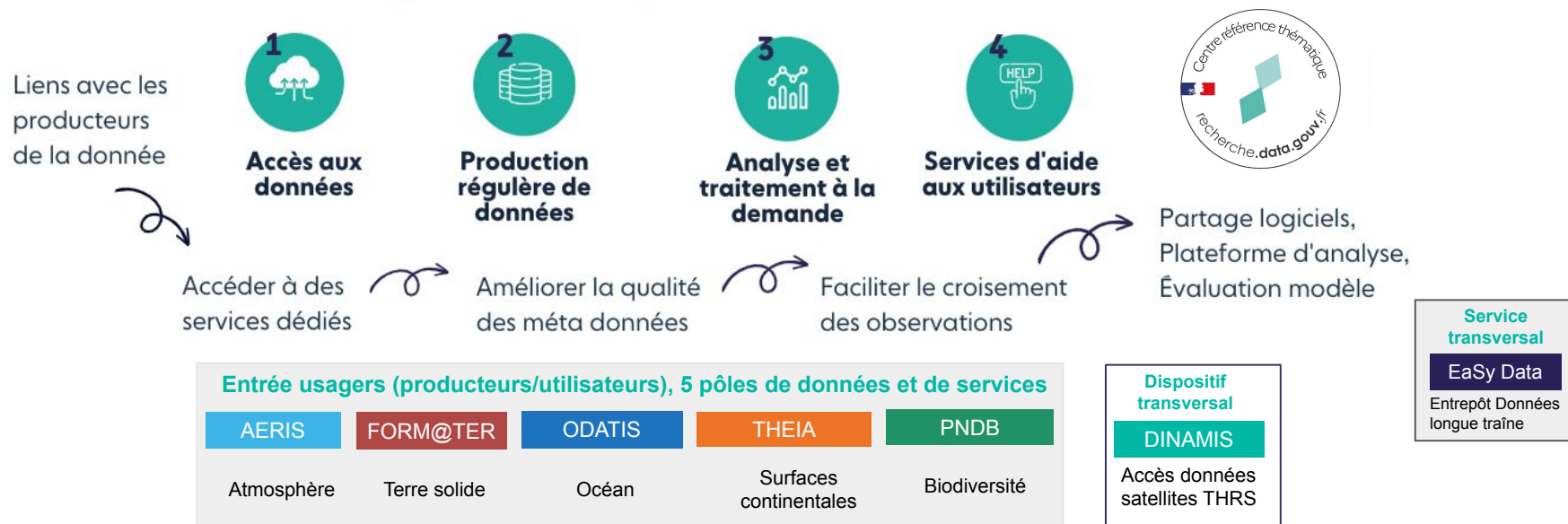
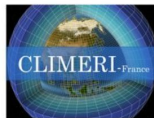


L'IR DATA TERRA propose des **services** autour des données d'observation du **système Terre interopérables et interdisciplinaires** à tous les niveaux



GAIA DATA, Equipex+ PIA 3 un projet structurant



DATA
TERRA



Pôle National
de Données de Biodiversité



anr®



8

Années pour la
réalisation du projet



21

Organismes
partenaires



62 M €

Budget Total (dont
16, 2 millions
d'euros de l'ANR)



400

Ingénieurs,
scientifiques et
experts

Contribution aux initiatives
internationales et européennes
en **appui aux politiques
publiques de développement
durable**

