# Metadata

## Exercise : using the Ecological Metadata Language - EML through Galaxy-Ecology

July 2025

Yvan Le Bras

Scientific and technical coordinator

@PNDB

@DataTerra

Olivier Norvez

Animation coordinator

@DataTerra

@PNDB

@THEIA
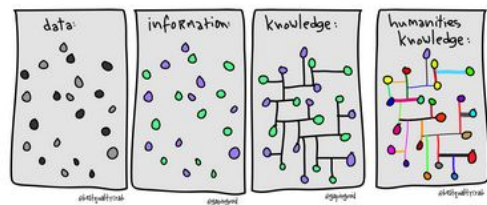
https://biodiversitydata.github.io/
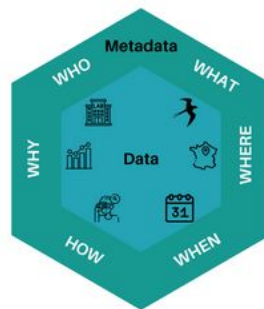
# Tables of contents

- Reminders
  - context and issues
  - EML
  - Galaxy ecology
- Connect to Galaxy Ecology Europe
- Creating an EML metadata record from EML metadata template files in text format
- Creation of EML Assembly Line metadata template files
- Creating tab-delimited text metadata template files from an EML record
- Creating an EML record from an ISO19115 metadata record
- Creation of an EML record for genetic data
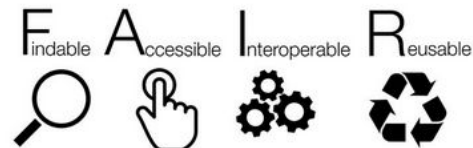- Perspectives

# Reminder : context and issues

- **Heterogeneity** (data types, origin, standards) & **diversity** of "objects" to be linked together[1]
- **Loss of information** over time[2]
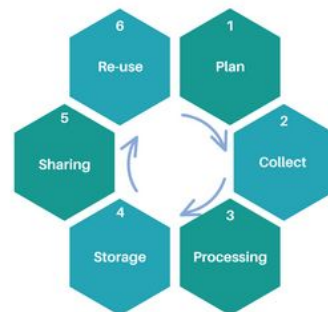- Toward a better **open science** and **reproducibility** [3][4]



WISDOM
KNOWLEDGE
INFORMATION
DATA

data: information: knowledge: humanities knowledge:



BIODIVERSITY (META)DATA DEFINITION

Metadata
WHO    WHAT
WHY    Data    WHERE
HOW    WHEN

CONTINIUM BETWEEN DATA AND METADATA



Findable  Accessible  Interoperable  Reusable

Cycle de vie des données

6 Re-use    1 Plan
5 Sharing    2 Collect
4 Storage    3 Processing

1. Page (2016)
2. Michener et al. (1997)
3. Powers & Hampton (2018)
4. Genkins et al. (2016)

https://biodiversitydata.github.io/

FRB
FONDATION POUR LA RECHERCHE SUR LA BIODIVERSITÉ

CESAB
CENTRE FOR THE SYNTHESIS AND ANALYSIS OF BIODIVERSITY

# Reminder : Ecological Metadata Language



M.B. Jones et al., 2006 https://doi.org/10.1146/annurev.ecolsys.37.091305.110031 + https://eml.ecoinformatics.org/

| Record-level Terms | Dublin Core terms, institutions, collections, nature of data record | Simple Darwin Core (flat) |
|---|---|---|
| Occurrence | evidence of species in nature, observers, behavior, associated media, references. | |
| Event | sampling protocols and methods, date, time, field notes | |
| Location | geography, locality descriptions, spatial data | |
| Identification | linkage between Taxon and Occurrence | |
| Taxon | scientific names, vernacular names, names usages, taxon concepts, and the relationships between them | |
| GeologicalContext | geologic time, chrono-stratigraphy, biostratigraphy, lithostratigraphy | |
| ResourceRelationship | explicit relationships between identified resources (e.g., one organism to another, taxon to location, etc.) | Generic Darwin Core (relational) |
| MeasurementOrFact | measurements, facts, characteristics, assertions, references | |

