

Annex 4

The Import module of DaRWIN

1. INTRODUCTION	2
2. METHODOLOGY	2
2.1 Reworked importation procedure	2
Figure 1: The new import procedure in 3 steps using XLS templates	3
2.2 Lower complexity when importing external specimen data	4
2.3 Additional importation templates and verification interfaces	5
3. INFRASTRUCTURE	7
4. RESULTS AND RECOMMENDATIONS	7
4.1 The DaRWIN side changes	7
4.1.1 Taxonomy	7
4.1.2 Localities	10
4.1.3 Lithostratigraphy	11
4.1.4 Embedded multimedia files	13
4.1.5 Links to remote multimedia files	14
4.2 Input Templates	15
4.2.1 Need for an import tool	15
4.2.2 RMCA Excel template	16
a. Sheets	16
b. Buttons	18
c. Forms	19
d. Export	22
4.2.3 LibreOffice template	24
a. Sheets	24
b. Buttons	25
c. Forms	25
d. Export	29
e. Taxonomy check	29
4.3 Integration of previous databases and Import of data in DaRWIN	29
4.3.1 Mapping of the RBINS MISTA database	30
4.3.2 The RBINS Geology Collection	30
4.3.3 The RBINS Paleontology Collection	30
4.3.4 RMCA zoology	30
4.3.5 Mapping of RMCA wood biology data	31

1. INTRODUCTION

The existing import module of DaRWIN was developed by RBINS in the framework of a previous project. The procedure was based on a huge XLS file, exporting an XML file with all data (Sampling location, Taxonomy and specimens data). This XML file was then imported by DaRWIN with several levels of data checking.

This was extremely complex thanks to the size of the XLS file and the number of fields and frustrating for users because the import was always blocked somewhere.

The import procedure is nevertheless extremely useful to add new specimens to the existing database(s). Curators and research scientists already have many specimens encoded in smaller databases or use XLS files. They know how to use spreadsheets which are common softwares for users.

The manual encoding of data into the DaRWIN database is estimated to be between 5000 and 10000 specimens / year / encoder. This process is thus very slow and not efficient.

The import procedure from existing databases and/or spreadsheets allows to import up to 10 times more specimens in the same time period for a trained FTE “import” encoder. This is why we decided to completely review the import processes.

2. METHODOLOGY

2.1 Reworked importation procedure

The workflow, database logic and interfaces to import external data from files into DaRWIN have been extensively reworked and expanded.

The import is now divided in 3 steps:

- Taxonomy
- Location(s)
- Specimens data

This segmented procedure allows to simplify the validation of the data. For Taxonomy and Locations files, it is possible to use external data validation using web services of authority databases.

A reference manual on the current importation procedure, written by Marielle Adam, is available on the GitHub repository of the project:

https://github.com/naturalsciences/natural_heritage_darwin/blob/STABLE_2020/doc/import%20user%20manual.docx

Number of the action	Description of the action	KPI	Deliverable	2015	2016	2017	2018
SO1_OO1_A2	To select the optimal tools for collection management (with easy import functionalities)	An analysis is available in the first half of 2016		X	X		

	A	B	C	D	E	F
	phylum	class	order	family	genus	species
1	CHORDATA	ACTINOPTERYGII	SILURIFORMES	Claridae	Tanganikallabes	Tanganikallabes alluapera Wright & Bailey, 2012
2	CHORDATA	ACTINOPTERYGII	SILURIFORMES	Claridae	Tanganikallabes	
3	CHORDATA	ACTINOPTERYGII	SILURIFORMES	Claridae	Tanganikallabes	Tanganikallabes mortiauxi Poll, 1943
4	CHORDATA	ACTINOPTERYGII	SILURIFORMES	Mochokidae	Synodontis	
5	CHORDATA	ACTINOPTERYGII	SYMBRANCHIFORMES	Mastacembelidae	Mastacembelus	Mastacembelus tanganicae (Günther, 1893)
6	CHORDATA	ACTINOPTERYGII	SILURIFORMES	Claridae	Lophiobagrus	Lophiobagrus cyclurus (Worthington & Ricardo, 1936)
7	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Neolamprologus	Neolamprologus toae (Poll, 1949)
8	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Aulonocranus	Aulonocranus dewindti (Boulenger, 1899)
9	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Chalinochromis	Chalinochromis brichardi Poll, 1974
10	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Cyprichromis	
11	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Eretmodus	
12	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Neolamprologus	Neolamprologus fasciatus (Boulenger, 1898)
13	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Neolamprologus	Neolamprologus tetrocephalus (Boulenger, 1899)
14	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Neolamprologus	Neolamprologus toae (Poll, 1949)
15	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Neolamprologus	Neolamprologus niger (Poll, 1956)
16	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Neolamprologus	Neolamprologus niger (Poll, 1956)
17	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Telmatochromis	Telmatochromis bifrenatus Myers, 1936
18	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Telmatochromis	Telmatochromis dhonti (Boulenger, 1919)
19	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Pseudosimochromis	Pseudosimochromis curvifrons (Poll, 1942)
20	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Petrochromis	Petrochromis fasciatus Boulenger, 1914
21	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Petrochromis	
22	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Uphthamotilapia	Uphthamotilapia nasuta (Poll & Matthes, 1962)
23	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae	Aulonocranus	Aulonocranus dewindti (Boulenger, 1899)
24	CHORDATA	ACTINOPTERYGII	PERCIFORMES	Cichlidae		

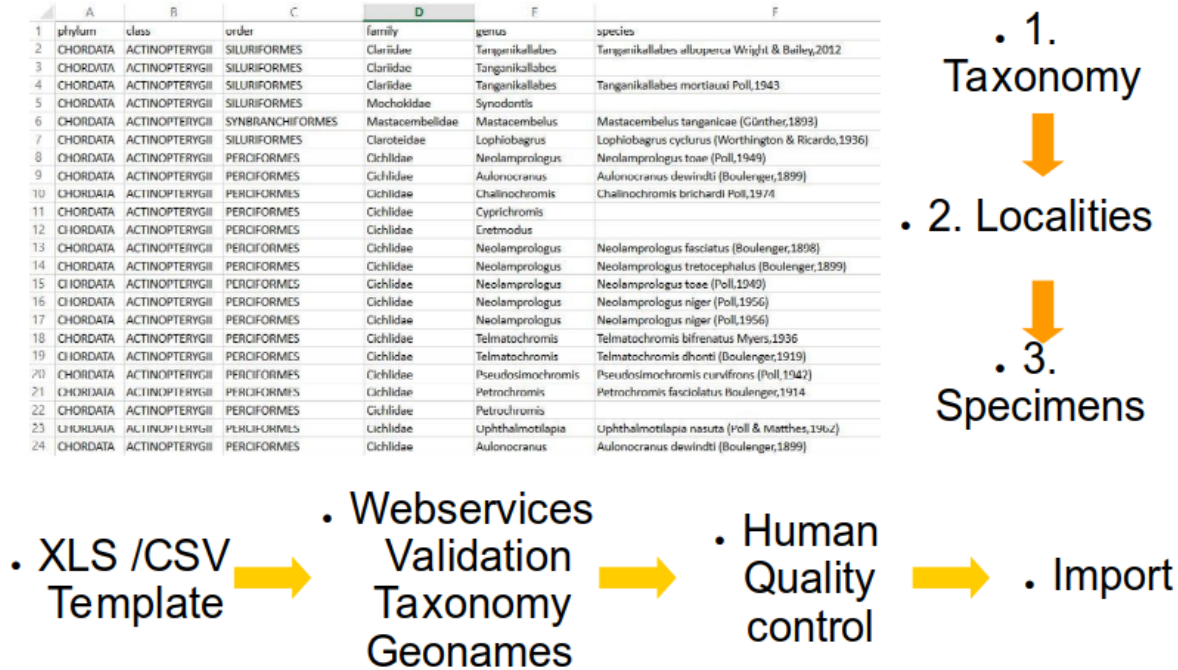


Figure 1: The new import procedure in 3 steps using XLS templates

As described in the reference manual, the complete workflow to import specimen data in DaRWIn is now splitted into three steps:

1. Importing the taxonomy
2. Importing localities (collecting stations)
3. Importing specimen data

Specimen data has to be imported last:

- The binding between the specimen and the taxa is done by the *full scientific name* (word containing the scientific name and the authorship information, without considering the rank)
- The binding between the specimen and the locality is done via the station number, which is then supposed to be unique. It is possible to disambiguate duplicate station numbers for one specific record or the whole dataset in the validation interface.