Accès simplifié aux
données multi-sources

**Interopérabilité des
services**


www.gaia-data.org


**constellation
satellites
optique, radar...**


**Observations
long-terme
In-Situ**




**Variables
bio-geochimiques**


**Modèles,
algorithmes...**

Développement de services
**sur le cycle complet de la
donnée**

Se baser sur les capacités,
institutions, structures et
ressources existantes

**Approches
multidisciplinaires intégrées**
pour l'utilisation des données
de recherche d'observation de
la Terre


MINISTÈRE
DE L'ENSEIGNEMENT
SUPÉRIEUR
ET DE LA RECHERCHE

 **Ouvrir
la science !**



INFRASTRUCTURE GAIA DATA

Architecture de haut niveau



Les outils favoris de l'utilisateur capables d'utiliser les API des collections



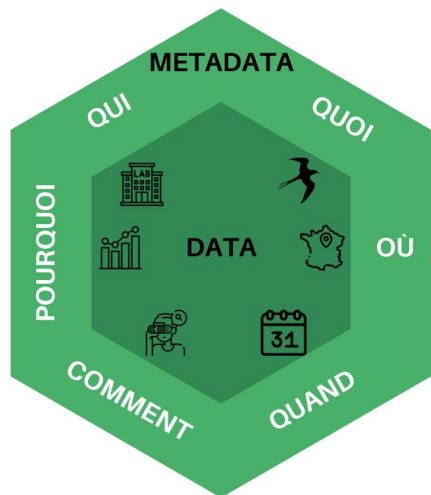
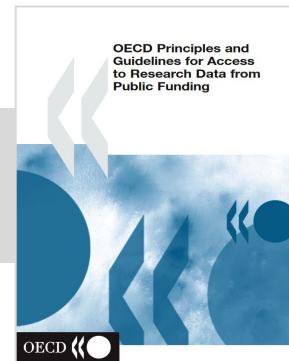
Virtual Research Environment proposé par GAIA DATA pour l'exploration des données



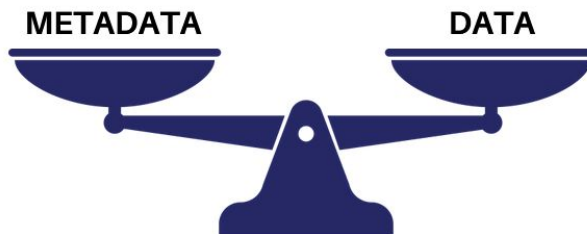
GAIA Dataspace Espace de données pour un projet, un centre, ...

Comprendre les données par les métadonnées : définitions