J. Wieczorek et al., 2012 https://doi.org/10.1371/journal.pone.0029715
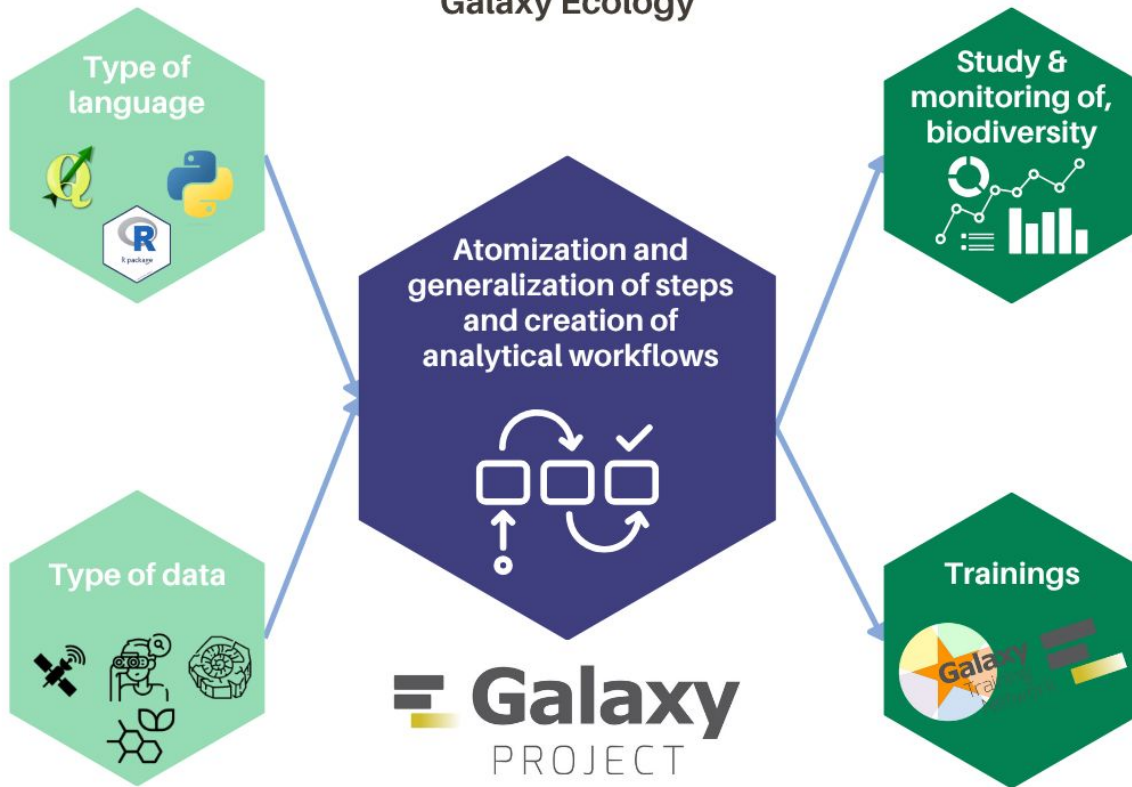
metadata        data

Raw vs. specific and derived information
Knowledge of the standard, its formalism, its restrictions
Need to know the time for standardization
Diversity of data types