2.2 Lower complexity when importing external specimen data

Another major change in DaRWIN has been the simplification of the procedures and workflow to import external data into DaRWIN. The initial version only accepted input data based on the ABCD XML schema for collection data. Flat data has to be first written on specific Excel files that feature a Visual Basic macro generating the XML file. This made the adaptation of the template procedure very complex, any modification had to be implemented in three different systems (the XML parser in darwin, the stored procedure in XML and the macro in Visual Basic). The macro and XML parser were very complicated to test and debug. The complexity in time (the parsing of XML documents is time consuming, such as the conversion of the Excel data into XML) and space (XML files are much larger than flat CSV files) was also needlessly high. The usage of XML introduced bugs related to the syntax of the XML document that were harder to diagnose: the intermediate data structure handed by the PHP server had to be serialized on the hard drive and syntactically analyzed for each problem. Besides, the XML ABCD format was not exposed to a public web service on the Internet as the corresponding web service BioCASE works through a database connection: the files could not be reused for any other usage. XML could have been interesting to handle controlled vocabularies and external thesaurus, but this functionality is not present in DaRWIN. A flat tabular structure was converted into an hierarchical structure, before being converted again the other way round into the flat structure of DaRWIN data model. The Excel files were also complicated to handle, having more than 100 columns, the name of each of them being case sensitive, and had a constraining order, while the imported data for specimen most of time contain just about 15 to 30 meaningful fields (label number, locality information scientific name).

Finally, the ABCD parser could work only on specimen data, and custom PHP parsers would have to be developed for any other type of data (people, localities etc...) . DaRWIN's team was also contacted in April 2019 by a team of scientists from the University of Rwanda for a JRS Biodiversity project . While they were interested in the possibility of integrating and cleaning external data into a reference collection database, this part of the system had to be simplified, to ease its documentation and usage for external users.

We decided to remove the XML part and parser and to replace it by a parser for CSV files that would follow the following concepts:

1. The data would be tab-delimited
2. The column name would follow a controlled vocabulary
3. Each column would be optional. The system could import specimens having already a collection number, or assign a new value (as numeric sequence) of this collection number is missing.
4. The column names would be case insensitive
5. Their order could be free

6. Geographical coordinates could be inserted in DMS (degree minutes seconds) while being converted in decimal degrees if they followed a consistent text pattern (eg "10°30' 45" W/E" or "N/S 14° 45')

These changes have been integrated gradually, first by keeping the XML existing parser and generating in-memory XML files inside of DaRWIN. This part has been finally removed and replaced by a PHP parser directly filling the staging table from the tab-delimited values. The import process was using four steps (creating of the Tab-delimited file, importing the data into the SQL *staging* tables, checking duplicate, and integrating the data in the normalized part of DaRWIN) instead of 6 (filling the template, generating the XML, parsing it, filling SQL *staging* tables, checking duplicate, and integrating the data in the normalized part of DaRWIN).

Finally, the jobs to import the data, that are asynchronous and background console operations, have been linked to the web interface of DaRWIN, while they were previously only available via Bash or DOS instructions (often provided by a SSH connection). This limited the importation tasks to IT-trained staff, while the new procedure is available for any user having the access rights to DaRWIN. However this introduced a moderate security risk, as the web interface has to execute shell commands. This risk can be mitigated by checking and controlling the type of parameters passed to the command as argument (limiting them as numeric values) .

This part of the work was surely the most complex and harder to test amongst the Darwin tasks within the framework of NaturalHeritage, as the completeness of imported data had to be cautiously checked, on more than 100 columns that could be combined in different ways. These checks could not be automated. The scalability (ie, ability to work on a great number of records) of the procedure had also to be verified, which implied huge amounts of data to be produced (which is sometimes almost as complex as developing the application) and lengthy test operations. We should actually have defined more rigorous and standardized test procedures. However, this kind of data is placed outside of the scope of unit testing which is easier to automate (which are more targeted and specific, but do not correspond to this scenario as they compare the behavior and backward-compatibility of new versions of a programme to a reference and stable behaviour). The import speed can be estimated by 2500 to 5000 records/hour for importing data from the source file and approximately the same duration for the check (detection of duplicates). Integration of the data in DaRWIN is faster (5000 to 10000 records/hour)

2.3 Additional importation templates and verification interfaces

Once the importation workflow has been simplified and the XML part removed in design pattern, a design pattern that could be reused with other types of data was made available. This design pattern followed a three steps approach, each of them corresponding to asynchronous background operations:

1. Importation of tab-delimited data into the staging part
2. Iterative checks to detect and remove duplicated either by
 - Creating new values and doing batch updates or...
 - Choosing an existing value and doing batch update on the others records
3. Integrating the cleaned data into the normalized part of DaRWiN . This is also an iterative task that can be done on parts of the dataset, and that can be sprawled, interrupted and resumed on several sessions

The figure 2 describes the importation workflow for specimen data, that checks duplicates, disambiguate homonymes, and missing data, and correct then by single or batch updates on:

1. Peoples
 - Collectors
 - Identifiers (for taxonomic, or geological attributions)
 - Donators
2. Institutions
3. Taxonomic identification (if the taxa is missing it can be created)
4. Expedition
5. Sampling locations (using their station code as link, that becomes a mandatory field)

It is also possible to enforce or disable a unicity constraint on the main specimen code, while uploading the file.

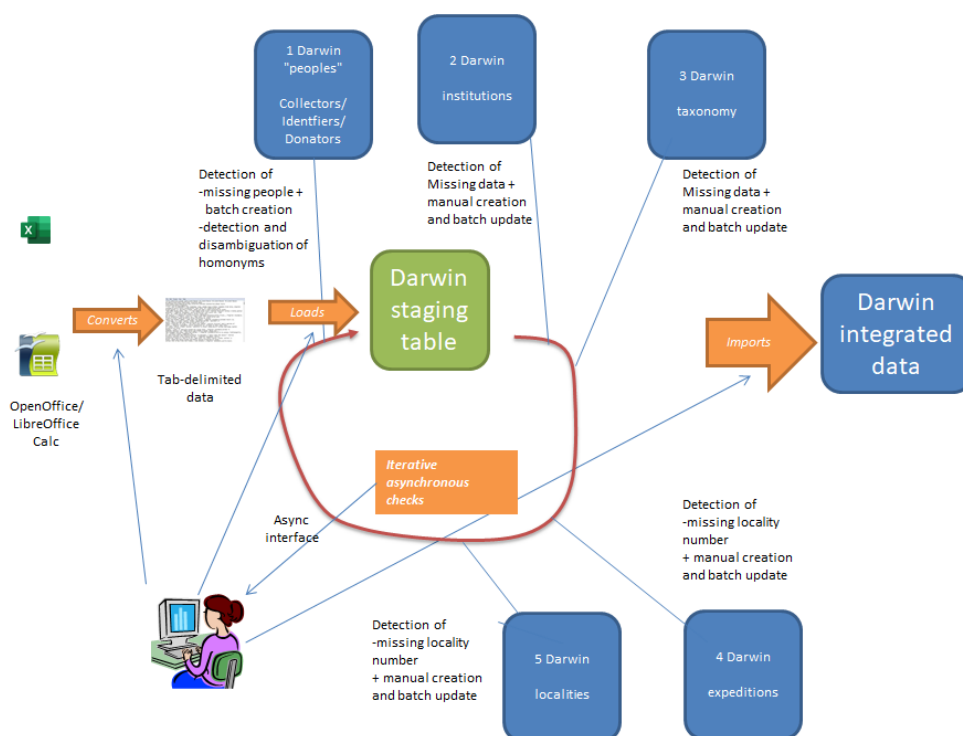


Figure 2. Importation workflow for specimen data

This template has been adapted to several other content types that can now be uploaded from tab-delimited files into DaRWIN and verified. The most important development effort didn't reside in the programming of the import logic, but in the development of visualization and data-management interface allowing users to access the staging tables and clean data from the Internet. A semaphore mechanism had to be implemented, to notify the users whether import tasks successfully ended or not, and to report errors (PHP and PostgreSQL exception are thrown to the web interface giving the status of an import to ease debug and correction of data, as they often give information about syntactical issues).

3. INFRASTRUCTURE

The adopted procedure simplifies the infrastructure needs as no specific server is requested by the new procedure. The complete process can be realized with Open Source technologies as the templates were developed for the proprietary Microsoft Office Excel but also for the Open Source LibreOffice suite. The main difference was in the programming which is in Visual Basic for the Excel macros and in Basic for the LibreOffice version.

4. RESULTS AND RECOMMENDATIONS

4.1 The DaRWIN side changes

4.1.1 Taxonomy

DaRWIN initial versions already featured a template mechanism to upload taxonomic hierarchies, but it was using a custom XML schema derived from ABCD (also requiring a Visual Basic macro) and had no validation interface allowing the user to check and validate data from the web interface.

