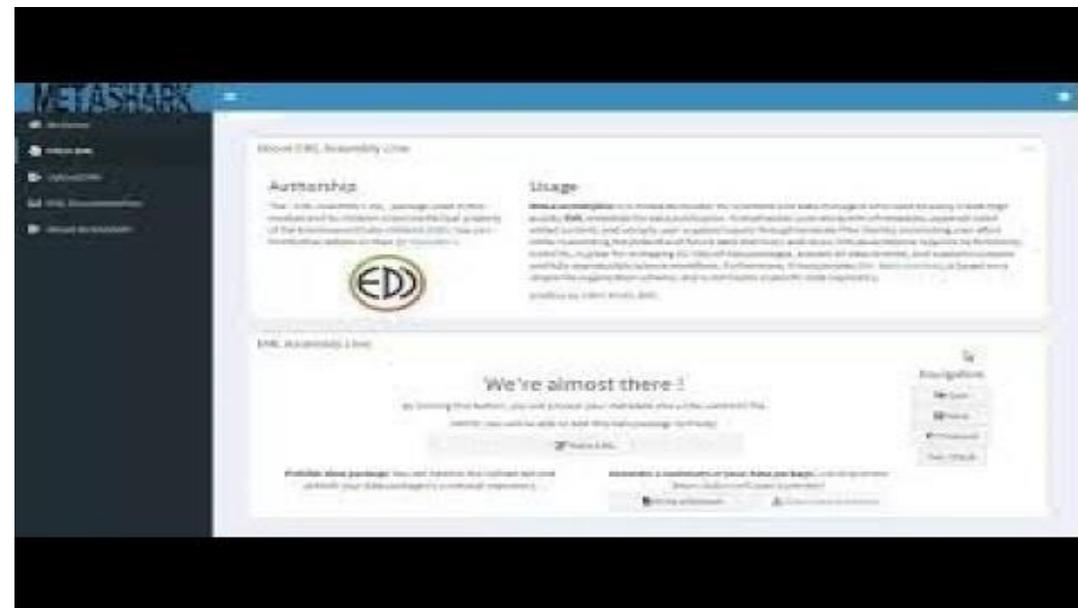# Démonstrations !!!!!!!



**Structuration**

**Accessible** **FAIR**

**Publication**
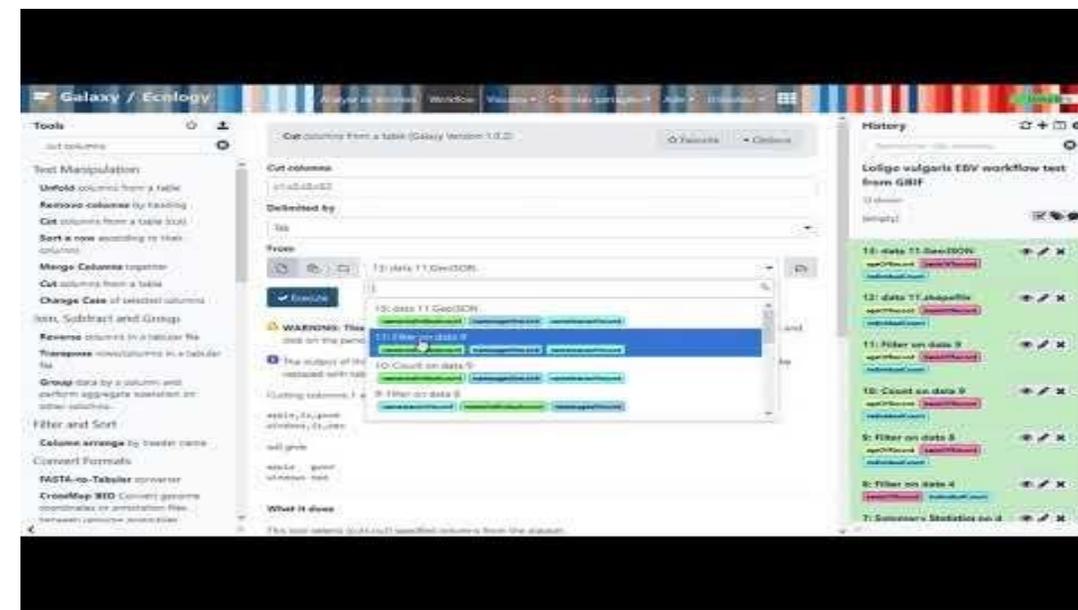
**Collaboratif**

**Reproductible**

**Analyse**

Démo saisie et publication
données / métadonnées
MetaShARK
https://youtu.be/OVViSMzRGtw