**COMPLEMENTARITY OF THE TWO APPROACHES**

< EML />

https://biodiversitydata.github.io/

# Reminder : Galaxy-Ecology



Galaxy Ecology

Type of language

Type of data

Atomization and generalization of steps and creation of analytical workflows

Galaxy PROJECT

Study & monitoring of, biodiversity
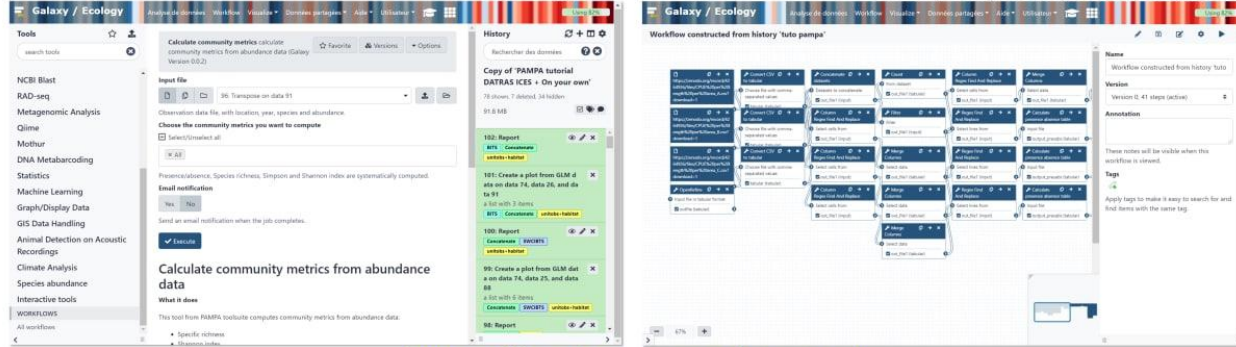
Trainings

See more
- https://www.pndb.fr/pages/galaxy-ecology0/
- Royaux et. al., 2025. Guidance framework to apply best practices in ecological data analysis: lessons learned from building Galaxy-Ecology. Gigascience. https://doi.org/10.1093/gigascience/giae122

https://biodiversitydata.github.io/

# Reminder : Galaxy-Ecology



Use scripts on a cluster ──────→ with 0 programming skills! ←── Create, edit, share, reuse workflows

https://ecology.usegalaxy.eu/

Create & share interactive visualizations,
Even deploy and share Jupyter notebooks or R shiny apps

https://biodiversitydata.github.io/

# Reminder : Galaxy-Ecology



https://biodiversitydata.github.io/

# Connect to Galaxy Ecology Europe

To create a Galaxy Europe account: https://ecology.usegalaxy.eu/login/start?redirect=None  or https://usegalaxy.eu/login/start?redirect=None

To access training resources: https://usegalaxy.eu/join-training/eml

Then go to Galaxy Ecology Europe: https://ecology.usegalaxy.eu/

FRB
FONDATION
POUR LA RECHERCHE
SUR LA BIODIVERSITÉ

CESAB
CENTRE FOR THE SYNTHESIS AND ANALYSIS
OF BIODIVERSITY

# Context and issues

- Rich metadata production (F2 of the FAIR principles, Wilkinson et al., 2016)
- Easy creation/modification of metadata by researchers
  - Metadata inference
  - Data attributes
  - Taxonomic coverage
  - Geographic coverage
  - Personal information
  - Semantic annotations
  - No need to know the specifications of the target standard
  - No coding skills required
  - Web-enabled
  - Possibility of incremental and collaborative development
- Human-machine exchange
- List of minimum criteria chosen by the PNDB

**good Level of FAIRness**

- Open data (CC-BY 4.0 compatible with Etalab)
- Mandatory license
- Direct link to download raw datasets
- Thematic scope (All biodiversity including paleo- and archaeo-biodiversity)
- Geographic scope (Data produced by France)
- Temporal coverage (at least one data acquisition date)
- Abstract
- Title, authors and contacts
- Acquisition framework (at least via a text field)
- DOI / unique identifiers
- taxonomic coverage (if taxa are present)
- keywords related to the Thesaurus
- Data attributes (Dictionary of data attributes with units and descriptions)
- Semantic annotation (Keywords and attribute names, unlimited usable resources)

F indable  A ccessible  I nteroperable  R eusable

PNDB
Pôle National
de Données de Biodiversité

Ouvrir
la science !

