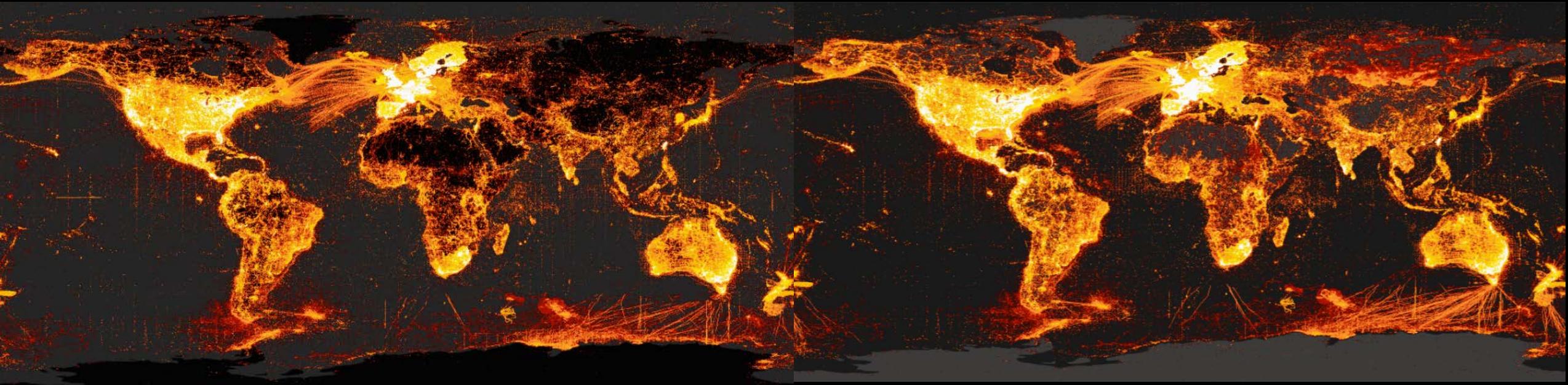


2018



2020



981,464,491



1,602,089,347



## **DNA-DERIVED DATA IN GBIF GUIDANCE AND TRAINING**

**Dmitry Schigel | Scientific officer**

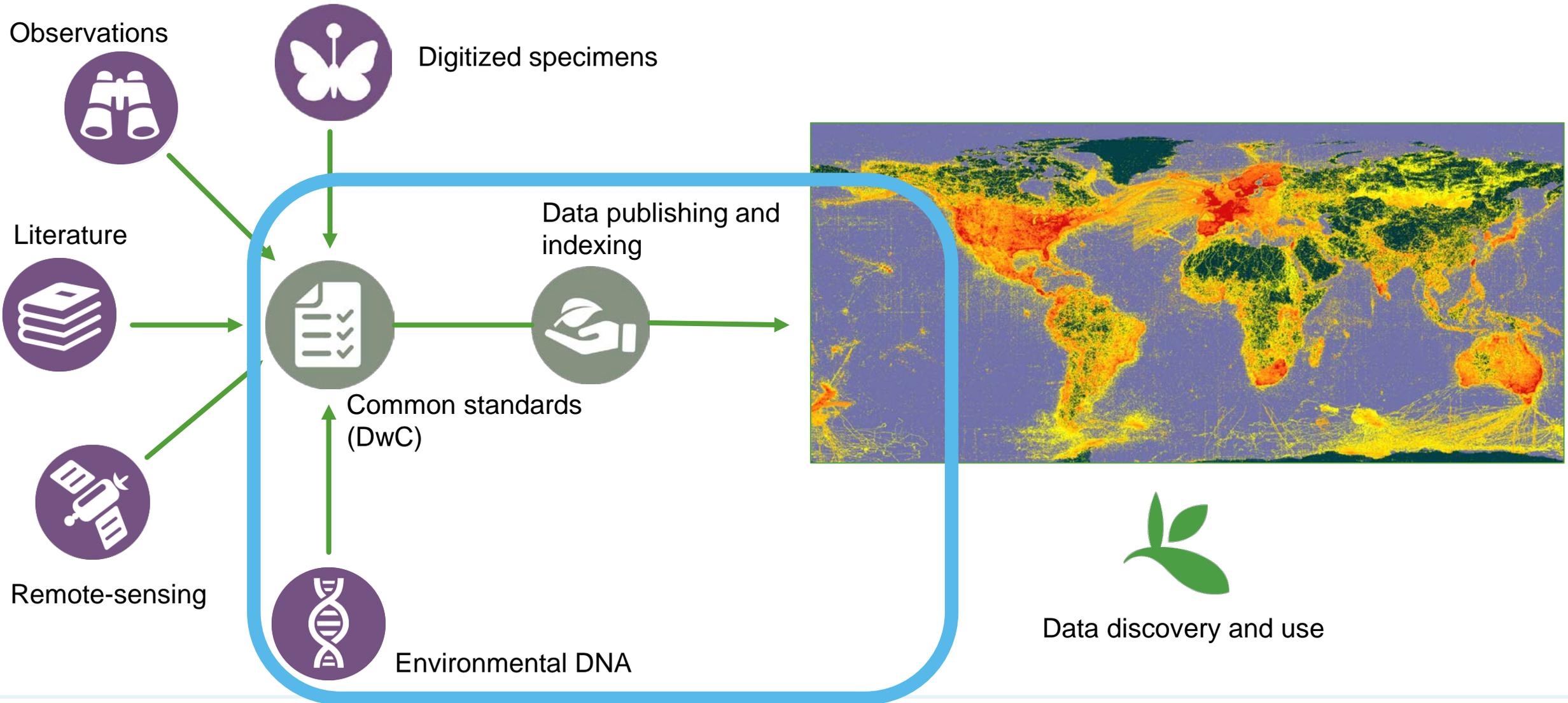


**GBIF**

Global Biodiversity  