Démo portail de données et métadonnées basé sur Metacat/metacatui/Geoserver
https://youtu.be/STwsYDHEt2A

Démo plateforme d'analyse de données de biodiversité via Galaxy Europe
https://youtu.be/HelAHggX6D4

# MetaShARK

# Au début était l'EML

Ecological Metadata Language (by knb )

# Au début était l'EML

Ecological Metadata Language (by knb )

Solutions techniques existantes

- MetaCAT ( → DataONE )
- Morpho (écriture métadonnées)

# Au début était l'EML

Ecological Metadata Language (by knb )

Solutions techniques existantes

- MetaCAT ( → Data⊗NE )

- Morpho (écriture métadonnées)

  - **MetaShARK** =

    EML Assembly Line (by EDI ) + Travail de stage 2019

# Documentation

- Localisation

- Interactions

- Recherche

- Utilisation

- Infos techniques

# Saisie

## Etape 1: sélection d'un projet

- Structuration en *data packages:* données + métadonnées

- Gestion de la propriété intellectuelle

# Saisie

**Etape 2: téléchargement des fichiers à décrire**

- Fichiers tableurs (pour l'instant)

- Description de chaque fichier

# Saisie

Etape 3&4: attributs

– Typage + explication des attributs (=variables)

– Définition des modes possibles des variables catégorielles

– "Custom units"

# Saisie

Etape 5: couverture géographique

- Définition en employant des colonnes du jeu de données

- Définition manuelle, site par site

# Saisie

Etape 6: couverture taxonomique

- Définition en employant des colonnes du jeu de données

- Connection à des ressources terminologiques (GBIF, ITIS)

# Saisie

## Etape 7: personnel & rôles

- Remplissage automatique par moissonnage



- Précision du rôle par personne

# Saisie

Etape 8: "détails"

– Abstract

– Méthodes

– Durée

– Mots-clés

– Informations additionnelles

# Saisie

## Etape 9: écriture

- – Génération de .xml "EML-valide"

- – Possibilité de proto-data paper (*emldown*)

# Résultats

- Data package téléchargeable

- *emldown* succinct

# Résultats

- Data package téléchargeable

- *emldown* succinct

- Mise en ligne possible vers MetaCAT

  **Exemple** *Study of specific diversity of phytoplankton populations*, F. Rigault-Jalabert (lien)

# WIPs

- EML Annotation: connection aux ontologies

- Mise en ligne: réparation en cours

- Galaxyfication

- Conversion vers d'autres standards (DwC, SINP, ISO 19k+, …)

- Sessions personnalisées + profilage des utilisateurs par communauté

# Galaxy-E : une instance de Galaxy dédiée à l'analyse de données en Ecologie

Coline Royaux

# Galaxy

**Galaxy PROJECT**

Plateforme web pour le partage et le traitement des
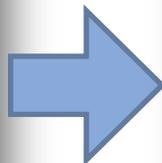données de recherche

Permet un accès facilité à l'analyse cloud et au Calcul Haute
Performance par l'interfaçage de n'importe quel langage informatique

### Quatre principes clés

✔ Accessibilité ✔ Reproductibilité ✔ Transparence ✔ Peer review

# Galaxy

# Le workflow Galaxy PAMPA



**Calculate community metrics** calculate community metrics from abundance data (Galaxy Version 0.0.1)

**Input file**

No tabular dataset available.

Observation data file, with location, year, species and abundance.

**Choose the community metrics you want to compute**

Select/Unselect all

× All

Presence/absence, Species richness, Simpson and Shannon index are systematically computed.