FRB
FONDATION
POUR LA RECHERCHE
SUR LA BIODIVERSITÉ

CESAB
CENTRE FOR THE SYNTHESIS AND ANALYSIS
OF BIODIVERSITY

# Creating an EML metadata record from EML metadata template files in text format

Click on Galaxy History with the input data: https://ecology.usegalaxy.eu/u/ylebras/h/eml-assembly-line-tape-1--cration-deml

Steps

- **Import the data, then click on the "home" icon at the top of the page and/or refresh the page.**
- **File review / What information is there (by looking at just the file names first, then browsing the contents of the metadata templates)?**

# Creating an EML metadata record from EML metadata template files in text format



≪ History: Tuto EML Assembly Line sans MetaShARK fin

## Metadata templates original

a list with 16 datasets

⬇ Download

1: taxonomic_coverage.txt 👁

2: abstract.md 👁

3: additional_info.md 👁

4: attributes_02_Ref.txt 👁

5: attributes_datafile_1.txt 👁

6: attributes_datafile_2.txt 👁

https://biodiversitydata.github.io/

# Creating an EML metadata record from EML metadata template files in text format

Taxonomic coverage

# Creating an EML metadata record from EML metadata template files in text format

Attributes



*The missing value code explanation is missing.*

# Creating an EML metadata record from EML metadata template files in text format

Categorical variable

# Creating an EML metadata record from EML metadata template files in text format

Entities to describe the spatial information of raster and vector files

# Creating an EML metadata record from EML metadata template files in text format

Geographic coverage



## 12: geographic_coverage.txt

3 lines

format **tabular**, génome de référence **?**

| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| geographicDescription | northBoundingCoordinate | southBoundingCoordinate | eastBoundingCoordinate | w |
| Biscay | 46.0253 | 46.0253 | -4.8054 | - |
| Channel | 46.4696 | 46.4696 | -4.57 | - |

westBoundingCoordinate is missing

https://biodiversitydata.github.io/

# Creating an EML metadata record from EML metadata template files in text format

Keywords

# Creating an EML metadata record from EML metadata template files in text format

Contacts of persons (data producer, PI, data manager,...)

# Creating an EML metadata record from EML metadata template files in text format

Steps
- Import the data, then click on the "home" icon at the top of the page and/or refresh the page.
- File review / What information is there (by looking at just the file names first, then browsing the contents of the metadata templates)?
- **EML creation from EML Assembly Line template files using the Galaxy Make EML tool**
  - Provide the data package title "EML dataset creation from EML Assembly Line files"
  - Temporal coverage "2021-01-01" to "2021-12-12"
  - Select the data collections of type "dataTable", "spatialRaster", and "spatialVector"
  - The license is automatically CC-BY 4.0 compatible with Etalab 2.0 open license.
    - make eml /
    - What was written in the metadata?
    - EML validated?

# Creating an EML metadata record from EML metadata template files in text format

# Creating an EML metadata record from EML metadata template files in text format

Steps
- Import the data, then click on the "home" icon at the top of the page and/or refresh the page.
- File review / What information is there (by looking at just the file names first, then browsing the contents of the metadata templates)?
- EML creation from EML Assembly Line template files using the Galaxy Make EML tool
- **MetaSHRIMPS: FAIRness assessment of metadata and creation of a draft data paper**
  - [https://ecology.usegalaxy.eu/root?tool_id=interactive_tool_metashrimps](https://ecology.usegalaxy.eu/root?tool_id=interactive_tool_metashrimps)  on the generated EML
  - Execute => Draft of data paper + FAIR assessment
    - The draft "Data Paper" allows you to see the different metadata elements presented in the form of a static web page, a bit like in a data/metadata catalogue but without waiting for sending to the catalogue administrators and putting online!
    - We can see that some attributes have definitions added manually, while others have been automatically populated by EML Assembly Line
      - We could modify them afterwards to improve this
    - You can download an editable version of the draft data paper, to have a shareable or modifiable word version that can serve as a basis for writing an internal document or for wide distribution such as a data paper.
    - FAIRness score 30 success / 10 failure / 5 warning
      - On the "warning" side, we note a point on the size of the summary which would facilitate the "findable" aspects, another on the definition of an attribute "Present.Surface.pH" which only contains 3 words and would improve the "reusable" aspects.
      - On the "failure" side, it's worth noting the lack of semantic annotation.