"Les données de recherche sont définies comme des enregistrements factuels utilisés comme sources primaires pour la recherche scientifique, et qui sont généralement acceptés dans la communauté scientifique comme nécessaires pour valider les résultats de la recherche." **OCDE, 2007.** <https://www.oecd.org/sti/inno/38500813.pdf>



STANDARDISER LES {MÉTA}DONNÉES



Informations brutes vs. spécifiques et dérivées
Connaissance du standard, son formalisme, ses restrictions
Nécessité de connaître le temps pour la standardisation
Diversité des types de données

COMPLÉMENTARITÉ DES DEUX APPROCHES

"Les métadonnées, que l'on peut définir simplement comme « des données sur les données », sont un moyen de nommer les choses et de représenter les données et leurs relations." **Christine L. Borgman, 2020.** <https://books.openedition.org/oep/14692>



Comprendre les données par les métadonnées : *définitions*

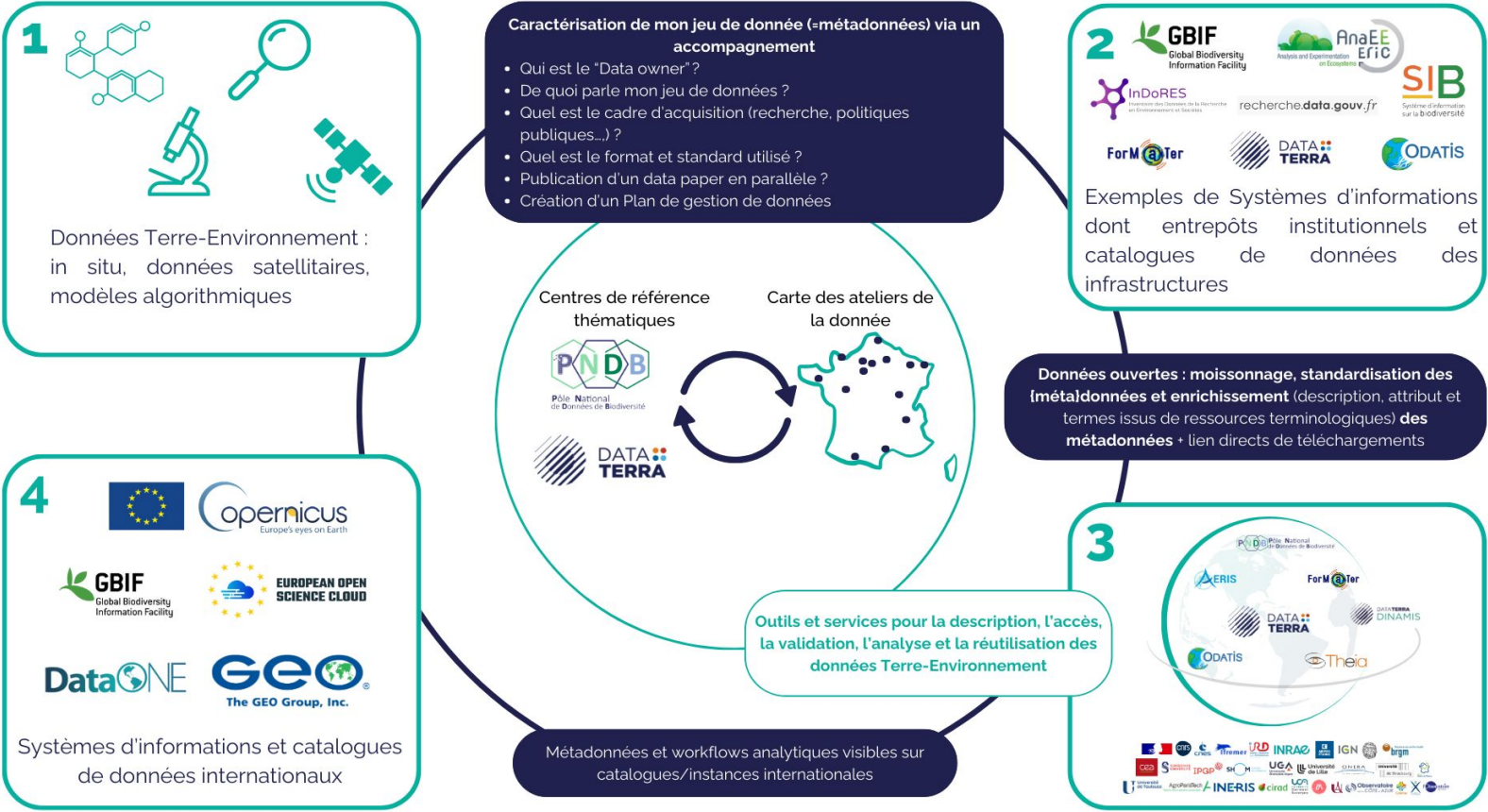
La [gouvernance du PNDB](#), a validé une liste de métadonnées indispensables afin d'avoir un degré de "FAIRitude" (cf. [principes FAIR](#)) minimal et relativement élevé