Besides, the concept of "parallel" taxonomies (or taxonomical metadata) had been introduced in DaRWIN, allowing to publish different hierarchies for the same taxon and annotate their scientific accuracy. The initial importation mechanism used also one SQL transaction (either all data could be imported or none, error or taxonomic conflict cancelling the whole job).

The taxonomic template, and a substantial part of the database logic in the *staging* part of DaRWIN, needed therefore to be reworked. A template for a tab-delimited file has been defined, where users can provide a list of scientific names with upper ranks and authors within each row. The higher rank provided for each row is free (it can be the phylum, the order, the family or others...) but has to be already created in DaRWIN,

which should build the complete descending taxonomic tree from the higher taxonomic level to the lower one.

For each attempt to create taxon, 4 types of operations can be detected:

1. The taxon is missing and could be successfully created in DaRWIN (his parent exists both in the file and the system)
2. The taxon is missing but couldn't be created (the parent in the file cannot be created in DaRWIN)
3. The taxon has another parent in DaRWIN, for the considered parallel taxonomy, at least one of its parent has another hierarchy (Upper level conflict with DaRWIN)
4. The taxon is present twice or several times in the file with different hierarchies, that contradict themselves (Upper level conflict within the file)

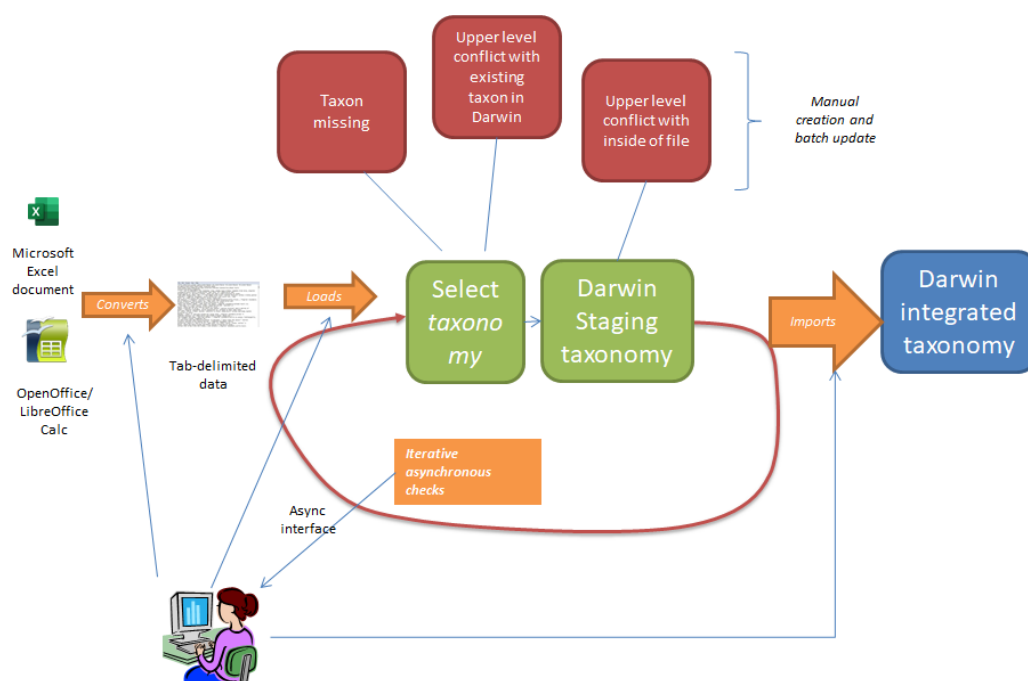


Figure 3. Importation workflow for taxonomic data

A validation interface has been developed (see figures 4 and 5). It features a pager with global statistics on the imported records, allowing the user to navigate through results. Each page displays 1000 rows.

- A simple color code (green rows for imported results and orange for errors, makes it more readable and allows users to rapidly identify issues).
- Each row contains a field where taxonomic hierarchies in DaRWIN and in the imported files are displayed as paths separated by “/”, allowing a rapid comparison.
- A button opens a modal window allowing you to manually create the taxon, and/or to correct the hierarchy of an existing taxon in DaRWIN.

- Buttons placed at the bottom of the web page allows to launch the check and integration again, and to change the parallel taxonomy that has been chosen in the import procedure. This allows users to handle problem without reimporting the tab-delimited file, and in several internet sessions

Count all : 2872
Current page : 1 / 3
Page : 1 [go](#)

Message	Count
imported_taxon	939
taxonomic_conflict	42
taxonomic_hierarchy_already_exists	12
taxon_to_be_created_without_suitable_parent	7
Total :	1000

All data:

Message	Count
imported_taxon	2665
taxonomic_conflict	152
taxonomic_hierarchy_already_exists	31
taxon_to_be_created_without_suitable_parent	24
Total :	2872

[Download all](#) [Download unimported](#) [Recheck and reimport](#)

Figure 4. Validation interface for taxonomic import (pager and statistics)

nautilus.rhins.be/darwin/backend/

Darwin Help

natural_heritage_darwin/imp...

nautilus.rhins.be/darwin/backend_dev.php/import/viewUnimportedData?id=552

Debug toolbar 1.5.12-dev

Config config

View Layer view

Log logs

Memory 4095.5 KB

Time 244 ms

SQL queries 15

Ctrl

Count all: 63

Current page: 1 / 1

Page: 1 / 1

go

Message	Count
imported_taxon	42
taxonomic_hierarchy_already_exists	21
Total:	63

All data

Message	Count
imported_taxon	42
taxonomic_hierarchy_already_exists	21
Total:	63

[Download all](#) [Download unimported](#) [Backcheck and reimport](#)

id	name	level	ref	name_cluster	imported	import_exception	compare hierarchies	import
295388	Anthomyidae	family	1	FALSE	taxonomic_hierarchy_already_exists		<div>Staging hierarchy (Anthomyidae family)</div> <div>Darwin hierarchy (Eucaryota (domain)/Animalia (kingdom)/Arthropoda von Siebold, 1848 (phylum)/Hexapoda Blainville, 1816 (sub phylum)/Insecta Linnaeus, 1758 (class)/Pterygota Gegenbaur, 1878 (sub class)/Ooptera Linnaeus, 1758 (order)/Brachycera (sub order)/Muscomorpha Sharp, 1894 (infra order)/Anthomyidae (family)</div>	Create taxon
295389	Macrochis	genus	1	TRUE	imported_taxon		<div>Staging hierarchy (Anthomyidae family)/Macrochis (genus)</div> <div>Darwin hierarchy (family)</div>	
295391	Rhynchosomops	genus	2	TRUE	imported_taxon		<div>Staging hierarchy (Anthomyidae family)/Rhynchosomops (genus)</div> <div>Darwin hierarchy (family)</div>	

Type here to search

15:40 23/12/2019

Figure 5. Validation interface for taxonomic import (hierarchy viewer)

4.1.2 Localities

It has been decided to create a specific importation template for localities, in order to keep in line with the normalization of localities described above.

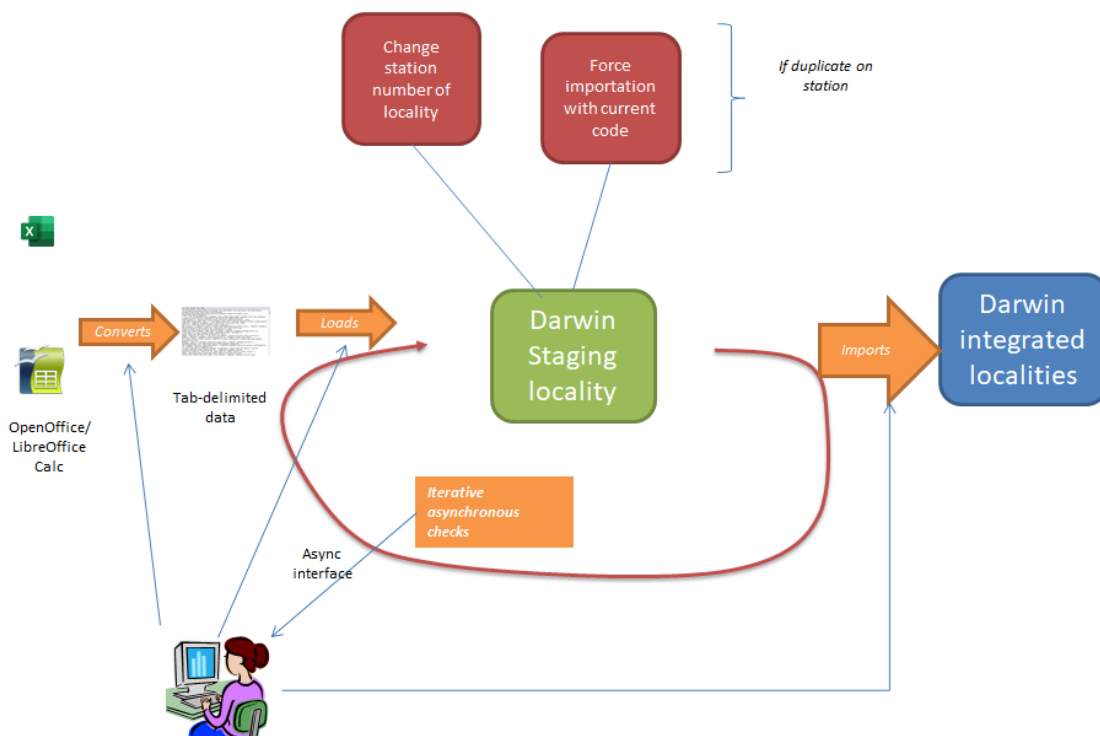


Figure 6. Importation workflow for locality data

4.1.3 Lithostratigraphy

A template has also been developed to import additional lithostratigraphic classifications, for fossil or mineralogical collections.