Information Facility

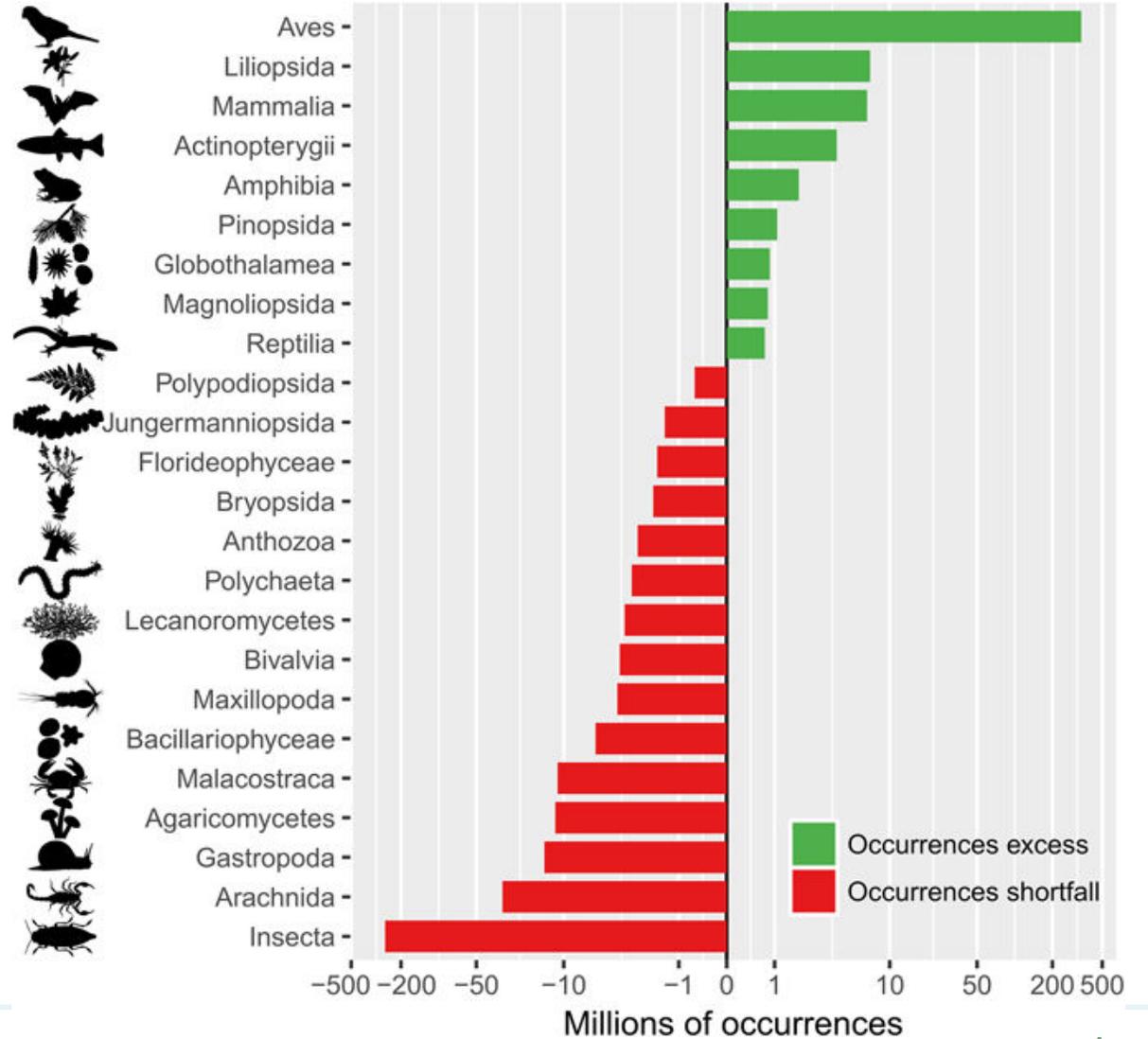
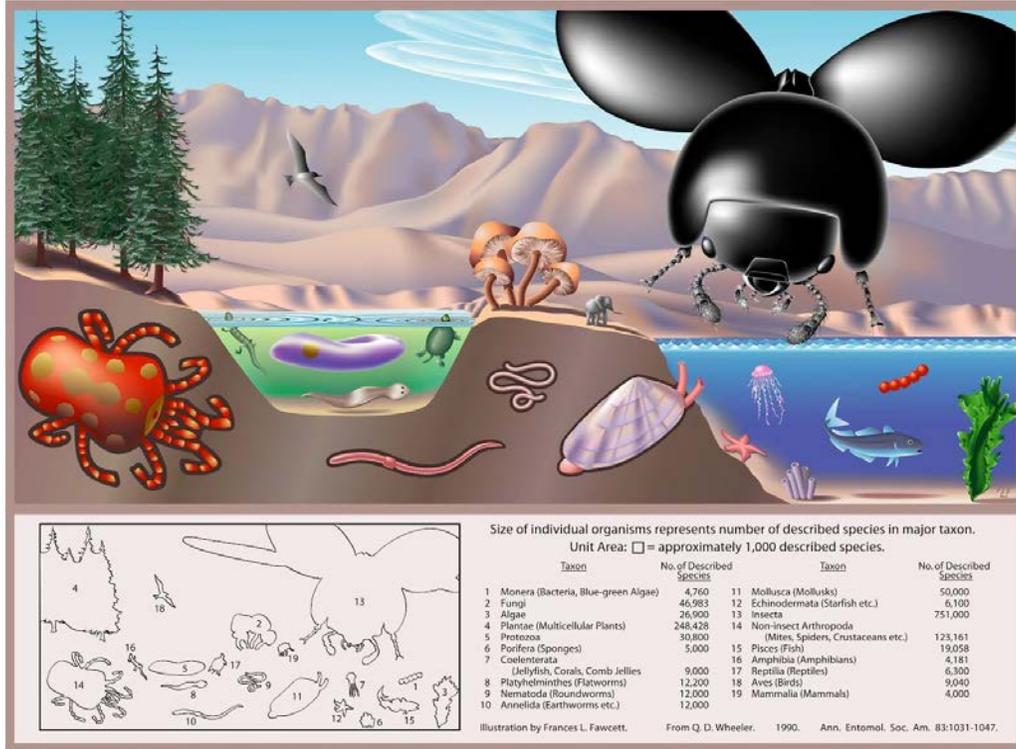
October 2020

# A WINDOW ON EVIDENCE ABOUT WHERE SPECIES HAVE LIVED, AND WHEN



# GLOBAL BIODIVERSITY VS. DIGITALLY AVAILABLE DATA

Image: FL Fawcett in Wheller Ann. Entomol. Soc. Am. 1990



Troutet et al. Nature Scientific Reports 2017

**DATA STREAMS IN GBIF**

**DNA-based evidence**

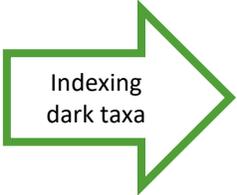
natural history collections

human observations etc.

