

The Catalogue of Life Standard Dataset

Version 7, 23rd September 2014

The Catalogue of Life plans to deliver a standard set of data for every known species. This document presents a simple description of the standard dataset which is both the core knowledge set of the Catalogue of Life and around which processes and protocols are designed. This standard dataset is used in many contexts and includes minimum content of data transmitted between components of the programme, and the minimum content of data transmitted in public products. These data are drawn from an array of participating taxonomic databases: Global Species Databases (GSD) – databases containing worldwide coverage of all the species within one taxon, Regional Species Databases (RSD) – databases containing regional coverage of the species within a taxonomic group, or Thematic Species Databases (TSD) – databases containing a subset of species based on special biodiversity or social themes. In this document we will use the name ‘Source Database’ for GSDs, RSDs and TSDs.

The Catalogue of Life has defined **14 field groups to be the standard set of data** for each species (or infraspecific taxa).

- 1. Accepted Scientific Name** linked to **Reference(s)** (obligatory)
- 2. Synonym(s)** linked to **Reference(s)** (obligatory, where available)
- 3. Common Name(s)** linked to **Reference(s)** (obligatory, where available)
- 4. Classification above genus, and up to the highest taxon in the database** (obligatory, where available)
- 5. Distribution** (obligatory, where available)
- 6. Life zone** (obligatory, where available)
- 7. Current and Past Existence** (obligatory, where available)
- 8. Additional Data** (optional)
- 9. Latest taxonomic scrutiny** (obligatory)
- 10. Reference(s)** (obligatory, where available)
- 11. Taxon Globally Unique Identifier** (obligatory, where available)
- 12. Name Globally Unique Identifier** (obligatory, where available)
- 13. Catalogue of Life LSID** (obligatory)
- 14. Source Database** (obligatory)

Some of the source databases additionally supply subspecies or varieties. The same dataset is used for each of these. Also, all information from field groups # 1, 2, 3, 5, 6, 7, 9 & 10 for infraspecific taxa should be given for both the species and infraspecific taxa (i.e. the ‘replicated’ system of TDWG Plant Names Standard).

Additional information is available either within the appropriate Source Database, or through hyperlinks to other databases.

1. Accepted Scientific Name (obligatory)

The Accepted, Valid or Correct scientific name (terminology for this name varies between the Codes of Nomenclature, in the Catalogue of Life we use the term ‘Accepted’) currently accepted for the species as a taxon. There should be exactly one per species. Two variants of NameStatus are possible in databases: ‘Accepted name’ or ‘Provisionally accepted name’.

‘**Accepted name**’ is the name currently accepted for the species by the compiler or editor of dataset as a quality taxonomic opinion.

‘**Provisionally accepted name**’ is the name currently accepted for the species by the dataset compiler, but with some element of taxonomic or nomenclatural doubt.

Content: a) Accepted Name of species

Genus | SubGenusName (where appropriate) | SpeciesEpithet |
AuthorString | Sp2000NameStatus | Reference(s)

b) Accepted Name of infraspecific taxon

only subspecies for taxa under the Code of zoological nomenclature; only subspecies, varieties and forms for taxa under the Code of botanical nomenclature:

Genus | SubGenusName | SpeciesEpithet | AuthorString |
InfraspeciesMarker (where appropriate) | InfraspeciesEpithet |
InfraspeciesAuthorString | Sp2000NameStatus | Reference(s)

In the case of Virus Names (i) the Genus is placed in the Genus field, and (ii) the polynomial species name is placed in the species epithet field. Virus species names have no official author.