It has been decided to make the lithostratigraphic scale dynamically updatable, but to keep the existing chronostratigraphic frozen, as the chronostratigraphic scale is global and stable, while the lithostratigraphic scale is dependant from the location and less standardized.

The template columns are:

supergroup
group
formation
member
layer
sub_level_1
sub_level_2

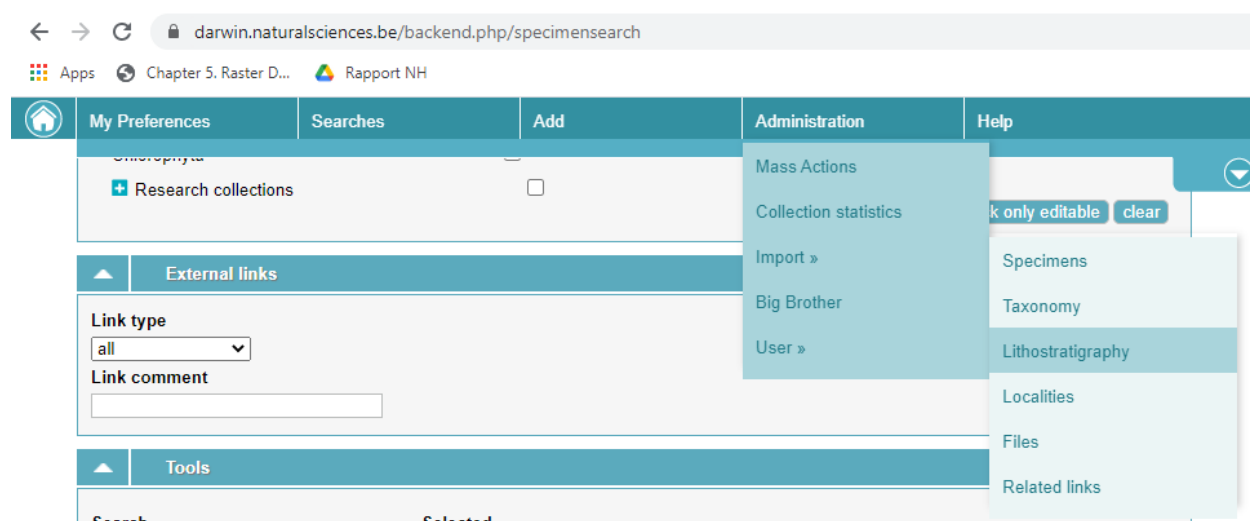


Figure 7. Access to the lithostratigraphy import

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	group	formation														
2	Kibarien moyen	Actuel														
3	Kibarien moyen	Anversien														
4	Kibarien moyen	Anversien														
5	Kibarien moyen	Anversien; ou Messinian (Bolderien) ID Age: 9														
6	Kibarien moyen	Bartonien (Asschien)														
7	Kibarien moyen	Bartonien (Wemmelien)														
8	Kibarien moyen	base miocène moyen														
9	Kibarien moyen	between lower and middle Eocene														
10	Kibarien moyen	Bolderien														
11	Kibarien moyen	Bruxellian														
12	Kibarien moyen	Bruxellian B1 (facies dit panisellen)														
13	Kibarien moyen	Bruxellian sup.														
14	Kibarien moyen	Bruxellien														
15	Kibarien moyen	Campinien q2														
16	Kibarien moyen	Coblenzien, Siegenien sup.														
17	Kibarien moyen	Continental Landenian														
18	Kibarien moyen	Continental Landenian (L2)														
19	Kibarien moyen	Continental Landenian L2														
20	Kibarien moyen	Couches de passage de l'Oligocène au Miocène														

Figure 8. Example of lithostratigraphic import

The template for specimens contains several fields to link specimens to geological classifications. These are:

a. Paleontology / chronostratigraphy

GeologicalEpoch
GeologicalAge
GeologicalAge3

b. Lithostratigraphy

lithostratigraphyGroup
lithostratigraphyFormation
lithostratigraphyMember
lithostratigraphyBed
lithostratigraphyInformalName

c. Mineralogy (identification level)

mineralologicalIdentification
mineralologicalIdentifier
mineralologicalIdentificationYear
mineralologicalIdentificationMonth
mineralologicalIdentificationDay

4.1.4 Embedded multimedia files

A template was also developed to embed multimedia files, that are stored as Media files in the filesystem of the DaRWIN server (hosting the PHP backend) and associated to specimen records.

These files are available in the backend part of DaRWIN which is password protected. This may be relevant for data that may not be publicly disclosed to (.e.g. Material for future paper).

This import consists of two files:

- A Zip file, which is uploaded and unzipped on the server, that contains the files
- A tab-delimited file, that has to be called *meta.txt*, which describes the compressed files

The fields of *meta.txt* are:

Field name	Field Description	Mandatory
UnitID	The main specimen code of the associated specimen	x (or UUID)
filename	the name of the file in the associated ZIP	x
title	The title (caption) of the file	
description	Free text description of the file content	
sub_type	Sub-type of the file	
mime_type	mime type of the file (to help the client to choose the appropriate viewer or player)	x
technical_parameters	the technical parameter of the file (pixel resolution, sampling rate etc...)	
internet_protocol	the internet protocol of "external_uri"	

field_observations	field parameters of the object (e.g. water temperature, salinity etc...)	
external_uri	link to an external resource describing or completing the data	
uuid	the UUID of the associated specimen in DaRWIN	x (if no UnitID provided)

4.1.5 Links to remote multimedia files

Finally, a fifth template was created to associate remote multimedia documents available on the Web to existing DaRWIN, as links that are batch-created. This can be images, sounds, description of the specimen in an on-line publication, related specimens in other databases, links to DNA sequences in GenBank...

DaRWIN also features an IIIF client (Mirador) which is synchronized with this template. It can be minked to the “virtualcol” platform which itself gets the UUID identifiers from DaRWIN. This template allows exchanging data between the two systems.

Field from this template are the followings:

Field name	Description	Values(ex)
UnitID	The main collection code of the specimen in DaRWIN	INV.2090
UUID	The uniform unique identifier of the specimen in DaRWIN (if no UnitID provided)	89f7383b-87c2-47a5-946c-1032bef0ae73
URL	The URL of the resource to link	
Type	The type of link (abbreviate)	DNA; IIIF ; CITES; Nagoya
Comment	Link description (searchable in DaRWIN)	

4.2 Input Templates

4.2.1 Need for an import tool

The DaRWIn web interface is very complete and allows you to enter all data needed concerning specimens. It's very useful when 1 or 2 specimens are to be encoded but it's time consuming if there are a lot of specimens. There is also a need to be online and connected to the DaRWIn server to enter data.

For data coming from outside, there is also no common template: data may come from text files, excel, databases in various formats in data structure and data format. So a tool was needed to work offline (at home or in field work) and to import lots of specimens in one step.

As spreadsheets are known by nearly everyone working with data, this kind of file has been chosen. DaRWIn can also use csv files to import lots of data and a first version of an import template has already been done in excel.

- The new template will use the possibility to create forms above the spreadsheet to more easily enter data that are sometimes spread in many sheets of a workbook.
- Export to csv is easy and automatically done by buttons in the spreadsheet and the generated files can be used to be imported in DaRWIn without other treatment.
- A link to a tool to check taxonomy has also been added in the spreadsheet.

The advantages of this tool and of the forms are an easy way to fill in data, and a view of all the data in one screen, so data can be checked easily for completeness and integrity.

Two versions of the template have been developed:

- a first one in Microsoft Excel and VisualBasic for users using the desktop version of Microsoft Excel on Windows and Mac. This template is not working with Office 365 online or with libreOffice or previous versions of Excel.
- a second template in Calc and basic for users using the Open Source LibreOffice suite on Windows, Mac and Linux OS.