**Sequence names**

“taxon hypotheses”, OTUs

- UNITE: SH
- BOLD: BINs
- SILVA
- ...



**DNA occurrences**



Latin names

# DATA STREAMS IN GBIF

## DNA-based evidence

natural history collections

human observations etc.

### Sequence names

"taxon hypotheses", OTUs



### DNA occurrences



### Latin names

*sequencing*

individuals / isolates

species mixes (samples)

Amplicon  
(barcoding)  
*one or more  
marker genes*

Genomic  
*shotgun  
DNA*

Transcriptomic  
*shotgun  
RNA*

Amplicon  
(metabarcoding)  
*one or more  
marker genes*

Metagenomic  
*shotgun  
DNA*

Metatranscriptomic  
*shotgun  
RNA*



**DATA STREAMS IN GBIF**

**DNA-based evidence**

natural history collections

human observations etc.

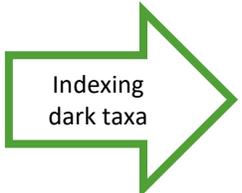
**Sequence names**

“taxon hypotheses”, OTUs

UNITE: SH   BOLD: BINS   SILVA   ...

**DNA occurrences**

Latin names



*sequencing*

individuals / isolates

species mixes (samples)

Targeted detection  
qPCR

Norway   ?

Amplicon  
(barcoding)  
*one or more  
marker genes*

INSDC / ENA

BOLD occurrences   UNITE   ...

Genomic  
*shotgun  
DNA*

Transcriptomic  
*shotgun  
RNA*

Amplicon  
(metabarcoding)  
*one or more  
marker genes*

MGnify

UNITE   BIOWIDE

Metagenomic  
*shotgun  
DNA*

Metatranscriptomic  
*shotgun  
RNA*

MGnify

... (dashed box)

**DATA STREAMS IN GBIF**

**DNA-based evidence**

natural history collections

human observations etc.

CHECKLISTS

CHECKLISTS

**Sequence names**

“taxon hypotheses”, OTUs

UNITE: SH   BOLD: BINs   SILVA   ...

**DNA occurrences**

**Latin names**



*sequencing*

individuals / isolates

species mixes (samples)

Targeted detection  
qPCR

Norway   ?

Amplicon  
(barcoding)  
*one or more  
marker genes*

INSDC / ENA

BOLD occurrences   UNITE   ...

**OCCURRENCE DATASETS**

Genomic  
*shotgun  
DNA*

Transcriptomic  
*shotgun  
RNA*

Amplicon  
(metabarcoding)  
*one or more  
marker genes*

Metagenomic  
*shotgun  
DNA*

Metatranscriptomic  
*shotgun  
RNA*

MGnify   MGnify

UNITE   BIOWIDE   ...