https://biodiversitydata.github.io/

# Creating an EML metadata record from EML metadata template files in text format

Steps
- Import the data, then click on the "home" icon at the top of the page and/or refresh the page.
- File review / What information is there (by looking at just the file names first, then browsing the contents of the metadata templates)?
- EML creation from EML Assembly Line template files using the Galaxy Make EML tool
- MetaSHRIMPS: FAIRness assessment of metadata and creation of a draft data paper
- **How to improve the files and get a better EML?**
  - Added annotation (EML report) + modified abstract and attributes (MetaShRIMPS report)

FRB
FONDATION
POUR LA RECHERCHE
SUR LA BIODIVERSITÉ

CESAB
CENTRE FOR THE SYNTHESIS AND ANALYSIS
OF BIODIVERSITY

# Creating an EML metadata record from EML metadata template files in text format

Steps
- Import the data, then click on the "home" icon at the top of the page and/or refresh the page.
- File review / What information is there (by looking at just the file names first, then browsing the contents of the metadata templates)?
- EML creation from EML Assembly Line template files using the Galaxy Make EML tool
- MetaSHRIMPS: FAIRness assessment of metadata and creation of a draft data paper
- How to improve the files and get a better EML?
- **Edition de fichiers de template EML Assembly Line pour améliorer le score de FAIRitude**
  - editing the "abstract.txt" file to obtain a summary of more than 100 words, of the "attributes_Present.Surface.pH.txt" file to expand the description of the "Present.Surface.pH" attribute ("surface present pH" instead of "ph" for example).
  - Creation of an "annotations.txt" file to add a keyword to the metadata record, "is about" "biodiversity"
  - => We recreate an EM
  - => We re-evaluate the "FAIRness" of our metadata via MetaShRIMPS (33 successes / 9 failures / 3 warnings)

# Creating an EML metadata record from EML metadata template files in text format

You can create a workflow from your history by clicking on the "gear" in the top right corner of the history, then "extract workflow".

https://ecology.usegalaxy.eu/u/ylebras/w/workflow-constructed-from-history-tuto-eml-assembly-line-tape-1--cration-deml



=> But how do you create these famous template files?

# Creation of EML Assembly Line metadata template files

Galaxy history with input data:
https://ecology.usegalaxy.eu/u/ylebras/h/eml-assembly-line-tape-2--cration-des-fichiers-modeles-de-mtadonnes

- Galaxy EML Assembly Line tools for generating:
  - Lists of attributes
    - For data tables
    - For raster GIS files
    - For vector GIS files
      - Surval data from the Quadrige database retrieved on November 30, 2021 for the REPHY program
  - Geographic coverage
    - If this does not appear to work, carefully review the error message (by clicking on the "insect" icon in the bottom left of the data file preview) to determine the source of the problem and suggest a solution.

# Creation of EML Assembly Line metadata template files

- Galaxy EML Assembly Line tools for generating:
  - Lists of attributes
  - Geographic coverage
  - **Taxonomic coverage**
  - Regarding the Temporal Coverage => as previously mentioned, this must be provided when creating the EML, along with the title
- We are seeing missing information in the metadata template, which explains why not all data files are populated in the EML.
  - Find and fill in the missing information by navigating to the "Warning message" file.