- **Données ouvertes** (CC-BY 4.0 compatible Etalab)
- **Licence obligatoire**
- **Lien direct de téléchargement** des jeux de données brutes
- **Périmètre thématique** (Toute la biodiversité y compris la paléo- et archéo-biodiversité)
- **Périmètre géographique** (Données produites par la France)
- **Couverture temporelle** (à minima une date d'acquisition de données)
- **Résumé**
- **Titre, auteurs et contacts**
- **Cadre d'acquisition** (a minima via un champ texte)
- **DOI / identifiants uniques**
- **couverture taxonomique** (si présence de taxons)
- **mots clés en lien avec Thesaurus**
- **Attributs des données** (Dictionnaire des attributs de données avec unités et descriptions)
- **Annotation sémantique** (Mots-clés et noms d'attributs, ressources utilisables illimitées)

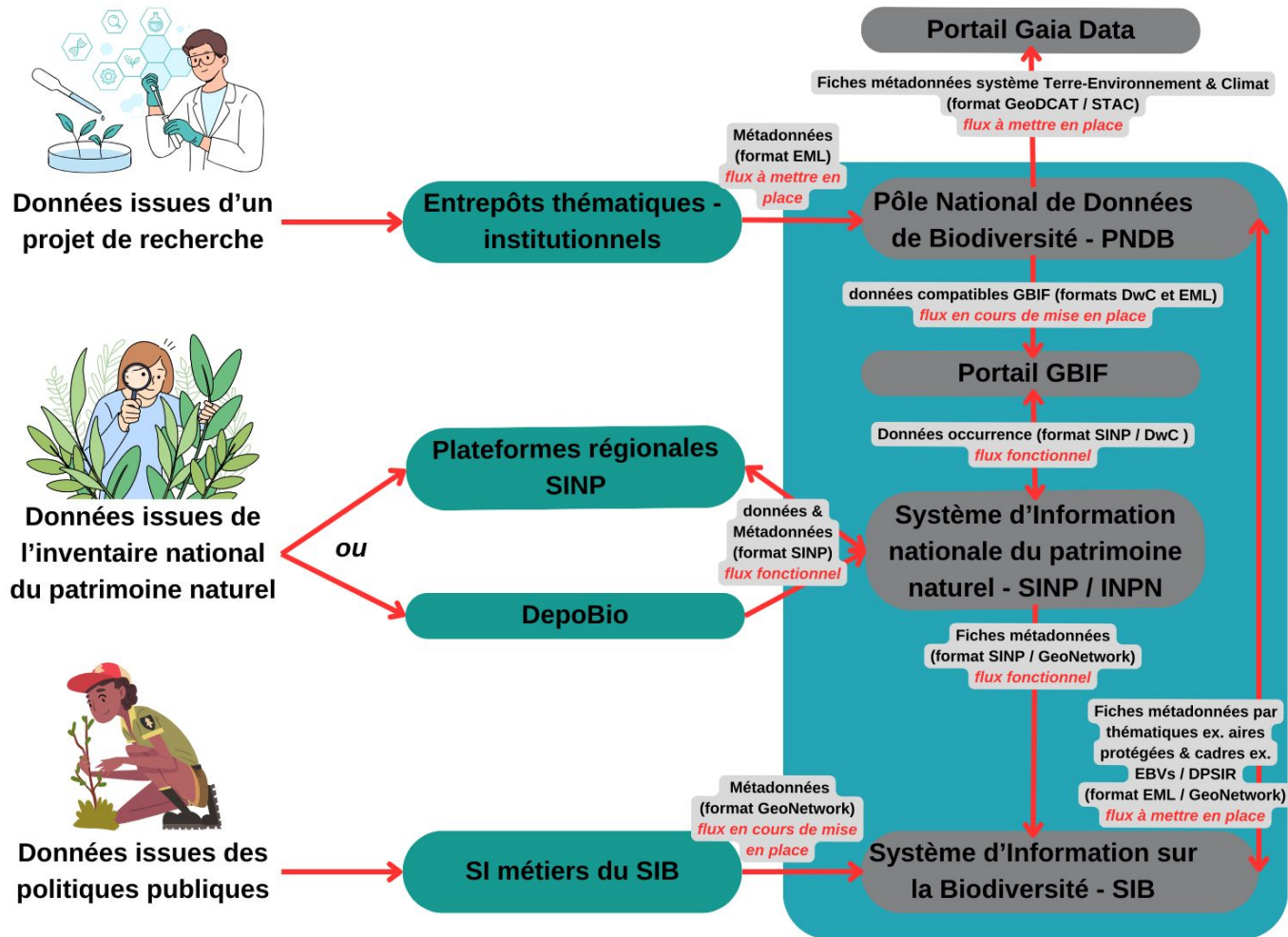


Ces informations minimales sont donc communes à toutes les fiches de métadonnées qui seront présentes dans le [catalogue du PNDB](#)