The template offers the possibility to import data in DaRWIn both at RMCA and RBINS. The template may be filled as a simple spreadsheet, by filling in each sheet one after the other but as there is a risk to write data at the wrong place if we choose the wrong line on a sheet, it's better to use the forms that gather all fields of the same line in simple forms.

4.2.2 RMCA Excel template

a. Sheets

Data are splitted into several sheets:

Code, location, DNA, ecology, taxonomy, counts_storage and acquisition.

	A	B	C	D
1		Field form	MRAC user form	Export to Darwin
2		Specimen code	Secondary code	Collection
3				
4				
5				
		Code	Location	DNA
			Ecology	Taxonomy
			Counts_Storage	Acquisition

Figure 9. Different sheets of the template

Code

	Specimen code	Secondary code	Collection	Entered by	Description - Notes
1	SP19-005				
2	SP19-005				

DNA

A	B	C	D	E	F	G	H	I
	Code	Specimen info				Fin-clip info		Notes
	Specimen code	Location Code	Tag number	DNA box	Tube number	Horizontal position	Vertical position	DNA notes
1	SP19-005	Loc-001	19004	23	25	2		6 form wing
2	SP19-005	Loc-001	19004	23	25	2		6 form wing

Ecology

A	B	C	D	E	F	G	H	I
	Code							Ecology parameters
	Specimen code	Location Code	Water temperature	Hour(HH:MM)	pH	mV	Conductivity (µS/cm)	O2 dissolved (%)
1	SP19-005	Loc-001	20	0.524305556	7	45	52	10
2	SP19-005	Loc-001	20	12:35	7	45	52	10

J	K	L	M	N
				Notes
O2 dissolved (mg/l)	hPa	Air temperature	Relative humidity	Ecology notes
6	1020	21	82	very warm
6	1020	21	82	very warm

Acquisition

A	B	C	D	E	F	G	H	I
	Code	Acquisition information				Acquisition dates		Notes
	Specimen code	Location Code	Type	From	Day	Month	Year	Acquisition notes
1	SP19-005	Loc-001	Donation	Mr X	25	12	2012	cadeau
2	SP19-005	Loc-001	Donation	Mr X	25	12	2012	cadeau

Figure 10. Columns of sheets Code, DNA, Ecology, Acquisition

Location

A	B	C	D	E	F	G	H
Code		Location names and description					
Specimen code	Location Code	Continent	Country	State-province	Municipality	Exact site	
1 SP19-005	Loc-001	Africa	Madagascar	Anta	Antanarivo	near Antanarivo	
2 SP19-005	Loc-001	Africa	Madagascar	Anta	Antanarivo	near Antanarivo	

I	J	K	L	M	N	O	P	Q	R
DMS coordinates								Decimal coord.	
Degrees N/S	Minutes N/S	Seconds N/S	N/S	Degrees E/W	Minutes E/W	Seconds E/W	E/W	Latitude	Longitude
12	25	5	S	15	45	12	E		
12	25	5	S	15	45	12	E		

S	T	U	V	W	X
GPS		Collecting information			
GPS Weight Points	Day/Night_catch	Altitude(m)	Collectors	Collecting method	Expedition project
none	Day	1200	Merlijn	Apstein net	Merlijn et al
none	Day	1200	Merlijn	Apstein net	Merlijn et al

Y	Z	AA	AB	AC	AD	AE
Collecting dates						Notes
Start day	Start month	Start year	End day	End month	End year	Locality notes
10	11	1955	11	11	1955	loc notes
10	11	1955	11	11	1955	loc notes

Figure 11. Columns of sheet Location

Taxonomy

A	B	C	D	E	F	G	H	I	J	K	L	M
Code		Taxonomy										
Specimen code	Location Code	Temp. species field name	Kingdom	Phylum	Class	Order	Family	Genus	Species	Subspecies	Author and year	
1 SP19-005	Loc-001	temp taxon	Animalia				Sphingidae	Nephele	densoi		(Keferstein, 1870)	
2 SP19-005	Loc-001	temp taxon	Animalia				Sphingidae	Nephele	densoi		(Keferstein, 1870)	

N	O	P	Q	R	S
Identification			Type	Notes	
Identifier	Day	Month	Year	Type	Taxonomy notes
Jimh	10	12	1956	Specimen, Lectotype	no notes for taxo
Jimh	10	12	1956	Specimen, Lectotype	no notes for taxo

Counts_storage

A	B	C	D	E	F	G	H	I	J
Code		Specimen info		Relationship					
Specimen code	Location Code	Sex	Life stage	Type	Parasite species	Parasite code	Host species	Host code	
1 SP19-005	Loc-001	female	chrysalis	Host	baobab	par001			
2 SP19-005	Loc-001	female	chrysalis	Host	baobab	par001			

K	L	M	N	O	P	Q	R	S	T	U	V
Amounts				Notes	Storage location						
Amount males	Amount females	Amount juveniles	Total number	Notes for amounts	Institution	Building	Floor	Room	Lane	Column	Shelf
1	1	2	2	not sure	RMCA RMCA RMCA	Palais de l'Afrique	Palais de l'Afrique	2 2 3	23 24	2 3 4	3 4 5
1	1	2	2	not sure	RMCA RMCA RMCA	Palais de l'Afrique	Palais de l'Afrique	2 2 3	23 24	2 3 4	3 4 5

W	X	Y	Z	AA	AB	AC	AD	AE
Status - Container - Type of medium								Notes
Status	Specimen part	Container ID	Container type	Container medium	Subcontainer ID	Subcontainer type	Subcontainer medium	Storage notes
Dry---good state---	Fresh---incomplete body head	719 719 720	Cabinet Cabinet Jar	alcohol	199 200	Box Jar	unknown alcohol	
Dry---good state---	Fresh---incomplete body head	719 719 720	Cabinet Cabinet Jar	alcohol	199 200	Box Jar	unknown alcohol	

Figure 12. Columns of sheets Taxonomy and Counts_storage

Two additional sheets are hidden for the user and contain more technical info:

- The first hidden sheet contains predefined lists that can be completed if necessary and that are used in the form combo boxes and lists:

A	B	C	D	E	F	G	H	I	J	K	L	M
Acquisition	Types	Stages	Sex	Countries	Sampling tools	Continent						
Donation	Specimen	adult	male	Algeria	Agassiz trawl	Africa						
Gift	Allotype	subadult	female	Angola	Amphipod Trap	Europe						
Seizure	Cotype	immature	hermaphrodite	Benin	Anchor	Asia						
Purchase	Epitype	juvenile	mixed	Botswana	Angling	North America						
Exchange	Holotype	nestling	undetermined	Burkina Faso	Apstein net	South America						
Loan	Isotype	chrysalis		Burundi	Argos buoy	Oceania						
Expedition	Lectotype	cocoon		Cabo Verde	Artificial substrate frame	Antarctica						
Mission	Neallotype	pupa		Cameroon	Aspirator							
Collect	Neotype	nymph		Central African Republic	Baited traps							
Internal work	Paralectotype	cyst		Chad	balance à crabes							
Excavation	Paratype	nauplii		Comoros	Beam trawl							
Trip	Syntype	caterpillar		Congo, Democratic Republic of the	Beating							
Undefined	Topotype	larva		Congo, Republic of the	Bell Planktometer							
	Voucher	fry		Cote d'Ivoire	Berlese extraction							
		yolk sac larva		Djibouti	Big bottom net							
		prolarva		Egypt	Big Petersen net							
		embryo		Equatorial Guinea	Big thin stramine net							
		newly-hatched		Eritrea	Big trawl with gaule							
		ovum		Eswatini (formerly Swaziland)	Big trawl with gaule and declining irons							
		undetermined		Ethiopia	Big trawl with gaule and thin against-bag							
				Gabon	Bongo net							
				Gambia	Bottle							
				Ghana	Bottle out of glass							

Figure 13. Technical sheet Lists

The second sheet contains info about the mapping with Darwin and Virtual Collections:

	A	B	C	D	E	F
1	Name excel	Field name	Sheet	order	Name darwin	Name Virtual collection
2	Specimen code	TB_Fieldcode	Code	2	UnitID	Code
3	Secondary code	TB_Sec_code	Code	3	additionalID	
4	Collection	TB_Collection	Code	4	collection	Collection in Institution
5	Entered by	TB_Label_createdby	Code	5	label_created_by	
6	General notes	TB_GeneralNotes	Code	6	Notes	Description
7	Full code	TB_FullCode_Joc	Location	2		
8	Location Code	TB_SamplingCode	Location	3	samplingcode	
9	Continent	CB_continents	Location	4	continent	
10	Country	CB_Countries	Location	5	country	Country
11	State-province	TB_Province	Location	6	Province	
12	Municipality	TB_Locality	Location	7	Municipality	
13	Exact site	TB_ExactSite	Location	8	exact_site	Location details
14	Degrees N/S	TB_Deg_NS	Location	9	LatitudeDMSDegrees	
15	Minutes N/S	TB_Min_NS	Location	10	LatitudeDMSMinutes	
16	Seconds N/S	TB_Sec_NS	Location	11	LatitudeDMSSeconds	
17	N/S	CB_N	Location	12	LatitudeDMS_N_S	
18	Degrees E/W	TB_Deg_EW	Location	13	LongitudeDMSDegrees	
19	Minutes E/W	TB_Min_EW	Location	14	LongitudeDMSMinutes	
20	Seconds E/W	TB_Sec_EW	Location	15	LongitudeDMSSeconds	
21	E/W	CB_E	Location	16	LongitudeDMS_W_E	
22	Latitude	TB_Latitude	Location	17	LatitudeDecimal	Coordinates
23	Longitude	TB_Longitude	Location	18	LongitudeDecimal	