<i>Where:</i>	Genus	= Latin genus name.
	SubGenusName	= Latin subgenus name.
	SpeciesEpithet	= second part of species name, Latin epithet
	AuthorString	= name of author(s), who described this species or published current combination (Style of authorstring depends on nomenclatural practices under different Codes)
	InfraspeciesMarker	= marker of infraspecific rank, where appropriate following Code regulations, for example, subsp., var., forma – for plants. (Presence and style of infraspecific markers depends on nomenclatural and taxonomic practices under different Codes)
	InfraspeciesEpithet	= third part of trinomial name, Latin epithet

InfraspeciesAuthorString	= name of author(s), who described this infraspecific taxon or published current combination (Style of authorstring depends on nomenclatural practices under different Codes; it could include year where appropriate)
Sp2000NameStatus	= the Catalogue of Life name status translated from source database: <ul style="list-style-type: none"> • Accepted <i>or</i> • Provisionally Accepted. CoL Data Submission Format also includes field GSDNameStatus, which retains the original status of the name in the source database.
Reference(s)	= just one reference that contains the original (validating) publication of taxon name or new name combination – Nomenclatural Reference , <i>or</i> one or more references that accept this species in the same taxonomic status, and with the same name – Taxonomic Acceptance Reference(s)

Example: Acacia | sieberiana | DC. | Accepted name | ReferenceID
(Accepted name record for Acacia sieberiana extracted from ILDIS database)

2. Synonym(s)

(obligatory, where available)

The list of Synonyms can include from 0 to many species or infraspecific names, which are given the Catalogue of Life synonymic status (Sp2000NameStatus). The three possibilities give the information sufficient for clear synonymic indexing, but do not give the full nomenclatural details, as these differ markedly in structure and context across different Codes. It is therefore necessary to ‘translate’ the very varied sorts of synonymic status in the source databases to create a uniform, accurate, but broad set of synonymic links for use in the Catalogue of Life.

(Category A) List of "**Synonyms**" - names which point unambiguously at one species (synonyms, in the CoL sense, include also orthographic variants and published misspellings)

(Category B) List of "**Ambiguous synonyms**" - names which are ambiguous because they point at the current species and one or more others e.g. homonyms, pro-parte synonyms (in other words, names which appear more than in one place in the Catalogue).

(Category C) List of "**Misapplied names**" - names that have been wrongly applied to the current species, and may also be correctly applied to another species.

Some synonyms of species can be trinomials, and have taxonomic rank of subspecies (in zoology), or subspecies, variety and form (in botany).

Content: Genus | SubGenusName (where appropriate) | SpeciesEpithet | AuthorString | Sp2000NameStatus | Reference(s) (obligatory)

or for trinomial synonyms (subspecies and varieties):

Genus | SubGenusName (where appropriate) | SpeciesEpithet | AuthorString | IntraspeciesMarker | IntraspeciesEpithet | IntraspeciesAuthorString | Sp2000NameStatus | Reference(s) (obligatory)

Where:

Genus	= as above
SubGenusName	= as above
SpeciesEpithet	= as above
AuthorString	= as above
Sp2000NameStatus	= the Catalogue of Life synonym status translated from source database as <ul style="list-style-type: none">• Synonym• Ambiguous Synonym• Misapplied Name CoL Data Submission Format also includes field GSDNameStatus, which retains the original status of the name in the source database.
Reference(s)	= as above

Examples:

Acacia | purpurascens | Vatke | Misapplied name | ReferenceID

Acacia | sieberiana | DC. | subsp. | vermoesenii | (De Wild.)Troupin | Synonym | Refs. #, #, #

Acacia | abyssinica | sensu auct. | Misapplied name | ReferenceID

(Synonym records for Acacia sieberiana extracted from ILDIS database)

3. Common Name(s) (obligatory, where available)

List of Common Names (also known as Vernacular Names) can include from 0 to many names.

Content: CommonName | TransliteratedName | Language | Country | Area (optional, where appropriate) | Reference(s)

*D. Eades, C. Flann, W. Addink, L. Abucay & Y. Roskov (after F. Bisby, Y. Roskov, T. Bourgoïn & D. Ouvrard, 2012)
The Catalogue of Life Standard Dataset, version 7.0, 23rd September 2014*

phylum, sector of Cercopoidea Organised On Line is attached as superfamily, sector of ILDIS World Database of Legumes is attached as one family). The Catalogue of Life requires each GSD to indicate the highest taxon that is given in the GSD, and to provide the classification beneath it down to species level.