✔ Execute

Calculate community metrics from abundance data

Communauté

Calcul des métriques de communauté

Données prétraitées

**1**

Calcul des métriques populationnelles

Espèces - population

**Calculate presence absence table** calculate presence absence table from observation data (Galaxy Version 0.0.)

**Input file**

No tabular dataset available.

Observation data file, with location, year, species and abundance.

**Email notification**

Yes   No

Send an email notification when the job completes.

✔ Execute

Calculate presence absence table from abundance data

# Le workflow Galaxy PAMPA



Calculate community metrics from abundance data

Compute GLM on community data with selected interest variables

**Communauté**

Calcul des métriques de communauté → GLM sur les métriques de communauté

**Données prétraitées**

**1**

**2**

Métrique ~ site + année + habitat

Calcul des métriques populationnelles → GLM sur les métriques populationnelles

**Espèces - population**

Calculate presence absence table from abundance data

Compute GLM on population data with selected interest variables

5

# Le workflow Galaxy PAMPA



Communauté

Calcul des métriques de communauté

GLM sur les métriques de communauté

Données prétraitées

**1**

**2**

Métrique ~ site + année + habitat

**3**

Représentation graphique des résultats de GLM

Espèces - population

Calcul des métriques populationnelles

GLM sur les métriques populationnelles

6

# Les Variables Essentielles de Biodiversité (EBV)



Kissling *et al.* 2017

# Tutoriel sur données de pêche (CPUE)

https://bit.ly/2VjDTaQ

# Compute and analyze Essential Biodiversity Variables with PAMPA toolsuite

By: 🧑 Coline Royaux

| Atlantic cod | European plaice | European flounder |
|---|---|---|
| *Gadus morhua* | *Pleuronectes platessa* | *Platichthys flesus* |

## Overview

**❓ Questions**
- How to evaluate properly species populations and communities biological state with abundance data?
- How does trawl exploited populations of Baltic sea, Southern Atlantic and Scotland are doing over time?
- How to compute and analyze Essential Biodiversity Variables (EBV) on abundance data?

**◎ Objectives**
- Upload data from DATRAS portal of ICES
- Pre-process population data with Galaxy
- Learning how to use an Essential Biodiversity Variables (EBV) scientific workflow from raw data to graphical representations
- Learning how to construct a Generalized Linear (Mixed) Model from a usual ecological question
- Learning how to interpret a Generalized Linear (Mixed) Model

**✔ Requirements**
- Introduction to Galaxy Analyses

**⧗ Time estimation:** 2 hours

**↗ Supporting Materials**
    📄 Datasets    ◀ Workflows    🌐 Available on these Galaxies ▾

**📅 Last modification:** Nov 19, 2020

Introduction

Upload and pre-processing of the data

Compute Essential Biodiversity Variables

# Introduction

This tutorial aims to present the PAMPA Galaxy workflow and how to use it to compute Essential Biodiversity Variables (EBV) from species abundance data and analyse it through generalized linear (mixed) models (GLM and GLMM). This workflow made up of 5 tools will allow you to process temporal series data that include at least `year`, "location" and species sampled along with abundance value and, finally, generate article-ready data products.

# Outils Galaxy-E **généralisable**

## Variables primaires smilaires

### Données Suivi Temporel des Oiseaux Communs

| | carre | annee | espece | abond |
|---|---|---|---|---|
| 1 | 2 | 2016 | ACCGEN | 0 |
| 2 | 2 | 2017 | ACCGEN | 0 |
| 3 | 2 | 2018 | ACCGEN | 0 |
| 4 | 2 | 2019 | ACCGEN | 0 |
| 5 | 183 | 2016 | ACCGEN | 0 |
| 6 | 183 | 2017 | ACCGEN | 0 |
| 7 | 183 | 2018 | ACCGEN | 0 |
| 8 | 183 | 2019 | ACCGEN | 0 |
| 9 | 2 | 2019 | ACCGEN | 0 |

