



**Indexing for Life**

# **D6.1 - Integrated access to CoL in GBIF**

Workpackage 6

Markus Döring

30<sup>th</sup> April 2013

Capacities Programme of Framework 7: EC e-Infrastructure Programme – Virtual Research  
Communities - INFRA-2010-2

Grant Agreement No:	261555
Project Co-ordinator:	Dr Alastair Culham
Project Homepage:	<a href="http://www.i4Life.eu">http://www.i4Life.eu</a>
Duration of Project:	36 months
Start Date:	November 2010
End Date:	October 2013



## Introduction

The Catalogue of Life has been included in various parts of the GBIF infrastructure to provide integrated access to the catalogue's taxonomy and names through the GBIF systems. An automatic GBIF workflow has been established that downloads and integrates the latest published catalogue on a regular basis.

## GBIF Background

The GBIF infrastructure is made up of various distinct components of which the most important ones are quickly described here. The Catalogue of Life has a role in each of these.

### Registry

A main role of the GBIF registry is to keep track of datasets and provide some basic metadata about them. It is meant to be a public service, but registered service URLs can be secured and don't need to be public. Currently there are 3 types of datasets that are managed:

- **Occurrence datasets:** Accessible species occurrence data such as specimens or observations. Currently there are roughly 11.800 registered occurrence datasets using various protocols such as DiGIR, TAPIR, BioCASE or Darwin Core Archives.
- **Checklists:** taxonomic, nomenclatural or other taxon oriented checklists published as Darwin Core Archives. This includes the Catalogue of Life.
- **External Metadata:** dataset descriptions originating from external metadata networks with the actual data being opaque to GBIF. Currently GBIF provides access to roughly 25.000 entries from the Knowledge Network for Biocomplexity (KNB) <sup>1</sup>

### Checklistbank

All registered checklists are indexed by GBIF and kept in a single store called Checklist Bank (CLB). Checklist Bank offers uniform webservices across all checklists and uses the GBIF Backbone (see below) as a way to move between names in the different checklists. To do so, every checklist record, if possible, is mapped to a taxon from the GBIF backbone.

### Taxonomic Backbone

The GBIF taxonomic backbone, often simply called the *nub*, is the reference taxonomy used by GBIF to link all occurrence and checklist records. It provides a standard classification most importantly used for browsing and calculating various metrics. It is a synthetic dataset which is updated from time to time (currently every couple of months) by a software which processes a prioritized selection of the checklists indexed in Checklist Bank. The backbone itself is also treated and stored as a regular checklist in CLB.

### Occurrence index

Nearly 400 million occurrence records have been indexed by GBIF at this point and they are all linked to the backbone taxonomy if some taxonomic identification was provided. This is done by using a public taxonomic matching service which handles fuzzy name matching, partial classification matching and is homonym aware. It is important to stress that no taxonomic names found in occurrences are incorporated into the GBIF backbone.

---

<sup>1</sup> <https://knb.ecoinformatics.org/>

## Portals

GBIF is about to release a new data portal later this year which provides a fully integrated access to the above listed components. As the current portal<sup>2</sup> is not maintained further, all work done for integrating the Catalogue of Life has focussed on the new portal. A first test preview will be available within the next weeks.

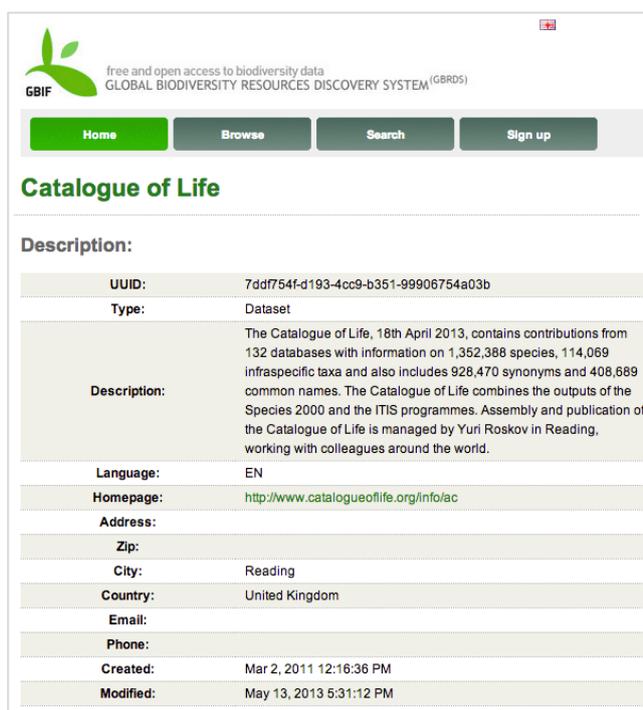
Until the new portal becomes available GBIF still updates the data presented in the current portal, but it often takes several months until data changes at the source actually become visible. An improved latency down to a week or even a day was a key motivation for building the new portal.

## Description of GBIF CoL processing workflow

### CoL in GBIF Registry

The entire catalogue of life is registered as a checklist with the dataset key 7ddf754f-d193-4cc9-b351-99906754a03b<sup>3</sup>, with the URL of the Darwin Core Archive provided by the CoL download service being registered as a service endpoint. In the GBIF registry all stored metadata details such as the title, description and contacts are automatically updated based on metadata found in the registered endpoints. In the case of CoL this would be an EML file bundled with the Darwin Core Archive. Currently this metadata file is missing from the downloaded archive, so manual curation is needed which is error prone and often lacking behind the actual version used. Future versions of the download service will provide the needed metadata for fully automated updates.