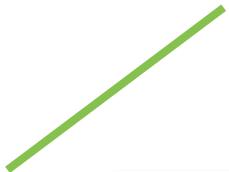
**SAMPLING EVENT DATASETS**

# OPERATIONAL TAXONOMIC UNITS



GBIF

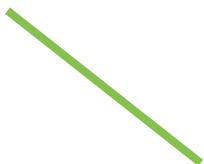
backbone taxonomy



SH ABC0001



OTU = SH,  
Species hypothesis



BIN DEF0002



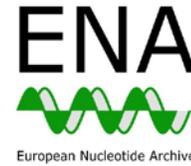
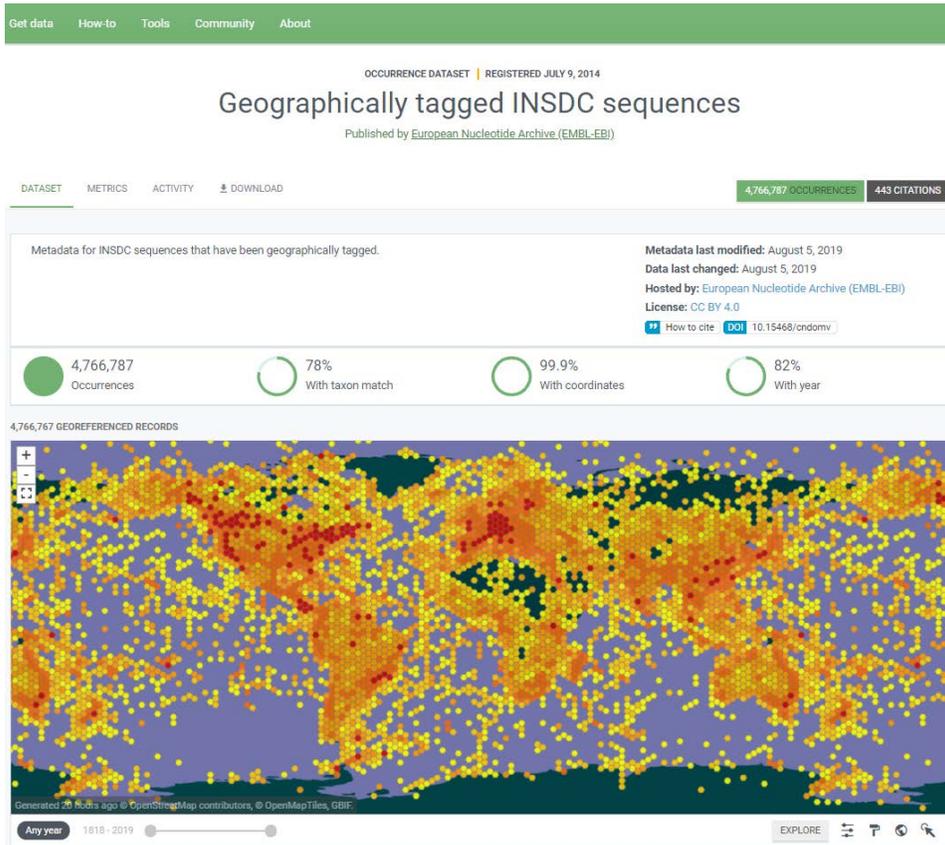
international  
BARCODE  
OF LIFE



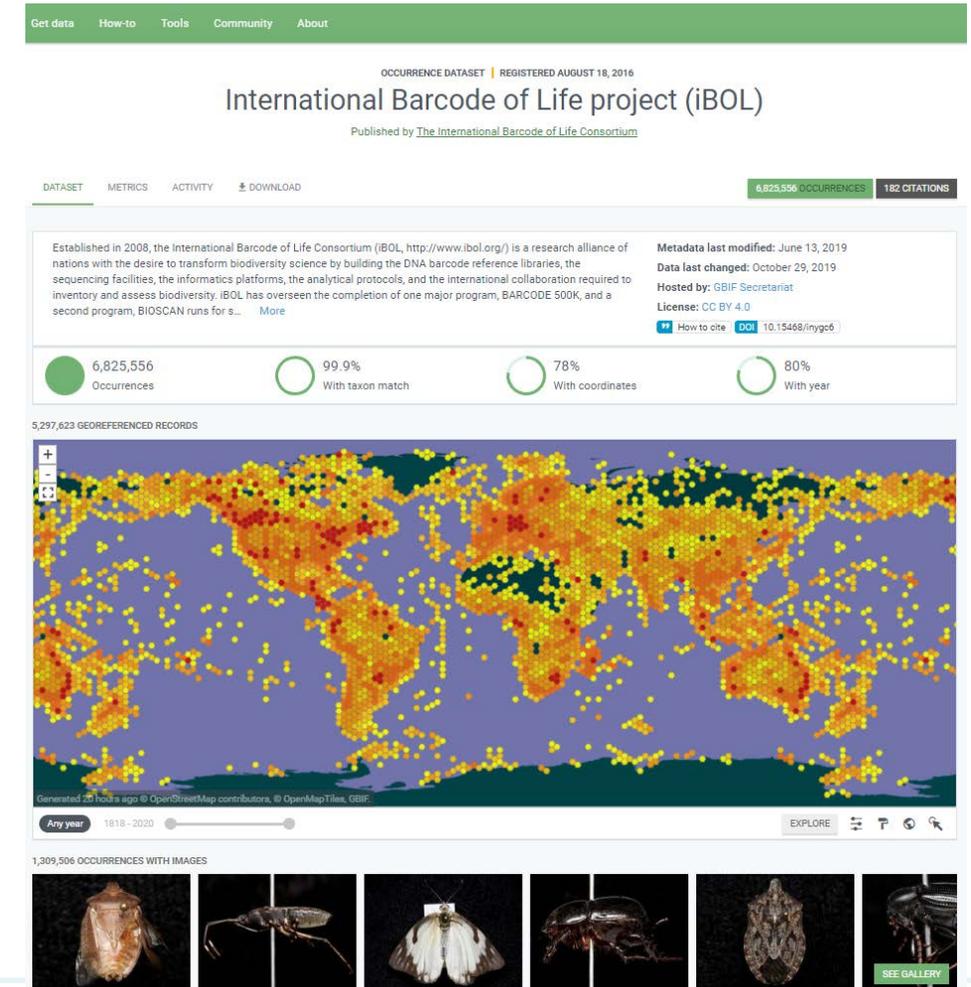
OTU = BIN,  
Barcode  
identification number