The Catalogue of Life management classification includes taxa of seven basic ranks only: **Kingdom – Phylum – Class – Order – Superfamily – Family – Genus**.

Content: Kingdom | Phylum | Class | Order | Superfamily | Family | Genus
**Incertae sedis* (“not assigned” as in the *Catalogue of Life*) taxa are also allowed in ranks of phylum, class, order, superfamily and family, but not in ranks of kingdom and genus.

Plus, Catalogue of Life Taxon LSID with every taxon in the classification.

Where:

Kingdom	= Latin scientific name of the kingdom that includes the specified phyla.
Phylum	= Latin scientific name of the division or phylum that includes the specified classes.
	If the taxon is not known then this must be stated (e.g. phylum, class, order, superfamily and family labelled <i>incertae sedis</i> (not assigned) in taxonomic treatments) and the next higher taxon must be given with its rank.
Class	= Latin scientific name of the class that includes the specified orders.
Order	= Latin scientific name of the order that includes the specified families or superfamilies for insects.
Superfamily	= Latin scientific name of the superfamily that includes the specified families (for insect groups only).
Family	= Latin scientific name of the family that includes the specified genera.
Genus	= Latin scientific name of the genus that includes this species.
CoL LSID	= CoL Taxon Matcher software issues CoL Global Unique Identifiers for every taxon recognised in the Catalogue of Life at the stage of optimisation of CoL database using the Life Science Identifier (LSID) system (http://sourceforge.net/projects/lsids).

Example: Orthoptera | Stenopelmatoidea | Gryllacrididae |
Abelona

(Example of classification above species Abelona frontalis (Order – Superfamily – Family – Genus) extracted from OrthopteraSF)

5. Distribution (obligatory, where available)

The Catalogue of Life can include structured distribution records in standard schemas as well as text information. Structured distribution records following standard schemas is the preferred choice.

Content: DistributionElement | StandardInUse | DistributionStatus

DistributionElement = TDWG code or Name of an area using one of the agreed distribution standards:
- for land areas: Updated TDWG Level 4 areas (preferred), or ISO 3-letter country codes
- for sea areas: Intersect of IHO's and EEZ areas (see: VLIZ (2010), Intersect of IHO Sea Areas and Exclusive Economic Zones (v5, 2009). The standard is available online at <http://www.vliz.be/vmdcdata/vlimar/downloads.php>)
- text up to 255 characters

StandardInUse = short name for each distribution standard; examples of expected values

- TDWG Level 4
- TDWG Level 3
- FAO_ISO *or*
- TEXT when providing free text distribution

DistributionStatus = multi-state descriptor code. Score for multiple states where more than one applies. Expected values

- Native
- Domesticated
- Alien
- Uncertain

Example: ANG-00 | TDWG4 | Native
BEN-00 | TDWG4 | Native
BKN-00 | TDWG4 | Native

(Distribution record of Acacia sieberiana extracted from ILDIS database)

Note: The proposed Catalogue of Life Standard does not include a source reference for this data. However, it is recommended to GSDs, as part of the ‘best practice’, that the GSD contains reference(s) linked to each area.

6. Life zone (obligatory, where available)

Including information about which broad life zone an organism live in allows a distinction to be made between marine, terrestrial or freshwater taxa.

Life zone is a single multi-state descriptor field, for which multiple values can be recorded. This field does not have references associated with it and is based on the Source Database author’s expert knowledge. The Life zones are: **Marine, Brackish, Freshwater, Terrestrial, Unknown** (Note: these categories corresponds with GISIN realms .).

- i) Species or infraspecific taxa occurring in more than one zone should have values in the field separated by commas.
- ii) Life zone for parasitic (or similar) species is the same as life zone for their host organisms.

Content: LifeZone

Where: LifeZone = a single multi-state descriptor according to this fixed and non-extensible coding (interoperable with the GISIN standard titled as ‘Realm’ by GISIN.):

- Marine
- Brackish
- Freshwater
- Terrestrial
- Unknown

Example: Freshwater, Brackish

(Life zone record of Neptunia plena extracted from ILDIS database)

7. Current and Past Existence (obligatory, where available)

This field group is aimed at distinguishing living species from extinct (no longer in existence, usually represented by fossils).