To help citation, every GSD included in the CoL is also registered as a constituent dataset of the catalogue on its own. Again the information stored for a GSD comes from an EML metadata file per GSD which is included in the download archive already, but which is currently still missing some key information **Error! Bookmark not defined.**



free and open access to biodiversity data  
GLOBAL BIODIVERSITY RESOURCES DISCOVERY SYSTEM (GBRDS)

Home Browse Search Sign up

### Catalogue of Life

**Description:**

<b>UUID:</b>	7ddf754f-d193-4cc9-b351-99906754a03b
<b>Type:</b>	Dataset
<b>Description:</b>	The Catalogue of Life, 18th April 2013, contains contributions from 132 databases with information on 1,352,388 species, 114,069 infraspecific taxa and also includes 928,470 synonyms and 408,689 common names. The Catalogue of Life combines the outputs of the Species 2000 and the ITIS programmes. Assembly and publication of the Catalogue of Life is managed by Yuri Roskov in Reading, working with colleagues around the world.
<b>Language:</b>	EN
<b>Homepage:</b>	<a href="http://www.catalogueoflife.org/info/ac">http://www.catalogueoflife.org/info/ac</a>
<b>Address:</b>	
<b>Zip:</b>	
<b>City:</b>	Reading
<b>Country:</b>	United Kingdom
<b>Email:</b>	
<b>Phone:</b>	
<b>Created:</b>	Mar 2, 2011 12:16:36 PM
<b>Modified:</b>	May 13, 2013 5:31:12 PM

Catalogue of Life has constituent **WTaxa: Electronic Catalogue of Weevil names (Curculionidae)** in the Catalogue of Life

Catalogue of Life has constituent **Xylariaceae: Home of the Xylariaceae** in the Catalogue of Life

Catalogue of Life has constituent **ZOBODAT: Zoological-Botanical Database (Vespoidea)** in the Catalogue of Life

Catalogue of Life has constituent **Zygomycetes** in the Catalogue of Life

The Catalogue of Life Partnership **owns** Catalogue of Life

HTTP Installation **serves** Catalogue of Life

**Endpoints**

 **DWC-ARCHIVE-CHECKLIST** - [http://www.catalogueoflife.org/DCA\\_Export/zip-fixed/archive-complete.zip](http://www.catalogueoflife.org/DCA_Export/zip-fixed/archive-complete.zip)

<sup>2</sup> <http://data.gbif.org/>

<sup>3</sup> <http://gbrds.gbif.org/browse/agent?uuid=7ddf754f-d193-4cc9-b351-99906754a03b>

## CoL in the GBIF Checklistbank

The entire catalogue as provided in the download archive is inserted on a regular basis into checklist bank. Until the new portal becomes available Checklist Bank offers a standalone user interface which can be browsed<sup>4</sup> and used through webservice<sup>5</sup>. A taxon page with the latest scrutinizer (given as *according to*) is shown in this screenshot:

The screenshot shows the GBIF Checklist Bank interface for the species *Sciurus vulgaris* Linnaeus, 1758. The page is titled 'Catalogue of Life - Sciurus' and includes a search bar and navigation tabs. A prominent 'BETA' stamp is visible. The main content area is divided into several sections:

- Classification:** Direct Parent: *Sciurus*; Kingdom: Animalia; Phylum: Chordata; Class: Mammalia; Order: Rodentia; Family: Sciuridae; Subfamily: ---; Genus: *Sciurus*; Subgenus: ---; Species: *Sciurus vulgaris*.
- Identifier:** urn:lsid:catalogueoflife.org:taxon:efbec818-29c1-102b-9a4a-00304854f820:col20130418 [LSID]; <http://www.catalogueoflife.org/annual-checklist/details/species/id/6905550> [URL]; 6905550 [SourceID].
- Related Sources:** unknown; Belgian Species List; Dyntaxa; Encyclopedia of Life; English Wikipedia Species Pages; Fauna Europaea; IUCN Red List of Threatened Species; Integrated Taxonomic Information System; Interim Register of Marine and Nonmarine Species; Mammal Species of the World, 3rd edition.
- Primary Biodiversity Data:** A map showing point observations for the species, primarily concentrated in Europe.
- Vernacular Names:** Eurasian red squirrel [en].
- Flicker Images:** A row of small images showing the squirrel in various poses.

## CoL and the GBIF Backbone

For GBIF by far the most important role of the Catalogue of Life is when building the GBIF backbone. The catalogue acts as the starting seed and also provides the entire higher classification above families as the only source. That means every species covered by the CoL will also be included one to one in the GBIF backbone. Because of this, GBIF does not make use of the i4Life crossmapping tools; every CoL taxon has its corresponding taxon in the GBIF taxonomy known in advance. It also means that it is crucial for GBIF that the CoL provides a comprehensive higher taxonomy, as GBIF does not augment the CoL taxonomy above family level. Currently CoL contributes 56% of all taxa in the GBIF backbone and 40% of all accepted family names.

<sup>4</sup> <http://ecat-dev.gbif.org/checklist/10>

<sup>5</sup> <http://ecat-dev.gbif.org/api/clb>