### Données de caméras sous-marines

| | UnitObs | rotation | codeSp | sexe | taille | classe_taille | poids | nb_ind |
|---|---|---|---|---|---|---|---|---|
| 1 | AS140155 | 3 | Hemifasc | -999 | -999 | P | -999 | 1 |
| 2 | AS140159 | 1 | Nasosp. | -999 | -999 | P | -999 | 3 |
| 3 | AS140159 | 3 | Gompvari | -999 | -999 | P | -999 | 1 |
| 4 | AS140160 | 3 | Gompvari | -999 | -999 | P | -999 | 1 |
| 5 | AS140099 | 2 | Parumult | -999 | -999 | P | -999 | 1 |
| 6 | AS140088 | 1 | Varilout | -999 | -999 | P | -999 | 1 |
| 7 | AS140088 | 2 | Gompvari | -999 | -999 | P | -999 | 2 |
| 8 | AS140041 | 1 | Nasosp. | -999 | -999 | P | -999 | 5 |
| 9 | AS140044 | 1 | Parumult | -999 | -999 | P | -999 | 4 |

### Données de pêche

| | Survey | Year | Quarter | Area | AphiaID | Species | LngtClass | CPUE_number_per_hour |
|---|---|---|---|---|---|---|---|---|
| 1 | BITS | 1991 | 1 | 22 | 126281 | Anguilla anguilla | 0 | 0.000000 |
| 2 | BITS | 1991 | 1 | 22 | 126281 | Anguilla anguilla | 720 | 0.009160 |
| 3 | BITS | 1991 | 1 | 22 | 126417 | Clupea harengus | 0 | 0.000000 |
| 4 | BITS | 1991 | 1 | 22 | 126417 | Clupea harengus | 80 | 0.075785 |
| 5 | BITS | 1991 | 1 | 22 | 126417 | Clupea harengus | 85 | 0.088277 |
| 6 | BITS | 1991 | 1 | 22 | 126417 | Clupea harengus | 95 | 0.037892 |
| 7 | BITS | 1991 | 1 | 22 | 126417 | Clupea harengus | 100 | 0.063293 |
| 8 | BITS | 1991 | 1 | 22 | 126417 | Clupea harengus | 105 | 0.012492 |
| 9 | BITS | 1991 | 1 | 22 | 126417 | Clupea harengus | 110 | 0.618357 |

### Données Vigie - Chiro

| | participation | Nuit | num_micro | groupe | espece | nb_contacts |
|---|---|---|---|---|---|---|
| 1 | 55de2cd52121b1000d27430e | 2015-07-26 | 0 | bat | Barbar | 1 |
| 2 | 55de2cd52121b1000d27430e | 2015-07-26 | 0 | bush-cricket | Barfis | 1 |
| 3 | 55de2cd52121b1000d27430e | 2015-07-26 | 0 | noise | noise | 5022 |
| 4 | 55de2cd52121b1000d27430e | 2015-07-26 | 0 | bush-cricket | Decalb | 5 |
| 5 | 55de2cd52121b1000d27430e | 2015-07-26 | 0 | bush-cricket | Tyllil | 18 |
| 6 | 55de2cd52121b1000d27430e | 2015-07-26 | 0 | bat | Nyclei | 1 |
| 7 | 55de2cd52121b1000d27430e | 2015-07-26 | 0 | bush-cricket | Phanan | 269 |
| 8 | 55de2cd52121b1000d27430e | 2015-07-26 | 0 | bat | Eptser | 5 |