Figure 14. Technical sheet Column_matching

b. Buttons

On the first sheet, on top of the sheet, are displayed 4 blue buttons:

A	B	C	D
	Field form	MRAC user form	Export to Darwin
			Check taxonomy

Figure 15. Buttons on first sheet

The 2 first buttons, “Field form” and “MRAC user form” are used to call 2 different forms. The first one contains only a limited set of fields and is intended to be used more in the field whereas the second one contains all the fields corresponding to every column of the 7 sheets. These forms will be described in paragraph 4.2.3.

The next button calls the export features for DaRWIN.

c. Forms

Field form contains basic fields that can be filled in in the field. It includes temporary code and taxonomy and mainly sampling location info. Specimen part may be mentioned and relationships with a host or parasite. Ecological data may also be entered, as well as some info about tissue taken for DNA and amount of specimens.

Majority of the fields are simple text fields and there are some comboboxes prefilled with lists. Some tests are done on data, to check the values. For example, values for coordinates are checked to have values in a range of values.

Figure 16.Field form

Once data are filled in, save the record by clicking on the button “Save record”. This action won’t save the file but will only send the data to the sheets.

If you want to create a new line in the sheets, click on “New record”. It will clear the form and the new data will be saved on a new line. If you want to copy an existing line, go to that line/record with the navigation buttons and click on Duplicate record: it will create a new line with the same data and you will have to change only the necessary data as the code.

The button “Clean content” empties all fields.

c.2 MRAC user form

Use of this second form is similar to the first one. Only the content is different because it contains all the fields.

Because the number of fields is more important and to keep readability of the form, it has been divided in 2 tabs, “General info” and “Secondary info”. General tab contains data about codes, taxonomy, sampling information. Taxonomy is much more detailed and complete taxonomy can be entered as well as the type, author and other info.

Figure 17. First screen of MRAC user form

“Secondary info” tab contains data about ecology, specimen parts, counts, relationship and acquisition. Acquisition is also new in regard to the field form and allows to mention the origin of a specimen other than a collect in the field.

Specimen parts are much more detailed: a container and subcontainer may be defined, with an ID, a type, a medium. Place of the container can be precisely given, as well as state of the specimen part. DNA being considered as a part, DNA info can be given in that section.

A new record has to be created for each part.

Specimen data entry

General info Secondary info Search

Ecology

Measure time:

Water

T (°C): pH: mV: Conductivity (µS/cm):

Dissolved O2 (mg/l): (%):

Air

T (°C): RH (%): P. Atm. (hPa):

Notes:

Specimen parts

Container

ID: Type: Medium:

Subcontainer

ID: Type: Medium:

Storage

Institution: Floor: Column:

Building: Room: Shelf:

Lane:

Specimen: Specimen state: Specimen usage: Specimen part:

Notes:

DNA

Tag number: DNA box:

Tube number:

Horiz. position: Vert. position:

DNA notes:

Counts

Total: ♂ ♀ Juv.:

Sex: Stage:

Notes:

Relationship

This specimen of

(UUID if different from part_of)

Acquisition

Type: Date: DD/MM/YYYY

From:

Notes:

New record Save record * : Mandatory fields

Clean content Duplicate record

Navigation 1/1 << < > >> Go to record n° Go

Figure 18. Second screen of MRAC user form

As for Field form, various combo boxes and lists are already filled with data to facilitate the work of the user.

Specimen data entry

General info Secondary info Search

Specimen info

Specimen code: Secondary code: Entered by: Collection:

Notes:

Taxonomy

Temporary field species name: Author:

Kingdom: Family: Identified by: Date: DD/MM/YYYY Det. St.: Type:

Phylum: Genus:

Class: Species:

Order: Subspecies:

Notes:

Sampling information

Code: Date: DD/MM/YYYY To: DD/MM/YYYY ☐ Day ☐ Night ☐ Dawn ☐ Dusk

Continent: Country: Province:

City: Exact site*:

Coordinates: ° ' " N S ☐ E W GPS weight points:

Latitude (Dec.): Longitude (Dec.):

Collectors:

Expedition: Location notes:

Collecting methods:

- Agassiz trawl
- Amphipod Trap
- Anchor
- Angling
- Apstein net
- Argos buoy
- Artificial substrate frame
- Aspirator
- Baited traps
- balance à crabes
- Beam trawl

New record Save record * : Mandatory fields

Clean content Duplicate record

Navigation 1/1 << < > >> Go to record n° Go

Figure 19. Examples of lists in the form

To make the navigation easier if there are a lot of records, a field “Go to record n°” has been added at the bottom of the form.

“Search” tab contains some fields to do a search in data. Search can be done on the most important fields of each section. It’s only a help to quickly find back one or more lines. If there are results, you can navigate through the results only with the navigation buttons. Click on “Reset” to go back to the whole set of data.

Figure 20. Search tab

d. Export

Data may be exported to DaRWIN.

Three csv files are generated: 1=taxonomy, 2=locations and 3=specimens. These 3 files are imported in DaRWIN in 3 consecutive steps. These files are generated automatically by clicking on the button “Export to Darwin” of sheet 1(Code). A popup window will ask you where to save the files.

Figure 21. Export popup window

d.1 Taxonomy check

The last button on the first sheet is “Check taxonomy”.

It launches a browser to display a web service allowing users to check the taxonomy of the csv file “taxonomy” exported with the button “Export to DaRWIN”. Taxonomy is checked against GBIF, IUCN, WoRMS, Fishbase:

Welcome to the Natural Heritage taxonomy checker

Mail :

Select Tab-delimited to upload: Darwin_im...._taxo.txt

Has header row : ☒

Column index of the name field (first = 1) :

Column index of the kingdom field (first = 1) [optional] :

DARWIN (RBINS): ☐

GBIF: ☒

GBIF (Vernacular names): ☐

IUCN: ☐

WORMS: ☐

Figure 22. Taxonomy check interface

finishedPAGE = 1 NB_PAGES = 1

Match type		Count
EXACT	OTHER	AUTHOR
2	2	2
MISSPELLING		2

[1](#) [Last](#)

class	genus	order	family	phylum	gbif_id	kingdom	species	gbif_url	gbif_rank	gbif_class	gbif_genus	gbif_order	gbif_author	gbif_family	gbif_phylum
	Nephele		Sphingidae		5124095	Animalia	densoi	http://api.gbif.org/v1/species/match?verbose=true&name=Nephele+densoi+	species	Insecta	Nephele	Lepidoptera	Keferstein, 1870	Sphingidae	Arthropoda
	Hippotion		Sphingidae		1862368	Animalia	gerion	http://api.gbif.org/v1/species/match?verbose=true&name=Hippotion+gerion+	species	Insecta	Hippotion	Lepidoptera	Boisduval, 1875	Sphingidae	Arthropoda
	Hippotion		Sphingidae		1862368	Animalia	gerion	http://api.gbif.org/v1/species/match?verbose=true&name=Hippotion+gerion+	species	Insecta	Hippotion	Lepidoptera	Boisduval, 1875	Sphingidae	Arthropoda

Figure 23. Taxonomy check results

4.2.3 LibreOffice template

One of the challenges of the Natural Heritage project is to promote the use of Open Source solutions. The Royal Belgian Institute of Natural Sciences decided to evaluate a free open source tool as the MS-Excel template is highly dependent on the MS Office version and request an A5 licence as the Office 365 online is not compatible with the Visual basic macros.

The choice was made to use LibreOffice with the use of basic as macro language. The option to enter media files data is also not yet developed in this version.

a. Sheets

The use of the LibreOffice template is the same as the excel template.

Data are spread into 5 sheets and a form allows us to fill all the sheets together, for one record. The number of sheets is smaller because the data are organized differently.