flux et stock de {méta}données, pour et par les communautés



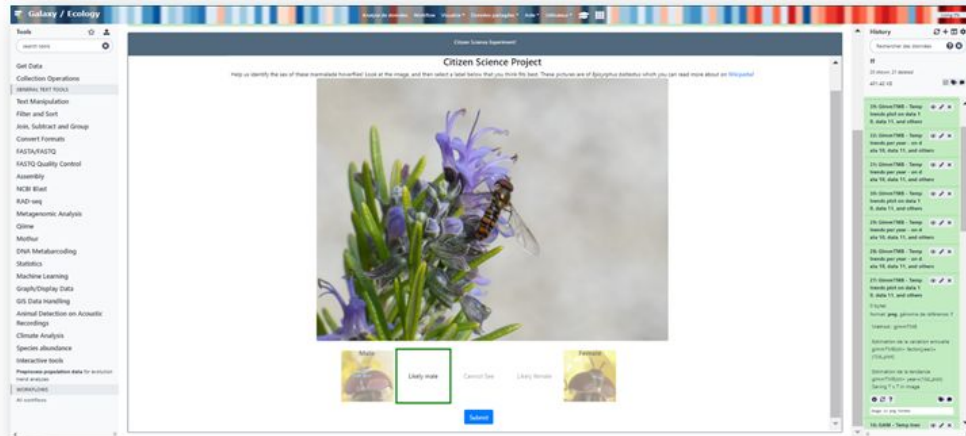
FLUX DES DONNÉES & MÉTADONNÉES DE BIODIVERSITÉ



Annotations et crowdsourcing : travaux préliminaires sous Galaxy



Crowdsourcing image



VIGIENATURE

Un réseau de citoyens qui fait avancer la science



Gamification of participatory science for training and education purposes



Production / SpiPoll Fly Logout

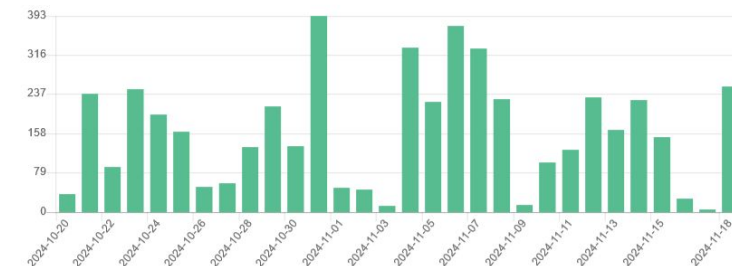
Overview

All time

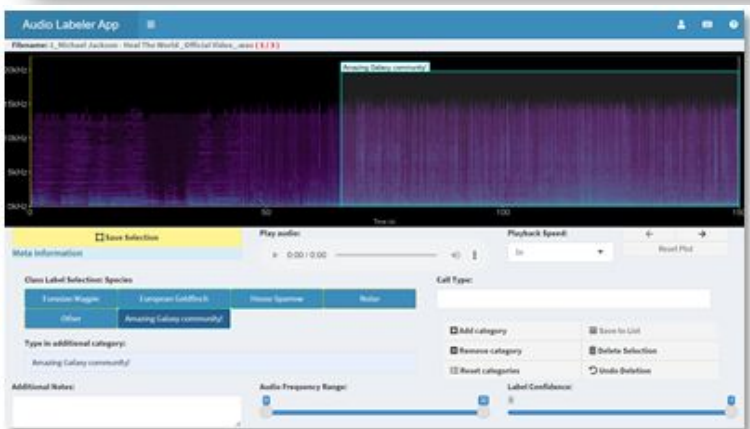


Classifications

Last 30 days



>3400 classifications par mois / >110 par jour



Annotation sons

Annotations et IA : travaux préliminaires sous Galaxy

Moover-inference Automated tracking, segmentation and feature extraction (Galaxy Version 1.0.0)

Tool Parameters

Input video *

69: test.mp4

accepted formats ▾

Use the default model or a custom model

default

Percentage *

50.0

Between 0 and 100, e.g. 50%

save predictions ?

No

extract info ?

No

save annotation ?

No

Run Tool

Help

Moover Inference Tool

This tool performs automated tracking, segmentation, and feature extraction on input videos. It utilizes the Moover-inference lib

Inputs:

- Input video:** Provide the video file on which the analysis will be performed.
- Use the default model or a custom model:** Select whether to use the default model or provide a custom one.
- Model weights:** If custom model is selected, upload the model weights file here.
- Percentage:** Set the confidence threshold percentage between 0 and 100 for the analysis.
- Save predictions?:** Choose whether to save the predictions generated by the tool.
- Extract info?:** Choose whether to extract additional information.
- Save annotation?:** Choose whether to save annotations.

Outputs: **Output file:** The results of the analysis will be provided as a zip file containing various outputs.