- Site
- Année
- Espèce
- Occurence

# Atomisation

Actuellement, en écologie…

Un script R pour chaque fichier de données
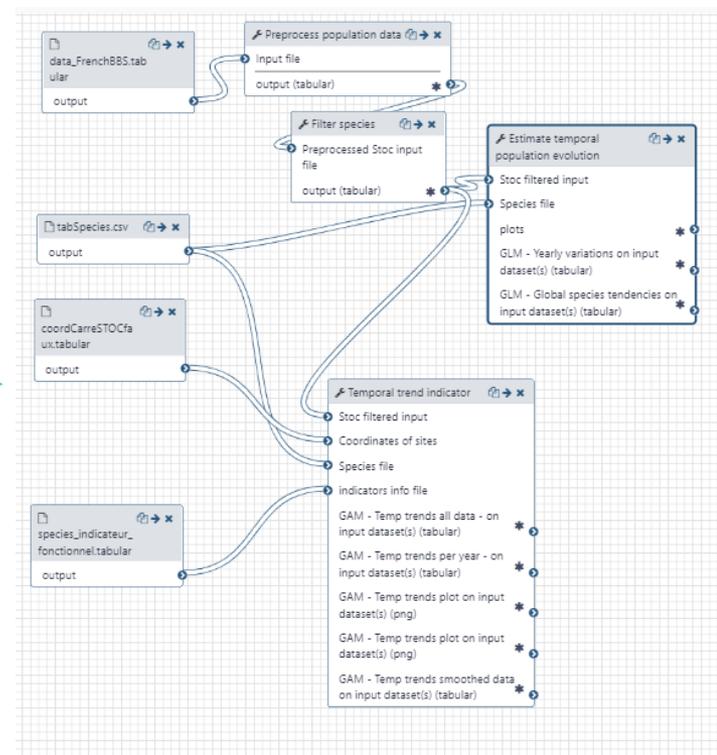
```
                    ,direction="wide")
  tab[is.na(tab)] <- 0
                                    #  filename <- "touverUnNom"
                                    #  chemin <- paste(rep,filename,sep="/")
                                    #  write.table(tab, chemin)
 colnames(tab) <- sub("nombre.","",colnames(tab))

    return(tab)
}

## sous jeux de donnees si choix d espece d annee ou d un pourcentage de carres
makeSousTab <- function(tab,vecSp=NULL,echantillon=1,
                        methodeEchantillon="carre",vecannees=NULL) {
    cat("  -- Fabrication du sous jeu de donnees --\n")
    flush.console()
   ## reduction de la table à certaine espèces
   if(!is.null(vecSp)) {
        cat("      selection",length(vecSp),"espece(s):\n -> ")
        cat(vecSp)
        cat("\n")
        tab <- data.frame(carre = tab$carre,annee = tab$annee,tab[,vecSp])
        colnames(tab) <- c("carre","annee",vecSp)
    }
   ## reduction de la table pour certaines annees
    if(!is.null(vecannees)) {
        tab <- subset(tab,annee>=vecannees[1] & annee <= vecannees[2])

    }
   ## reduction de la table par une proportion de carre suivie
    if(echantillon != 1) {
        if(echantillon < 1 & echantillon >0) {
            nbinit <- nrow(tab)
            if(methodeEchantillon == "global") {
                nb <- round(nrow(tab)*echantillon)
                cat("     echantillonage",echantillon*100,
                    "% des donnees par la methode",methodeEchantillon,"\n")
                cat(" -> conservation de",nb,"lignes sur",nbinit,"\n")
                flush.console()
                tab <- tab[sample(1:nrow(tab))[1:nb],]
            } else {
                if (methodeEchantillon =="carre") {
                    cat("     echantillonage",echantillon*100,
                        "% des carrees par la methode",methodeEchantillon,"\n")
                    nbcarreinit <- length(unique(tab$carre))
                    chat=sample(unique(tab$carre),
                        length(unique(tab$carre))*echantillon,replace=F)
                    cat(" -> conservation de",length(chat),"carrees sur",
                        nbcarreinit)
                    tab=subset(tab, subset = carre %in% chat)
                    cat(" (",nrow(tab)," lignes sur ",nbinit,")\n",sep="")
                } else {

                    stop("Methode d echantillonnage non reconnue")
```

# Atomisation

Actuellement, en écologie...

Un script R pour chaque fichier de données

# Atomisation

Actuellement, en écologie…

Un script R pour chaque fichier de données

Avec Galaxy…

Plusieurs scripts R atomisés pour analyser différents fichiers de données

# Atomisation

Actuellement, en écologie…

Un script R pour chaque fichier de données



Avec Galaxy…

Plusieurs scripts R atomisés pour analyser différents fichiers de données