Here are the columns of the 5 sheets:

Taxonomy

A	B	C	D	E	F	G	H	I	J	K
Code	Specimen code	Temp. species field name	Family	Author	Genus	Author	Species	Author	Subspecies	Author
Input form	test011	Felis à points notes	Echimyidae	Gray	Thrichomys	Trouessart, 1880	Thrichomys apereoides	(Lund, 1839)		
	test021	Big ape	Pongidae	Elliot	Pan	Oken, 1816	Pan troglodytes	(Blumenbach, 1775)		
Export to Darwin files	test031									
	test041									
Check taxonomy	test051									

L	M	N	O	P	Q	R	S	T	U	V
Identification	Type	Notes	Relationship	Host	Parasite species	Parasite code	Host species	Host code		
Yvrey Emmanuel	13	03	1955	Cotype	perhaps new species	Camptocaryus s	rmca1023s			
Yvrey Emmanuel	10	12	1953	Allotype	perhaps new species					
	10	12	1953	Allotype-Holotype	perhaps new species					

Specimen

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
Code	Specimen code	Secondary code	Collection	IG	Entered by	Media URL	General notes	Type	From	Day	Month	Year	Notes	Sex	Life stage	Males	Females	Juveniles	Total	Notes
	test011	B-023	coll1s	IG1s	JM1s	file1s	No general notes	Donation	Achille	08		1726	Given by director	hermaphrodite	subadult	2	14	3	20	2 males
	test021	C-019	coll2	IG2	JJ2	file2, file3	notes	Acquisition	Acena	10	12	2015	3 months after death of owner	Female	juvenile	1	1	1		Juvenile probably female
	test031																			
	test041																			

Location

B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S		
Code		Location names and description																DMS coordinates	
Specimen code	Station number	Continent	Country	Original country name	State/province	Municipality	Ocean	Sea	Exact site	Degrees N/S	Minutes N/S	Seconds N/S	Degrees E/W	Minutes E/W	Seconds E/W	E/W			
test011	locode011	Asia	Bahrain	old Bahrain	Miratus	Riffa	Southern Ocean	Bellingshausen Sea	center of country	1	1	1,280	S	1	1	1,500	W		
test021	locode021	Europe	Italy		Roma	Roma city	Atlantic ocean	Gulf of Guinea	center of country	12	4	50	S	5	53	47	E		
test031	locode031	Africa	Democratic Republi	Congo belge		matadi	Atlantic ocean		near road to Kinshasa										
test041	locode041	Europe	Belgium			brussels			test										
T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	AH	AI	AJ	AK		
Decimal coord.		Other infos																	
Latitude	Longitude	Origin of coordinates	Web service	Web service ID	Web service url	GPS Weight Points	Line or area in WKT	Coordinates accuracy(m)	Original coordinates	From	To	Accuracy	From	To	Accuracy	Height	Accuracy		
		Historical-Label				no points1	Point(4 2,5 311)	1,25	org in DMS1	100	120	5	10	12	1,23	44	0,2		
						no points2	Point(4 2,5 4)	100	org in DMS2	102	104	2	21	27	1	77	3		
-5.82566	13.46090	Calculated from web service	OpenStreetMap Nominatim	488235592	https://nominatim.org					10	20	2,00							
50.84656	4.35170																		

Figure 24.Columns of sheets Taxonomy, Specimen, Location

B	C	D	E	F	G	H	I	J	K	L	M	N
Code	Collecting info				Collecting time				Notes			
Specimen code	Collectors	Collecting method	Expedition project	Start day	Start month	Start year	Start hour	End day	End month	End year	End hour	Locality notes
test011	Abbas	Apstein net, Artificial substrate frame	Project Angola	11	1956	12:13	1957	14:16	locality near the sea			
test021	Aaron	Apstein net, Aspirator	Italy 2020	01	1956			12	1956			
test031	Abbas	Argos buoy, Aspirator	Congo 2020	01	2020	10:15	20	03	2020	12:00	Attention, coordinates retrieved from Openstreetmap based on city and country	
test041		Apstein net, Aspirator		01	01	2015		12	1956			
O	P	Q	R	S	T	U	V	W	X	Y	Z	AA
Ecology parameters												
Hour of measure(HH-MM)	Water temperature	pH	Salinity	Tot. susp. solids (mg/l)	Streamflow	mV	Conductivity (µS/cm)	O2 dissolved (%)	O2 dissolved (mg/l)	Air pressure (hPa)	Air temperature	Relative humidity
12:45	21	6	0.00	7.00	1200	1200	120	45	10	1020	35	85.00
12:00												85.00
15:45												85.00
12:00												85.00
AB			AC			AD			AE			AG
Ecology notes			Biogeog. realm			Terrestrial			Biosphere			Marine
found in forest	NA: Neartic		Biogeographic realm	Marine		Biosphere	Terrestrial		Biosphere	Freshwater		Marine
found in forest	NA: Neartic		Eastern Indo-Pacific			Tundra			Temperate floodplain rivers and wetlands			Littoral/Intertidal zone
found in forest	NA: Neartic		Eastern Indo-Pacific			Tundra			Temperate floodplain rivers and wetlands			Littoral/Intertidal zone
found in forest	NA: Neartic		Eastern Indo-Pacific			Tundra			Temperate floodplain rivers and wetlands			Littoral/Intertidal zone

A	B	C	D	E	F	G	H	I	J	K	L		
	Code		Storage location										
	Specimen code	Institution	Building	Floor	Room	Lane	Column	Shelf	Status	Specimen part	Container ID		
1	test011	RBINS RBINS RBINS RBINS RBINS RBINS	Aile des dinos aleb Main building	12 5 14	12 2 6	12 2 4	12 2 6	12 2 8	fresh—good state—Loan fresh—good state— ----- ----- head body leg arm part6	Container01			
2	test021	RMCA RMCA	CAPA1 Palais de l'Afrique3	112	112	112	112	112	fresh—good state—Loan fresh—good state— ----- ----- head body				
3	test031	RBINS RBINS RBINS RBINS RBINS RBINS	Aile des dinos aleb Main building	12 5	12 2 6	12 2 4	12 2 6	12 2 8	fresh—good state—Loan fresh—good state— ----- ----- head body leg arm				
4	test041	RMCA											
	M	N	O	P	Q	R	S	T	U	V	W	X	
	Status - Container - Type of medium			Subcontainer ID		Subcontainer type		Subcontainer medium		Notes			
	Container type	Container medium	Subcontainer ID	Subcontainer type		Subcontainer medium		Storage notes					
	Tag number	DNA box	Tube number	DNA		Horizontal position		Vertical position		Notes			
Microscopic slide Alkohol flask jar	alkohol algargl alkohl	SubID012 SubID22 None	Microscopic slide Alkohol flask jar	alkohol algargl alkohl		storage in boxes imported		tag011	box11	tube11	21	41	no details
Microscopic slide Alkohol flask jar	alkohol algargl alkohl	SubID012 SubID22	Microscopic slide Alkohol flask	alkohol algargl alkohl		storage in boxes imported		tag011	box11	tube11	21	41	no details
Microscopic slide Alkohol flask jar	alkohol <u>alugro</u> alkohl	SubID012 SubID22 None					tag011	box11	tube11	21	41	no details	

Figure 25. Columns of sheets Sampling and Storage

As for the excel template, it's much easier to enter data via the form. There is only one form here (no field form). So there are 3 buttons on the first sheet:

	A	B
1		Code
2	Input form	Specimen code
3		1 test011
4		2 test021
5	Export to Darwin files	3 test031
6		4 test041
7	Check taxonomy	5 test051

Figure 26. Buttons on the first sheet.

A first one to open the form, a second to export files to DaRWIn and the last one to check taxonomy.

Use of the LibreOffice form is the same as for the excel template. Navigation buttons allow you to go from one line/record to another and 4 buttons allow you to save, duplicate, clean and add a record (these buttons are displayed in this form as icons).

As in the excel template, data are displayed on several screens (5 screens here in place of 2 in excel). To go from one screen to the other, click on the big buttons on the right. A special button “Search” shows a screen where you can search data.

First screen is Taxonomy. It contains info about taxonomy, type, identification, interspecies relations.

The screenshot shows the first screen of the LibreOffice form. The main section is titled "Taxonomy" and contains several input fields: "Temporary taxon name" (Felis à points noirs), "Family" (Echimyidae), "Genus" (Thrichomys), "Species" (Thrichomys apereoides), "Subspecies", "Author" (Gray), and "Date" (13/03/1955). There is also a "Type" dropdown menu with options like Cototype, Epitype, Holotype, Isotype, Neotype, Paratype, Syntype, Topotype, and Voucher. A "Notes" field contains the text "perhaps new speciessss".

On the right side, there is a "Navigation" sidebar with buttons for "Taxonomy", "Specimen", "Location", "Sampling", and "Storage". The "Taxonomy" button is highlighted in green. Below the sidebar, there are icons for "Save", "New", "Duplicate", and "Clean form", and a "Search" button.