# INDIVIUDUAL SEQUENCES WITH COORDINATES

## European Nucleotide Archive: 4.8M records



## International Barcode of Life: 6.8M records





# NEW GUIDE: PUBLISHING DNA-DERIVED DATA THROUGH BIODIVERSITY DISCOVERY PLATFORMS

## Publishing DNA-derived data through biodiversity data platforms [Community review draft]

Anders F. Andersson · Andrew Bissett · Anders G. Finstad · Frode Fossey · Marie Grosjean · Michael Hope · Thomas S. Jeppesen · Urmas Kõljalg · Daniel Lundin · R. Henrik Nilsson · Maria Prager · Cecilie Svenningsen · Dmitry Schigel – Version 90868d9, 2020-10-15 12:04:43 UTC

This document is also available in [PDF format](#).

Mapping and data publishing

Cross-platform

About 40 pages long "cookbook"

- ❖ Introduction
- ❖ Categorization
- ❖ Mapping\*
- ❖ Visuals
- ❖ Future prospects
- ❖ Resources

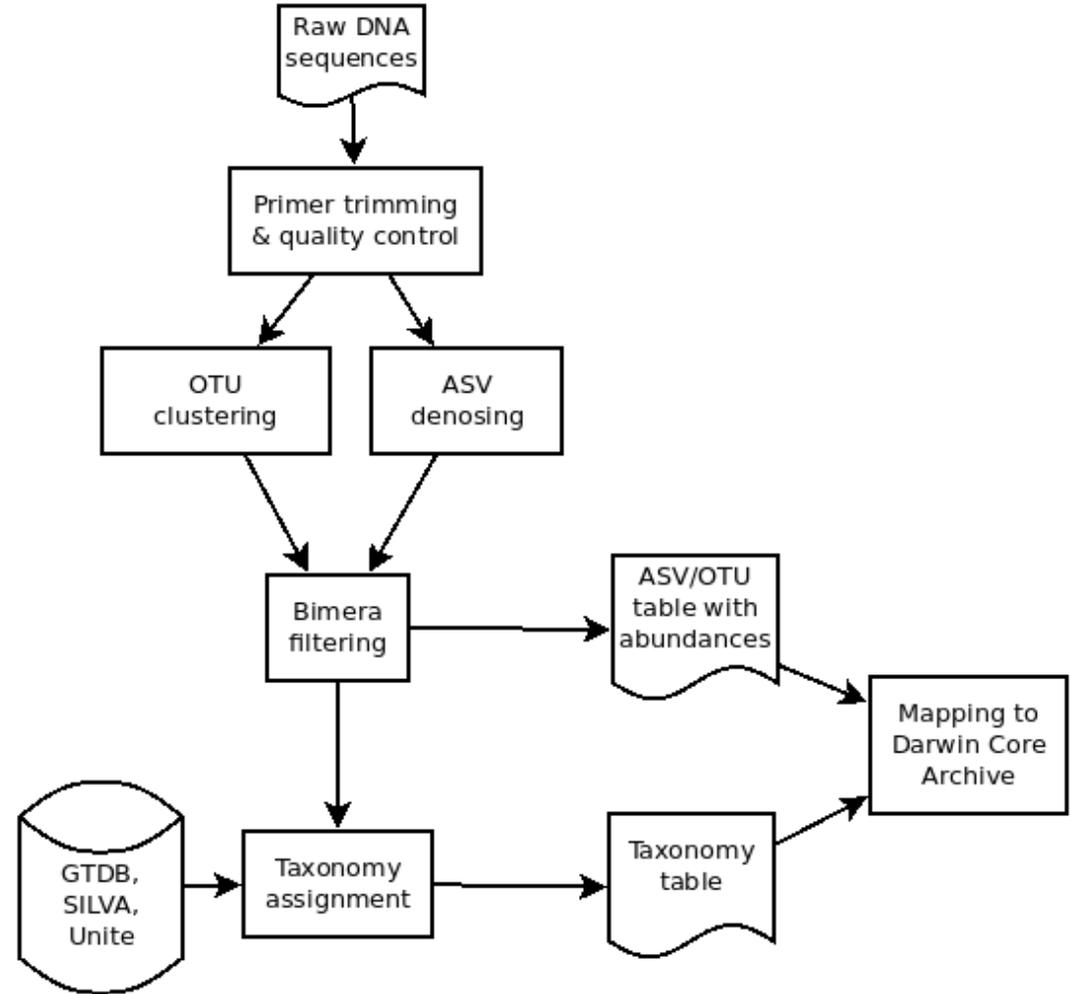
*\*Darwin Core and MIxS based*



# PUBLISHING SEQUENCE-DERIVED DATA: THE “LEARN” SECTION

## Introduction

- ❖ Rationale
- ❖ Audiences
- ❖ DNA derived occurrence data
- ❖ Biodiversity data
- ❖ Processing workflows
- ❖ Taxonomy of sequences