Content: IsExtinct | HasPreHolocene | HasModern

<i>Where:</i>	IsExtinct	= this field distinguishes taxa known as presently living on the planet, Value: 0 (FALSE), from species which are believed to be no longer living, Value: 1 (TRUE), or species with unknown status (Value: NULL). Default 0 (FALSE)
	HasPreHolocene	= means the taxon is known to have existed before the Holocene Era (more than 11,700 years ago). Value: 1 (TRUE) if the taxon is known to have occurred before the Holocene era. Value: 0 (FALSE) if this is not the case. Value: NULL if the status is unknown. Default 0 (FALSE)
	HasModern	= means the taxon is known to have occurred during the Modern era (less than 11,700 years ago). Value: 1 (TRUE) if the taxon is known to have occurred during the Modern era. Value: 0 (FALSE) if this is not the case. Value: NULL if the status is unknown. Default 1 (TRUE)

Examples: All living species without known fossils:
 IsExtinct = 0 (FALSE) | HasPreHolocene = 0 (FALSE) |
 HasModern = 1 (TRUE)

Ginkgo biloba (Jurassic tree, also with living specimens):
 IsExtinct = 0 (FALSE) | HasPreHolocene = 1 (TRUE) |
 HasModern = 1 (TRUE)

Raphus cucullatus (Dodo, extinct endemic to Mauritius island. The Dodo's appearance in life is evidenced only by drawings, paintings and written accounts from the 17th century. The last living specimen was recorded in 1662).
 IsExtinct = 1 (TRUE) | HasPreHolocene = 0 (FALSE) |
 HasModern = 1 (TRUE)

Stegosaurus stenops (Jurassic dinosaur):
 IsExtinct = 1 (TRUE) | HasPreHolocene = 1 (TRUE) |
 HasModern = 0 (FALSE)

Melanoplus spretus (Rocky Mountain Locust, commonly believed to be extinct, seen alive in Modern era, also exists as fossils):
 IsExtinct = NULL | HasPreHolocene = 1 (TRUE) | HasModern
 = 1 (TRUE)

8. Additional Data

(optional)

This field can contain free text up to 255 characters. It can contain information from one or several data fields from the source database (for example, type specimen, taxonomic comments, common name of the family, habit/life form, detailed ecology, host, etc.) as decided by the custodian of the source database. Unlike all other field groups, there is no intention to make these data compatible across taxa. It can therefore be distinctive or particular to the species supplied by one database.

Content: AdditionalData (*Free text, which might be structured with headings*)

Where: AdditionalData = free text up to 255 characters

Example: Type strain: strain ATCC 33244 = CFBP 3612 = CIP 105207 = ICPB EA175 = LMG 2665 = NCPPB 1846 = PDDCC 1850, "Pantoea ananatis corrig. (Serrano 1928) Mergaert et al. 1993, comb. nov."

(Additional Data record for Erwinia ananatis extracted from BIOS database)

9. Latest taxonomic scrutiny

(obligatory)

This field group should contain only one record of the Latest Taxonomic Scrutiny (LTS) of the species or infraspecific taxon in the source database. LTS includes (a) name(s) of the taxonomic expert or editor, who is responsible for the taxonomic concept accepted in the source database and (b) date when the expert or editor assessed the record. If the source database has multiple records, just the most recent should be passed to the Catalogue of Life. If the source database has no latest taxonomic scrutiny records, but is the work of one specialist or a small team, then for the whole database this field should show the scrutiny by that specialist or small team. Users of the Catalogue of Life content are obliged to cite the name of LTS specialist with each species taken from the Catalogue.

Content: LTSSpecialist | LTSDate

Where: LTSSpecialist = surname of taxonomic editor and initial(s) follows after the surname without additional separators (see example below).