# Creation of EML Assembly Line metadata template files

https://ecology.usegalaxy.eu/u/ylebras/w/copy-of-workflow-constructed-from-history-tuto-eml-assembly-line-tape-1--cration-deml



https://biodiversitydata.github.io/

# Creating tab-delimited text metadata template files from an EML record

- From the Kakila database's EML metadata record:
  https://data.pndb.fr/view/doi%3A10.48502%2F8bb5-pk85
  - Copy and paste the direct link URL to the metadata record into the
    Galaxy upload module / "paste/fetch data" field:
    https://pndb.fr/metacat/d1/mn/v2/object/doi%3A10.48502%2F8bb5-pk8
    5 Product analysis history
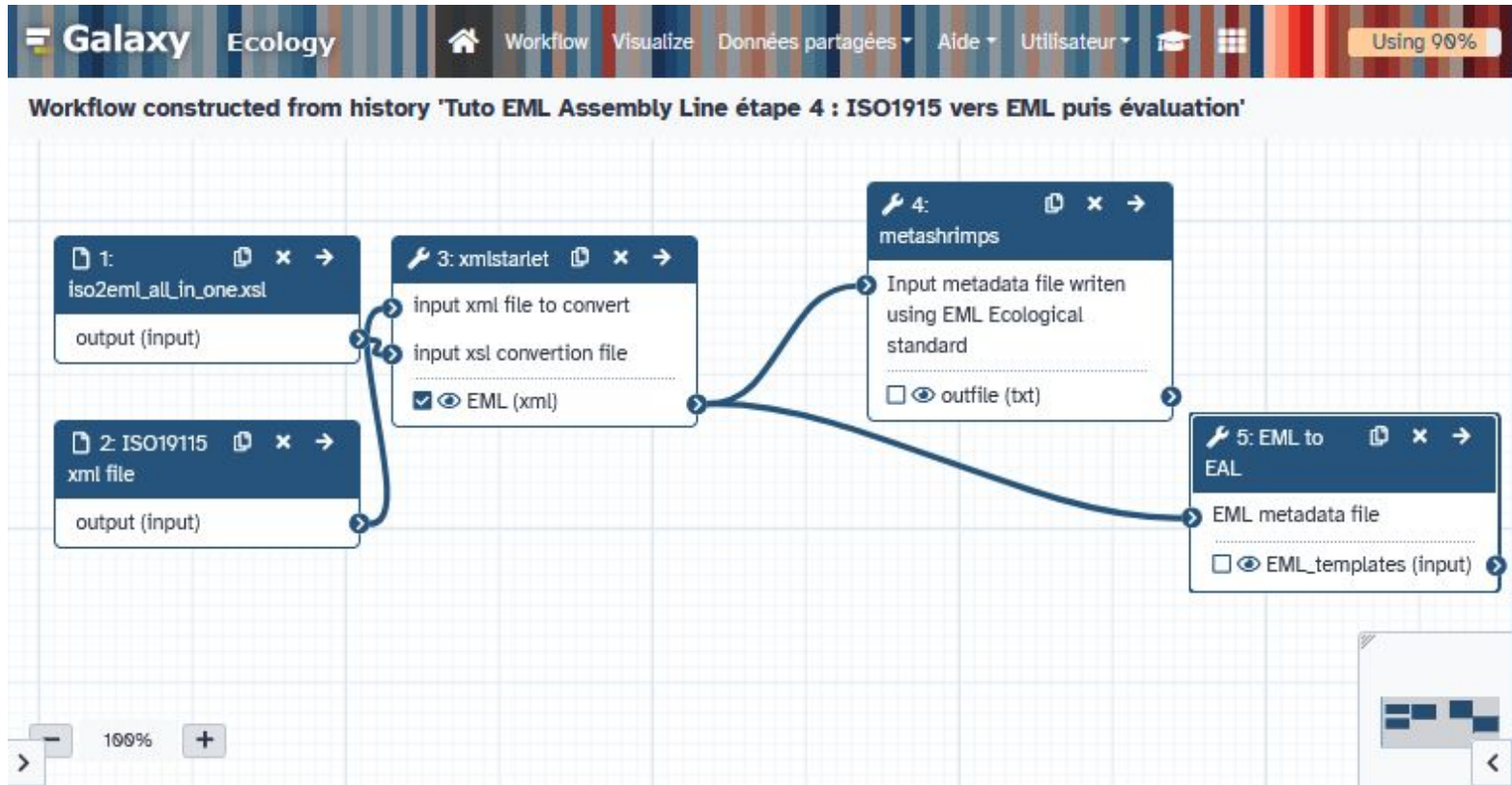- https://ecology.usegalaxy.eu/u/ylebras/h/tuto-eal-3-eml-to-eal