# PUBLISHING SEQUENCE-DERIVED DATA: THE “DO” SECTION

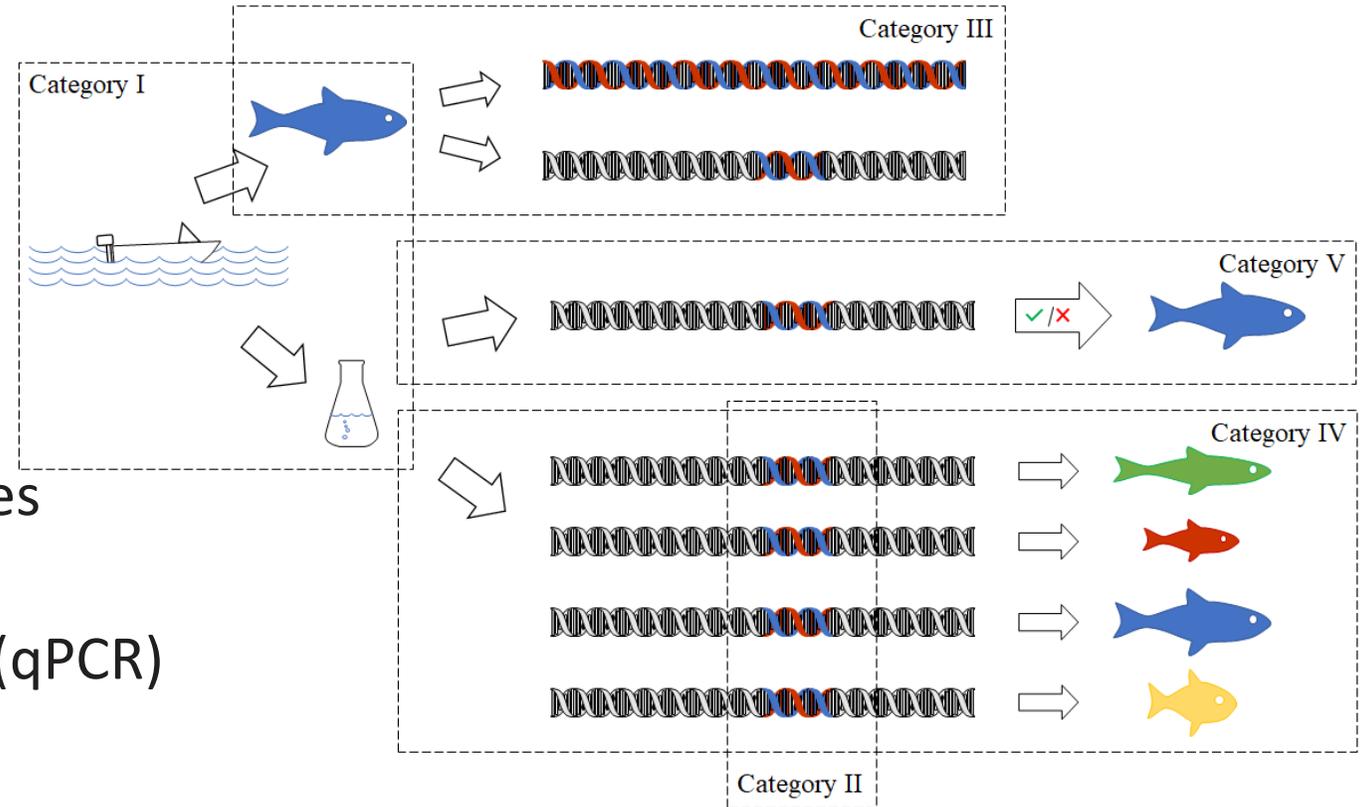
Data packaging and mapping

Categorisation of data

- I Sequence-based occurrences
- II Enriched occurrences
- III Targeted species detection (qPCR)
- IV Name references
- V Metadata-only

Mapping

Examples



# NEW TRAINING COURSE ON GBIF AND BOLD SKILLS – WITH BIODATA /GBIF NORWAY

Home UiO > Natural History Museum

For employees Norwegian website Search

UiO Natural History Museum

Home Visiting Research and collections About the museum

Research and collections

Research projects

BioDATA

Activities

People

## Data publication course in Tbilisi

Planned BioDATA data publishing course in Tbilisi to focus on publishing molecular data in GBIF. Notice that the course is postponed because of the COVID-19 travel recommendations - the dates will be updated.

Time and place: Jan 1, 2021, Tbilisi, Georgia  
Add to calendar



Red Cabbage Bug (*Eurydema ornata*) CC-BY-ds1001

**DATES**  
GBIF / BOLD course on data management skills was originally scheduled Tuesday 7 July - Friday 10 July (four full days from 9:00 till 17:00). This was immediately after CaBOL meeting: Tuesday, June 30 - Friday, July 3 and BioBlitz: Friday 3 July (evening) - Monday 6 July. Notice that the course is postponed because of the COVID-19 travel recommendations - the dates will be updated.

**SCOPE**  
Data management skills for publishing data through BOLD and GBIF data platforms. This is an observation/specimen -> published record course that does not include wet lab steps. Examples of the past, to be modified:  
2018 [https://www.forbio.uio.no/events/courses/2018/Data\\_mobilization\\_Baikal.html](https://www.forbio.uio.no/events/courses/2018/Data_mobilization_Baikal.html)  
2019 <https://sisu.ut.ee/publish-sequence-datasets>

**Organizer**  
BioDATA, GBIF and ForBio

Accelerating biodiversity research through DNA barcodes, collection and observation data

## Work in progress

*Why publish in GBIF?*

*Barcoding, data and biodiversity*

*Biodiversity databases in Georgia*

*Open data as a first-class research citizen*

*Databases as a research tool*

*Why do we need to identify species*

*Concept of DNA barcoding*

*Data in and data out: recognize and understand your data*

*Principles of data organization and personal data management*

*Data structure: standards*

*Data citation*

*Data exposure: why and when*

*Data papers*

*Barcode reference repositories*

*Quality control*

*Use case Bombus and legumes: uncovering pollination mysteries*

## KEY FACTS

- New GBIF course
- Tbilisi, Georgia
- 2021, 4 days
- Onsite or virtual
- GBIF and BOLD
- publish and use
- DNA-derived FAIR data
- Learn and practice

# THANK YOU

Dmitry Schigel

[dschigel@gbif.org](mailto:dschigel@gbif.org)