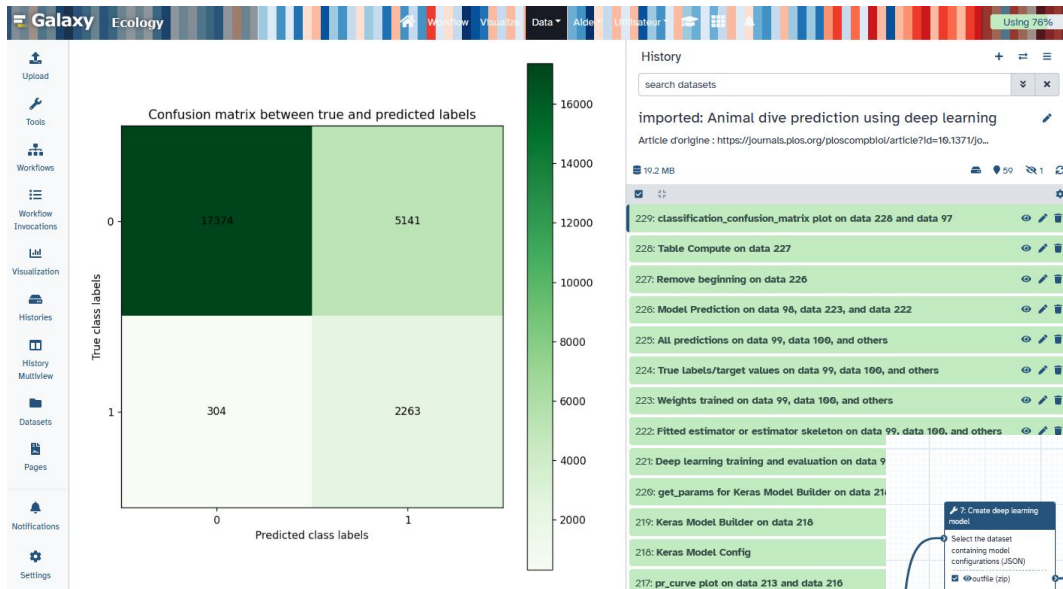
Annotation image ou vidéos



Automated segmentation /
tracking / feature extraction

Entrainement

Annotations et IA : travaux préliminaires sous Galaxy



PLOS COMPUTATIONAL BIOLOGY

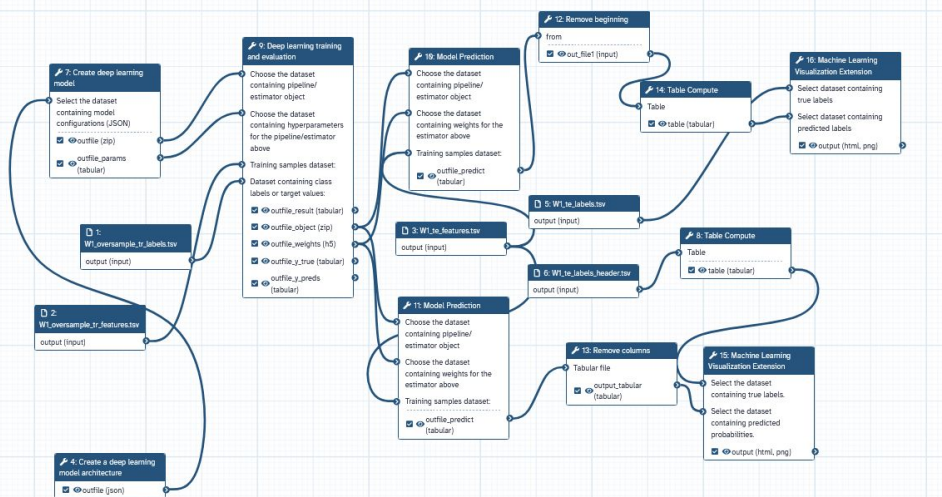
OPEN ACCESS PEER-REVIEWED

RESEARCH ARTICLE

Deep inference of seabird dives from GPS-only records: Performance and generalization properties

Amédée Roy, Sophie Lanco Bertrand, Ronan Fablet

Version 2 Published: March 11, 2022 • <https://doi.org/10.1371/journal.pcbi.1009830>

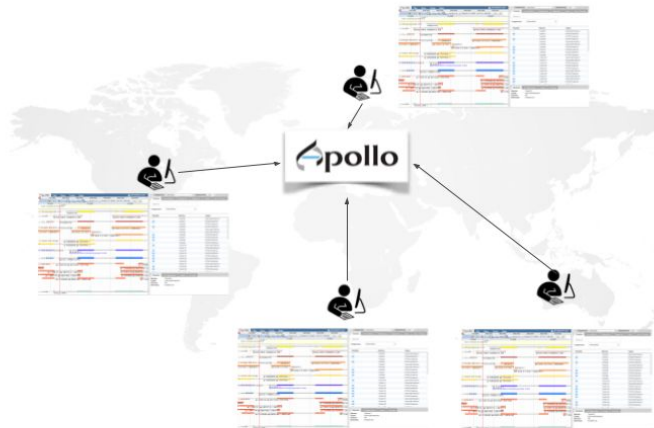


Annotations collaboratives : *Galaxy Genome Annotation as an example*



Galaxy Genome Annotation

A screenshot of the Galaxy web interface showing the annotation workflow for E. coli K12. The main panel displays a genomic track with various annotations, including gene models from Augustus and NCBI AnnotWriter. A context menu is open over a gene, listing options like 'pseudogene', 'tRNA', 'snRNA', etc. The right sidebar shows a 'History' panel with a list of tasks such as 'Annotate on data 23', 'Create or Update Organism on data 22', and 'JBrowse on data 21'. The top navigation bar includes 'Galaxy / Europe', 'Analyze Data', 'Workflow', 'Visualize', 'Shared Data', 'Help', and 'User'.



Classical “workflow”

1. Fetch Data
2. Analyse in Galaxy
3. Send to Apollo
4. Collaboratively Annotate → repeat

Using Galaxy, any user can upload its own organism, and start curating gene models collaboratively with other users.

Check out the [GTN training material](#) to learn how to use it!