# Creating an EML record from an ISO19115 metadata record

Galaxy history with input data:
https://ecology.usegalaxy.eu/u/ylebras/h/tuto-eml-assembly-line-tape-4--iso1915-vers-eml-puis-valuation-fairitude-et-bauche-de-data-paper

- Using the "xmlstarlet" tool to convert the SURVAL ISO metadata sheet accessible from this webpage https://surval.ifremer.fr/Donnees/Cartographie-Inventaire-du-reseau-REPHY#/metadata/aa8fe568-d2c0-4b53-a8bb-d9fcef2b5293 (direct link to metadata sheet: https://sextant.ifremer.fr/geonetwork/srv/api/records/aa8fe568-d2c0-4b53-a8bb-d9fcef2b5293/formatters/xml ) to EML using the conversion file "iso2eml_all_in_one.xsl"
- Creation of a "FAIritude" report and a draft "data paper" via "MetaSHRIMPS"
- Creation of metadata template files in tab-delimited text format via the "EML to EAL" tool

FRB
FONDATION
POUR LA RECHERCHE
SUR LA BIODIVERSITÉ

CESAB
CENTRE FOR THE SYNTHESIS AND ANALYSIS
OF BIODIVERSITY

# Creating an EML record from an ISO19115 metadata record
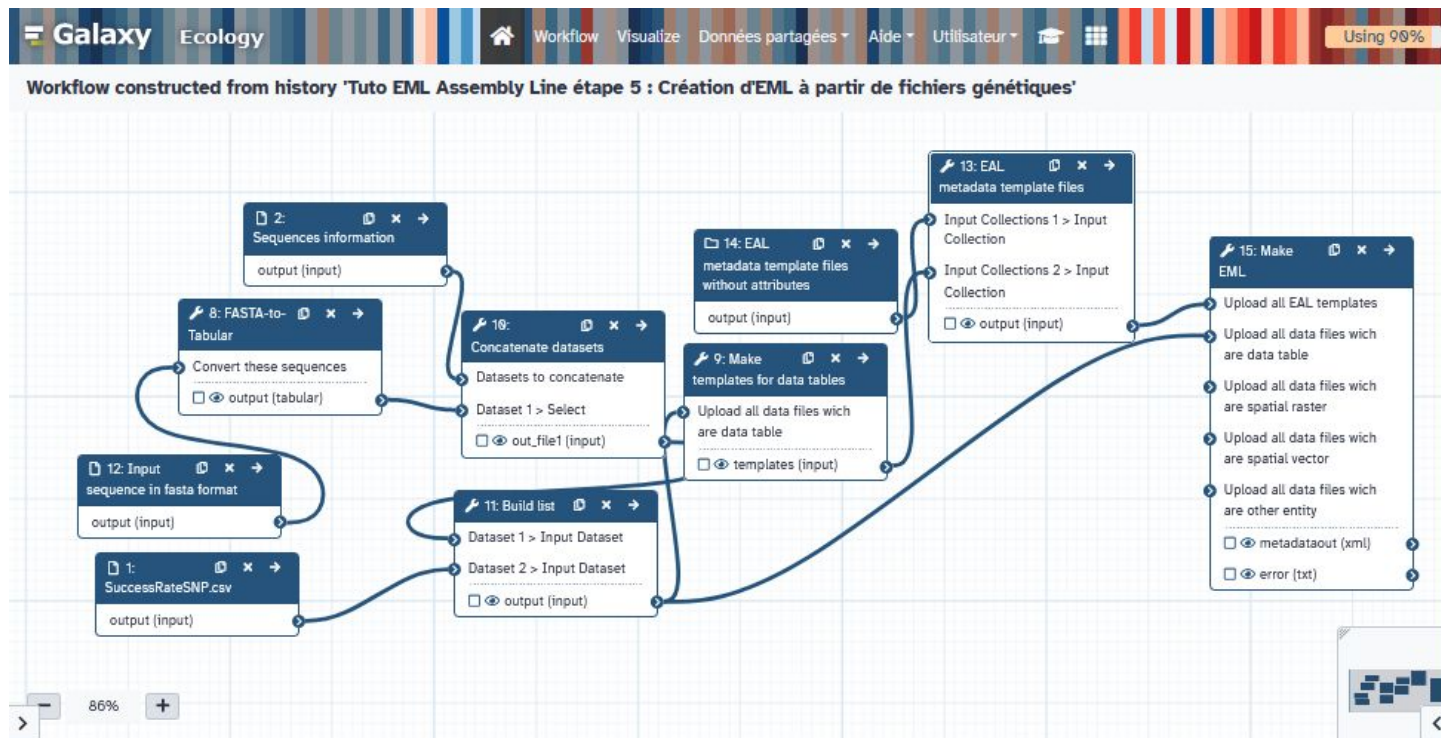
# Creation of an EML record for genetic data

