

How to create a Darwin Core Archive

Luke Marsden (lukem@met.no)

January 17, 2024

Contents

1	What is a Darwin Core	2
2	What is a Darwin Core Archive	2
2.1	Why do we use multiple CSV files?	3
3	Darwin Core terms and other controlled vocabularies	3
4	GBIF, OBIS and a network of data	3
5	Cores and extensions	4
6	Which columns are required in each extension?	5
7	How to create a Darwin Core Archive	5
7.1	Creating the CSV files	5
7.2	Creating a DwCA from your CSVs	6
8	Making your data available via SIOS	6
9	Citing your data in your paper	6

1 What is a Darwin Core

Darwin Core is a data standard originally developed for biodiversity informatics, though this has expanded to be useful for any type of data where you have data associated with one or a list of organisms. Darwin Core includes

- Darwin Core terms: A controlled vocabulary of terms - <https://dwc.tdwg.org/terms/>
- Darwin Core Archive: A FAIR data format

2 What is a Darwin Core Archive

A Darwin Core Archive (DwCA) is a self-describing dataset for taxonomic (species) data and sampling event data. It consists of one or more data tables (CSV files) and 2 XML files, one (meta.xml) that describes how the files are organised and a second (eml.xml) that provides the metadata describing the dataset as a whole. They are zipped together to create the Darwin Core Archive (DwCA) (Figure 1).

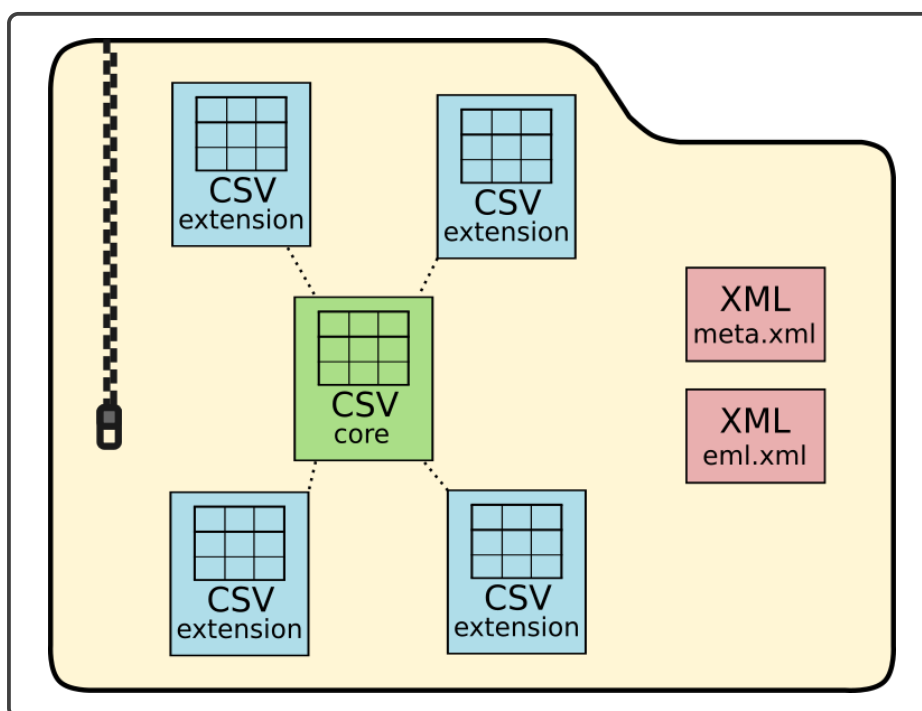


Figure 1: Visualisation of a Darwin Core Archive, portrayed using the star schema. The central event core can be surrounded by zero or many extension tables. It also contains a meta.xml file that describes what columns each CSV contains and links them to the term in a controlled vocabulary, and an eml.xml file that provides metadata that describes the dataset as a whole. They are zipped together to create a Darwin Core Archive.

Each CSV file contains a number of rows, and every row has its own unique ID (Figure 2). These IDs are used to link the CSV files together. Every row in every 'extension' CSV file must include the ID of one row in the 'core' CSV file. This is called the 'star schema' and you can think of the core as the centre of the star. Note that a DwCA can only include a single 'level' of extensions. In other words, it cannot include an extension to an extension file.

It has a single core in the centre of the star, for example event records, where a single row corresponds to a single event. This central core can then optionally be surrounded by extension tables, linked to the central core using this ID.

Event Core

eventID	samplingProtocols	eventDate	...
UUID 1	Bongo net 180 um	2021-09-24T12:17:14Z	
UUID 2	Bongo net 180 um	2021-09-27T09:28:04Z	

Can be hierarchical with eventID and parentEventID, e.g. expedition (parent) and sampling activities (children).