Annotations with arrows point to specific elements: "Navigation buttons" points to the navigation sidebar, "Buttons to change section" points to the sidebar buttons, "Buttons for actions on record: save, New, Duplicate, Clean form" points to the icons below the sidebar, and "Button to search" points to the "Search" button.

Figure 27. First screen of the LibreOffice form. Buttons.

Second screen is Specimen and contains info about codes, acquisition, sex, stage, counts.

The screenshot shows the second screen of the template, titled "Specimen". It contains three main sections: "Specimen info", "Acquisition", and "Sex - Stage - Counts".

Specimen info section includes fields for "Collection" (coll1s), "Entered by" (jims), "Media URL" (file1s), "I.G." (IG1s), "Main code" (test011a), and "Sec. code" (B-023). There is also a "General notes" field with the text "No general notes".

Acquisition section includes fields for "Type" (Donation), "Date" (08/1726), and "From" (Achille). There is also a "Notes" field with the text "Given by director".

Sex - Stage - Counts section includes fields for "Sex" (hermaphrodite), "Stage" (subadult), and "Counts" (Total: 20, Males: 2, Females: 14, Juveniles: 3). There is also a "Notes" field with the text "2 males".

On the right side, there is a "Navigation" sidebar with buttons for "Taxonomy", "Specimen", "Location", "Sampling", and "Storage". The "Specimen" button is highlighted in green. Below the sidebar, there are icons for "Save", "New", "Duplicate", and "Clean form", and a "Search" button.

Figure 28. Second screen of the template, Specimen

Third screen is Location and contains info about exact geographic place, coordinates, altitude, depth and more technical data.

Station number* loccode11

Exact site* center of country

Location info

Continent: Asia Ocean: Southern Ocean

Country* Bahrain Sea: Bellingshausen Sea

Original country name: old Bahrain

Province: Mirutus

Municipality: Riffa

Coordinates

Latitude/Latitude

Latitude: 1° 1' 1,280" N ☐ S ☒ Decimal

Longitude: 1° 1' 1,500" E ☐ W ☒ Decimal

Accuracy: 1,25 m

Origin of coordinates: Historical-Label

If more than 1 point, decimal coordinates in WKT: Point(4 2,5 311)

GPS weight points: no points1 Orig. coord.: orig in DMS1

Navigation: Record nr 1/13 << < > >>

Taxonomy

Specimen

Location

Sampling

Storage

Search

Figure 29. Third screen of the template, Location

A special function has been added in this LibreOffice version: it's possible to get the coordinates of a place, based on the country and the municipality. It works with a webservice as OpenStreetMap. When you have filled in the country and municipality, click on the button "Get coordinates". If it can find coordinates, they are written in decimal latitude and longitude and the origin of the coordinates is written below: calculated from OpenStreetMap.

Latitude/Latitude

Latitude: 0° 0' 0,000" N ☐ S ☐ Decimal: 41.89332

Longitude: 0° 0' 0,000" E ☐ W ☐ Decimal: 12.48293

Accuracy: 100 m

Origin of coordinates: Calculated from web service of OpenStreetMap Nominatim

Figure 30. Coordinates fields of the form.

A fourth screen is the sampling. It contains all data about how and when specimens were collected and also local ecological info, as well as larger ecological info in Biogeography.

Sampling

Expedition: Project Angola

Collector(s): Abbas

Sampling methods: Apstein net, Argos buoy, Artificial substrate frame

Date/Time: From 11/1956 12:13 to 1957 14:16

Notes: locality near the sea

Ecology

Measure Time: 12:45

Air: T°C 35, RH (%) 85.00, ATM (hPa) 1020

Water: Streamflow 1200, Tot. susp. solids (mg/l) 7.00, T°C 21, O2 dissolved mg/l 10, % 45

Parameters: Salinity 0.00, pH 6, Conductivity (mS/cm) 120, mV 1200

Notes: found in forest

Biogeography

Realm: Terrestrial and freshwater NA: Nearctic, Marine Eastern Indo-Pacific

Biome: Terrestrial Tundra, Freshwater Temperate floodplain rivers and w, Marine Littoral/Intertidal zone

Navigation: Record nr 1/13

Taxonomy, Specimen, Location, Sampling, Storage

Search

Figure 31. Fourth screen of the template, Sampling

The last screen is Storage. As for the excel, many parts may be added by clicking on the “Add a part”.

Part1

Add a part

Part

Specimen part: head1

Preparation: fresh

State: good state

Usage: Loan

Notes: storage in boxes imported

Storage

Institution: RBINS

Building: Aile des dinos

Floor: 1

Room: 1

Lane: 1

Column: 1

Shelf: 1

Container ID: Container01

Type: Microscopic slide

Medium: alcohol

Sub-container ID: SubID12

Type: Microscopic slide

Medium: alcohol

DNA tissue

Tag: tag011

Box: box11

Tube: tube11

Hor. position: 21

Vert. position: 41

Notes: no details

Navigation: Record nr 1/13

Taxonomy, Specimen, Location, Sampling, Storage

Search

Figure 32. Fifth screen of the template, Storage

There is a special screen to do a search in data. Search can be done on the most important fields of each section. It's only a help to quickly find back one or more lines. If

there are results, you can navigate through the results only with the navigation buttons. Click on “Reset filters” to go back to the whole set of data.

Figure 33. Search screen of the template

d. Export

Button “Export to Darwin files” on first sheet has the same function as in the excel, to export files that have to be imported in DaRWIN

e. Taxonomy check

Button “Check taxonomy” has the same function as in the excel template.

4.3 Integration of previous databases and Import of data in DaRWIN

The import procedure is a very important tool for the specimen data:

The total of records at the RMCA DaRWIN is 695.704 records, corresponding to 1.912.700 specimens. The import procedure allowed us to import 259.883 records which represent 37,3 % of the total of the database.

The total of records at the RBINS DaRWIN is 696.446 records, corresponding to 4.360.000 specimens of which 73.800 were imported with the excel template (10,6 %).

4.3.1 Mapping of the RBINS MISTA database

Data from the RBINS MISTA database (polar missions in the Antarctic), originally in Microsoft Access format have been mapped to the DaRWIn importation templates and would be ready to be imported into DaRWIn after check from the scientists. This database was chosen as a case study for the Natural Heritage project as it was a very complex database including specimens from different institutions.

In contrast to other collections, MISTA has taken more time because data were in a very complex database and it has been difficult to extract data correctly. MISTA also influenced the original development of DaRWIn, for example by putting dates in the geographical data because data of MISTA contain a majority of specimens caught in sea (Antarctica) and during cruises. A catch can begin at a point A at day J and end at a point B at day J+3, which is not the case with terrestrial specimens caught at a precise place. Stations and expeditions were also as important in the original database as specimens and because of this, it influenced the development of an import in 3 steps: taxonomy, localities and specimens. This allows to import first a list of stations based on cruises.

The data are now imported in a working study in DaRWIn and will be integrated in the main collections after a final checking of the data.

A total of 5.831 records of MISTA corresponding to 132.190 specimens were imported. This is the equivalent of 1 FTE of manual encoder during 1 Year.

4.3.2 The RBINS Geology Collection

The collections of Geology use Microsoft Access as a database. A previous attempt of import in RBINS DaRWIn was made in 2017 but the import was cancelled thanks to import errors in the validation of sampling locations.

It was now possible to import the data again using the 3 steps import procedure.

A total of 40.069 records are now available in the DaRWIn Collection Management System.

4.3.3 The RBINS Paleontology Collection

The collections of Paleontology also use Microsoft Access as a database.

Data was exported as XLS files and templates were prepared for further import in DaRWIn. The process is not yet completed as we need first to control the chronostratigraphy reference system existing in the main database.

More than 40.000 type specimens will be imported with this procedure in 2021.

4.3.4 RMCA zoology

Vertebrate data from the DataPerfect/Drosera system have been imported by using the template in 2019 and 2020.

This concerns the following collections:

- 21.216 records of Reptilia
- 16.069 records of Amphibians
- 16.000 records of Ornithology
- 451 records of ichthyology

Data from the Invertebrate collections were also imported:

- 1.276 records of trichoptera
- 28.936 records of Acari
- 185 records of Ephemeroptera
- 11.830 records of Crustacea
- 9.220 records of Myriapoda
- 1.541 records of Vermes
- 67.747 records of Coleoptera
- 1.837 records of Echinodermata
- 275 records of snails

4.3.5 Mapping of RMCA wood biology data

Data from the wood biology department have been successfully imported into DaRWIN in 2020, using the tab-delimited templates presented above. These data (83300 records) were originally conserved in offline Excel format.

Authors: Jean-Marc Herpers, Franck Theeten, Marielle Adam & Patrick Semal