Galaxy history with input data:
https://ecology.usegalaxy.eu/u/ylebras/h/tuto-eml-assembly-line-tape-5--cration-deml--partir-de-fichiers-gntiques

- Using SEANOE data DOI 10.17882/70546 https://www.seanoe.org/data/00593/70546/  Le Cam et al 2019
- From sequence data in FASTA format, identify a natural segmentation of sequence names into information elements, which we will call metadata elements, and create a tabulated file from this sequence file with one column for the sequence and as many columns as there are metadata elements.
- Looking at the FASTA file, we notice that there appear to be 7 metadata elements:
  - >RCL_P_1 ACATAACTTAGAGAAAGGAGCAGAGATCAGGGAAGGGGACAGCAACAAAG 2705094059 RCL_P_1_0_T_U_2705094059 23/07/2018 InfiniumII 1
- From the sequence file, we can therefore create a tabulated file using the FASTA TO TABULAR tool, ultimately obtaining a file with 8 columns.
- Next, you can add a header to the file using the Sequences Information history file and the "Concatenate dataset tail to end" tool.
- Then, you can create EAL metadata template files of type DataTable for this tab-delimited "SequencesSNP_Raja_clavata" file and the "SuccessRateSNP.csv" file, and then modify the content using the metadata provided in the "68857.txt" file, particularly regarding attribute descriptions.
- You can then use the content of the SEANOE metadata record "seanoe_metadata_export_20240312164552.txt" to populate the necessary information in the EAL template files. This step can be done entirely manually, but it is also possible to create a workflow to automate it using this file. If you choose this option, you will need to be or become a Galaxy Jedi, but this will allow you to share it with others and offer an automated workflow to go from a SEANOE txt quote file to EAL template files.

https://biodiversitydata.github.io/

# Creation of an EML record for genetic data

https://biodiversitydata.github.io/

# Perspectives

- **Transformation and Enrichment**
  - => Offers a simplified and collaborative way to share, convert, and enrich metadata
  - Identified need: Work on mappings and conversions between standards!
- **Automation by humans and/or machines**
  - => Use of the Galaxy API by external services, in particular
- **Populating data warehouses with rich metadata** without modifying the software used by the warehouses
  - => Using EAL template metadata files