Occurrence Extension

occurrenceID	eventID	scientificName	...
UUID 3	UUID 1	Boreogadus saida	
UUID 4	UUID 1	Gadus morhua	
UUID 5	UUID 2	Boreogadus saida	

Extended Measurement Or Facts Extension

measurementID	eventID	occurrenceID	measurementType	measurementValue	measurementUnits
UUID 6	UUID 1	UUID 3	Fork length	34	cm

Best practice: also include measurementTypeID to make the measurement machine readable

Figure 2: Example of a Darwin Core Archive with an Event Core and 2 extensions; an Occurrence Extension and a Extended Measurement Or Facts Extension. Not that this example does not include all the required columns for each core or extension.

For more information on DwCA, see <https://ipt.gbif.org/manual/en/ipt/2.5/dwca-guide>.

2.1 Why do we use multiple CSV files?

This is a common question, but we are not trying to make things needlessly complicated! This allows a many-to-one relationship to be logged, for example multiple species (occurrences) logged for a single sampling event. Hierarchical information can be stored clearly, so the data user can understand which samples were collected from the same sampling event.

This method is also more efficient. Certain metadata are consistent between the sampling event and all the samples collected from it, e.g. time, date, coordinates etc. These metadata can be logged only once for each sampling event in an event core. They therefore do not need to be included for every single sample, which would lead to a lot of duplication in some cases!

3 Darwin Core terms and other controlled vocabularies

Darwin Core includes a controlled vocabulary of terms for sharing information about biological diversity. Most (but not all) of the column headers you will include in your CSV files in your DwCA will be Darwin Core terms. Other terms will be from other controlled vocabularies.

We will come back to which terms you should include later.

4 GBIF, OBIS and a network of data

There are a number of data centres and services that manage Darwin Core data. GBIF (Global Biodiversity Information Facility) is the largest. Most (all?) of the other Darwin Core services make their data available via GBIF as well as their own platform. So if you publish your data with any of the below services, your data will also be available via GBIF.

Table 1: A network of Darwin Core data. This table is just a small selection of the services available.

Name	Website	Comment
OBIS	https://obis.org/	Only marine data. The OBIS community have been pushing the DwC standards forwards to build better functionality for scientific data.
iNaturalist	https://www.inaturalist.org/	For citizen science, nature enthusiasts and researchers. Offer some great apps like Seek that you can use on your mobile phone for taking pictures, identifying the organism and publishing the data https://www.inaturalist.org/pages/seek_app
Living Norway	https://livingnorway.no/	Norwegian ecological data project
Artsdatabanken	https://www.artsdatabanken.no/	Service for collecting, organizing, and disseminating data related to Norwegian flora and fauna

5 Cores and extensions

Which cores and extensions should you include in your DwCA?

For most scientific data, you should include an 'Event Core', where 1 row is 1 sampling event. This can also be hierarchical, using eventID and parentEventID (Figure 3).

https://rs.gbif.org/core/dwc_event_2022-02-02.xml

Event Core				
eventID	parentEventID	samplingProtocol	eventDate	...
UUID 1		Research cruise	2007-11-13/2007-12-05	
UUID 2	UUID 1	Bongo net 180 um	2007-11-27T09:28:04Z	
UUID 3	UUID 1	Bottom trawl	2007-12-02T19:01:32Z	

Figure 3: Example of a hierarchical event core. The research cruise is the parent of the children sampling activities. This is shown using the eventID and parentEventID columns.

You will most likely also have an occurrence extension, where one row is an observation of an organism or group of organisms of the same species.

https://rs.gbif.org/core/dwc_occurrence_2022-02-02.xml

Measurements related to the sampling event or occurrence (e.g. mass of your organism) go in an Extended Measurement or Facts extension.

https://rs.gbif.org/extension/obis/extended_measurement_or_fact_2023-08-28.xml

If you have material samples (soil, faeces) that you have taken you can log them using a Material Sample extension.

<https://rs.gbif.org/extension/ggbn/materialsample.xml>

If you have associated media you can refer to them using a Simple Multimedia Extension, and publish the media elsewhere. It is possible to publish these with GBIF too.

<http://rs.gbif.org/extension/gbif/1.0/multimedia.xml>

If you need to relate a row in one extension to a row in another extension, you can use the Resource Relationship Extension.

https://rs.gbif.org/extension/dwc/resource_relationship_2022-02-02.xml

And more extensions that are registered with GBIF can be found here.

<https://rs.gbif.org/extensions.html>

6 Which columns are required in each extension?

It can be difficult to know which columns to include and what the minimum requirements are. Unfortunately, this information is not easy to find online. You can contact the data centre that you are publishing your data with.

I have also developed the Nansen Legacy spreadsheet template generator for Darwin Core Archives that has requirements and recommendations for different cores and extensions included.

<https://www.nordatanet.no/aen/template-generator/config%3DDarwin%20Core>

If you notice any issues with this or want to see something else added, please get in touch, or better still, raise an issue on GitHub.

https://github.com/SIOS-Svalbard/Nansen_Legacy_template_generator

7 How to create a Darwin Core Archive

To create a Darwin Core Archive you first have to create the CSV files. Then, you can use the Integrated Publishing Toolkit (IPT), developed by GBIF, to turn them into a Darwin Core Archive.