LTSDate = date of record scrutiny (revision) in the source database; style specified by the custodian of Source Database; 'Year' is obligatory; 'Month' & 'Day' might be applied where available.

Example: Farjon A. | 2014

(Scrutiny record for Pinus strobus extracted from Conifer Database)

10. Reference(s) (obligatory, where available)

References should be linked with Accepted Scientific Names, Synonyms and Common names.

Content: Author(s) | Year | Title | Source | Reference Type

Where:

Author(s)	= author (or many) of publication
Year	= year of publication
Title	= title of paper or book
Source	= title of periodicals, volume number, and other common bibliographic details
ReferenceType	= taxonomic status of reference: <ul style="list-style-type: none">• Nomenclatural Reference NomRef (just one reference which contains the original (validating) publication of taxon name or new name combination <i>or</i>• Taxonomic Acceptance Reference(s) TaxAccRef (one or more bibliographic references, where the name is mentioned in the same taxonomic status (i.e. as a species or as a synonym) <i>or</i>• Common Name Reference(s) ComNameRef (one or more bibliographic references that contain common names)

Example: Ross, J.H. | 1979 | A conspectus of African Acacia | Mem. Bot. Surv. S. Afr. 44: 1-150 | TaxAccRef

(Reference record extracted from ILDIS database)

11. Taxon Globally Unique Identifier (obligatory, where available)

This field is provisional in the Catalogue of Life; it is intended for future use to support global “ecosystem” of biodiversity data and provide an interoperability between taxonomic databases, CoL and its corporate users.

Content: GSDTaxonGUID

Where: GSDTaxonGUID = single text field containing the single Globally Unique Identifier supplied for every taxon at species or infraspecific rank by the Source Database supplying this taxon as a concept. Taxon GUID suppose to reflect changes in the accepted concept.

Examples:

Notes

- i) This is for future use, once a pattern of Source Databases supplying Taxon Globally Unique Identifiers are established.
- ii) At present this leaves open the possibility that the Identifiers may be LSIDs or some other form of unique identifiers (for example, if Source Database does not provides Taxon GUIDs, the Catalogue of Life may populate unique identifiers which are in Source Database use).
- iii) It is assumed that the unique identifier received already contains some provenance indication, as is the case with LSIDs.
- iv) This field is intended for unique identification and provenance metadata tracking and possible use in the Catalogue of Life web-services. It is not necessarily for display in the public interface.

12. Name Globally Unique Identifier (obligatory, where available)

Content: GSDNameGUID

GSDNameGUID = one or more text strings containing one or more Unique Identifiers supplied for a single name at species or infraspecific rank by the Source Database supplying this name.

Examples: urn:lsid:ipni.org:names:30000959-2
(A plant scientific name from IPNI)

urn:lsid:indexfungorum.org:Names:213649
(A scientific name of a fungi from Index Fungorum)

Notes

- i) This is for future use, once a pattern of Source Databases supplying Name Globally Unique Identifiers are established.
- ii) This field is intended for unique identification and provenance metadata tracking and possible use in the Catalogue of Life web-services. It is not necessarily for display in the public interface.
- iii) There is no single authority for issuing unique identifiers for names of all organisms, although there is a perception that this task will be undertaken by existing the Nomenclator organisations. It is possible that we may receive more than one unique identifier per name, and may wish to store all of those supplied.

13. Catalogue of Life Taxon LSID (obligatory)

This field is reserved for CoL Taxon Life Science Identifiers. CoL Taxon Matcher software issues permanent CoL Global Unique Identifiers for every taxon concept recognised in the Catalogue of Life at the stage of optimisation of CoL database using the Life Science Identifier (LSID) system (<http://sourceforge.net/projects/luids>). Life Science Identifiers are persistent, locating independent and unique identifiers for specific data in the web. They should serve to track changes in the CoL records.

Attention: a replacement for current implementation of this field is under development.

Content: CoL Taxon LSID

Where: CoLTaxonLSID = a single text field containing the single Globally Unique Identifier supplied for every taxon at species or infraspecific rank by the Source Database supplying this taxon as a concept. Taxon GUID is supposed to reflect changes in the accepted concept.