7.1 Creating the CSV files

You could use the Nansen Legacy template generator to help you with this. This will help you create a spreadsheet with a separate sheet for each core and extension (Figure 4). The descriptions for each term will appear as notes each time you select a cell to help you fill it in. I advise that you export each sheet to separate CSV files before you proceed with creating the Darwin Core Archive.

Alternatively, you can create CSV files in any of way that you prefer. You can go to <https://rs.gbif.org/extensions.html> to see which cores and extensions you can use and which terms you can include in each. Note that it does not tell you which terms are required. This varies between data centres - GBIF has slightly different minimum requirements to OBIS for example for certain extensions.

7.2 Creating a DwCA from your CSVs

Once you have your CSV files, creating a DwCA is easy. You can using the integrated publishing toolkit (IPT), developed by GBIF, to create the DwCA and also publish it.

Here is a map of places that have an IPT installed <https://www.gbif.org/ipt>. Choose one, contact them and ask for login details to their node.

For most people in Norway, good choices are the Norwegian Marine Data Centre (operated by IMR) - datahjelp@imr.no, or GBIF Norway - helpdesk@gbif.no

8 Making your data available via SIOS

Data relevant to Svalbard should be available via the SIOS data access portal https://sios-svalbard.org/metsis/search?f%5B0%5D=dataset_level%3Alevel-1

If you have published your data with NMDC, they contribute to SIOS so the data will be made available via the SIOS. Just let them know that you want your data to be available via SIOS.

GBIF Norway need to link manually. You can use the metadata collection form on SIOS to do this. You will need to log in first. <https://sios-svalbard.org/metadata-collection-form>

9 Citing your data in your paper

Cite your paper just as you would cite any other scientific publication - in your list of references. You can also mention the data in a data availability statement if your chosen journal requires one, but this should be as well as (not instead of) including the data in your list of references.

The recommended citation can be seen on the landing page of the dataset in the data centre you chose to publish with (GBIF or NMDC most likely).

Configuration:

Darwin Core

Select

Sub-configuration:

Sampling Event

Select

Create spreadsheet templates that can be used to create Darwin Core cores and extension CSVs. Each sheet below represents one core or extension.

Tutorial video: <https://www.youtube.com/watch?v=DvYwnYXuPU>

Add Extended MoF Extension

Add Material Sample Extension

Add Resource Relationship Extension

Add Simple Multimedia Extension

Sheet: Event Core

The category of information pertaining to an action that occurs at some location during some time.
https://rs.gbif.org/core/dwc_event_2022-02-02.xml

Required	Recommended	Suggestions	Other
<input checked="" type="checkbox"/> eventID	<input checked="" type="checkbox"/> sampleSizeValue	<div>Suggestions</div>	<div>Add CF standard names</div>
<input checked="" type="checkbox"/> eventDate	<input checked="" type="checkbox"/> sampleSizeUnit		<div>Add Darwin Core terms</div>
<input checked="" type="checkbox"/> samplingProtocol	<input checked="" type="checkbox"/> parentEventID		<div>Add more fields</div>
	<input type="checkbox"/> samplingEffort		
	<input checked="" type="checkbox"/> locationID		
	<input checked="" type="checkbox"/> decimalLatitude		
	<input checked="" type="checkbox"/> decimalLongitude		
	<input type="checkbox"/> geodeticDatum		
	<input type="checkbox"/> footprintWKT		
	<input type="checkbox"/> footprintSRS		
	<input type="checkbox"/> countryCode		

Sheet: Occurrence Extension

The category of information pertaining to the existence of an Organism (sensu <http://rs.tdwg.org/dwc/terms/Organism>) at a particular place at a particular time.
https://rs.gbif.org/core/dwc_occurrence_2022-02-02.xml

Required	Recommended	Suggestions	Other
<input checked="" type="checkbox"/> occurrenceID	<input type="checkbox"/> taxonRank	<div>Suggestions</div>	<div>Add CF standard names</div>
<input checked="" type="checkbox"/> eventID	<input type="checkbox"/> kingdom		<div>Add Darwin Core terms</div>
<input checked="" type="checkbox"/> basisOfRecord	<input checked="" type="checkbox"/> individualCount		<div>Add more fields</div>
<input checked="" type="checkbox"/> scientificName	<input checked="" type="checkbox"/> organismQuantity		
	<input checked="" type="checkbox"/> organismQuantityType		
	<input type="checkbox"/> occurrenceStatus		

Generate

	A	B	C	D	E	F	G	H	I
1									
2	Required								
3	Recommended								
4	Other fields								
5	CF standard name								
6	Darwin Core terms								
7	Use 'paste special' / 'paste only' so not to overwrite cell restrictions								
8									
9	eventID	eventDate	samplingProtocol	sampleSizeValue	sampleSizeUnit	parentEventID	locationID	decimalLatitude	decimalLongitude
10									
11									
12									
13									
14									
15									
16									
17									
18									
19									
20									
21									
22									
23									
24									
25									
26									
27									

Event Core Occurrence Extension Conversions README

Figure 4: The Nansen Legacy template generator (a) and an example template (b)