Examples: Zorka angelinae | urn:lsid:catalogueoflife.org:taxon:653bc0d9-f89c-11e0-af7a-6c3dedecd876:col20121017

14. Source Database (obligatory)

This information (metadata) will be supplied for each source database only once, but it will be shown as a part of every record in the Catalogue of Life.

Content: DatabaseFullName | DatabaseShortName | DatabaseVersion | ReleaseDate | AuthorsEditors | TaxonomicCoverage | GroupNameInEnglish | Abstract | Organisation | HomeURL | Coverage | Completeness | Confidence | LogoFileName | ContactPerson

Where:

DatabaseFullName	= full title of Source Database; as supplied by the custodian
DatabaseShortName	= abbreviated or shortened memorable name of Source Database intended for easy use in day-to-day communications; as supplied by the custodian
DatabaseVersion	= database version (number or code, plus date, where Month and Year are obligatory) provided by the custodian; style specified by the custodian of Source Database

ReleaseDate	= original date (Year-Month-Date) of issue of the version for the Catalogue of Life.
AuthorsEditors	= name(s) of Source Database editor or author; as specified by the custodian
TaxonomicCoverage	= higher taxon(s), which represent taxonomic sector(s) covered by the database
GroupNameInEnglish	= English name of the taxon covered by the Source Database
Abstract	= standardised short database description (text of 50-70 words) for use in the Catalogue of Life supporting materials, such as the booklet published with the Annual Checklist on DVD
Organisation	= name of the Organisation that hosts the Source Database
HomeURL	= Uniform Resource Locator (Internet address) of Source Database home page
Coverage	= geographic coverage of the Source Database for the taxon <ul style="list-style-type: none"> • Global for a worldwide coverage • Regional for a geographically restricted coverage. If value is Regional, the region should be specified in brackets
Completeness	= percentage of completeness of species list of the taxon provided by the Source Database
Confidence	= quality of taxonomic checklist with values 1 to 5; quality is stated by the custodian. Confidence indicators are as follows: <p>1 - Caution! This data set does not contain well scrutinised taxonomic checklist, and in parts may be a list of taxonomically unvetted names only. However, it is used temporarily by the Catalogue of Life to fill major gaps as only available source at the time. See database abstract for more details.</p> <p>2 - Caution! This data set is a scrutinised taxonomic checklist, but it is incomplete and at an early stage of its development. See database abstract for more details.</p> <p>3 - This is a well-scrutinised taxonomic checklist, but it is restricted to a subset of species by geography (regional database), or sector of</p>

biological discipline (e.g. thematic database in particular ecological area, conservation, quarantine, pest and disease control, medicine or molecular biology, etc). This data set was included in the Catalogue of Life to fill gaps at lower levels of the taxonomic classification (e.g. species, genera) as temporarily solution. See database abstract for more details.

4 - This is a nearly complete and fully scrutinised taxonomic checklist with a good quality of expertise at the current stage of its development.

5 - This is a complete and fully scrutinised taxonomic checklist for an entire taxon with a high quality of expertise and frequent updates, which covers nearly all known species diversity in the taxon worldwide.

LogoFileName = name of the file containing the Source Database logotype. Digital image of logo should be supplied with first download of data. Technical requirements include minimum image size as 10x10 mm with 300 dpi

ContactPerson = name of contact person for the Source Database and email address; as specified by the custodian; for internal use by the CoL Secretariat only; this information will not be released to the public

Example: Cercopoidea Organised On Line | COOL | 4, Sep 2011 | 2011-09-07 | Soulier-Perkins A. | Cercopoidea | Froghoppers | COOL is a systematic... | Muséum National d'Histoire Naturelle, Paris, France | <http://rameau.snv.jussieu.fr/cool/> | Global | 95 | 4 | <http://www.catalogueoflife.org/images/databases/COOL.png> | Soulier-Perkins A.

(Source Database record extracted from COOL database)