

The background of the slide features a large, blue-tinted iceberg floating in a dark blue ocean under a sky filled with orange and yellow hues from a setting sun.

# Using a template to structure data: practical tips and tricks

BIODIVERSITY.AQ

Biodiversity data from the field to research

Yi-Ming Gan, Maxime Sweetlove, Anton Van de Putte

SCAR Antarctic Biodiversity Portal





Let us know if you prefer not to be recorded.

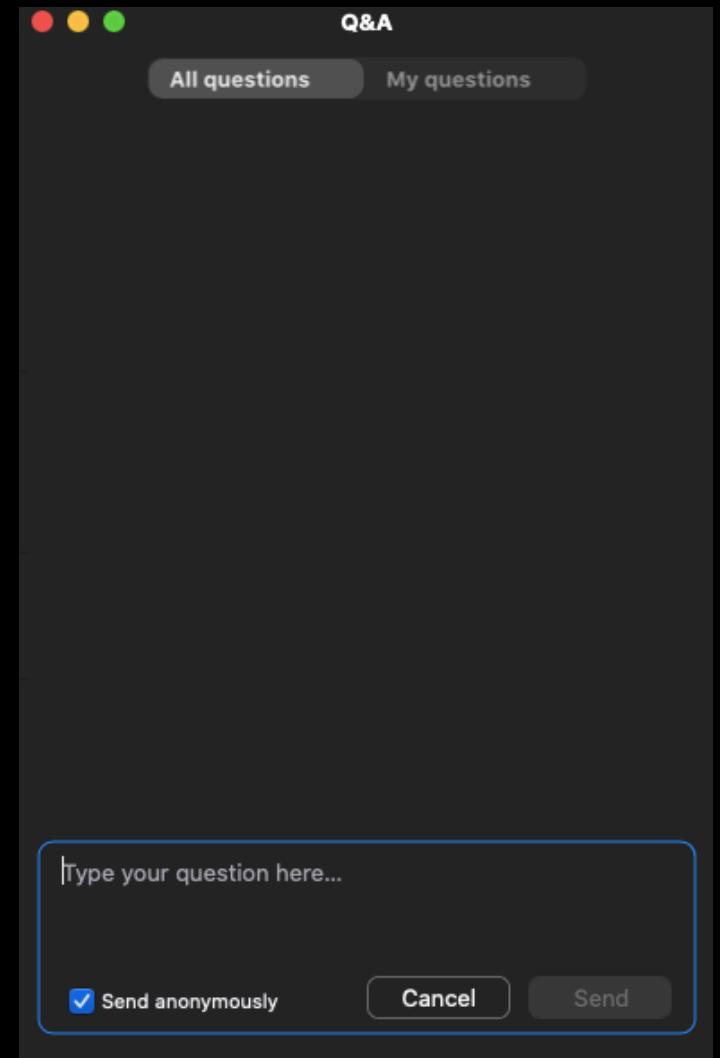
# Code of Conduct

- Be respectful
- We will follow the principles of the rOpenSci Code of Conduct
  - <https://ropensci.org/code-of-conduct/>

# Using Zoom in this webinar



Ask questions using Q&A feature  
or  
raise your hand



Mute

Chat

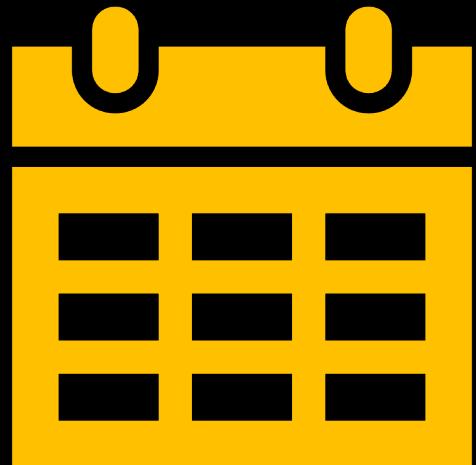
Raise Hand

Q&A

Leave

These webinars won't cover everything

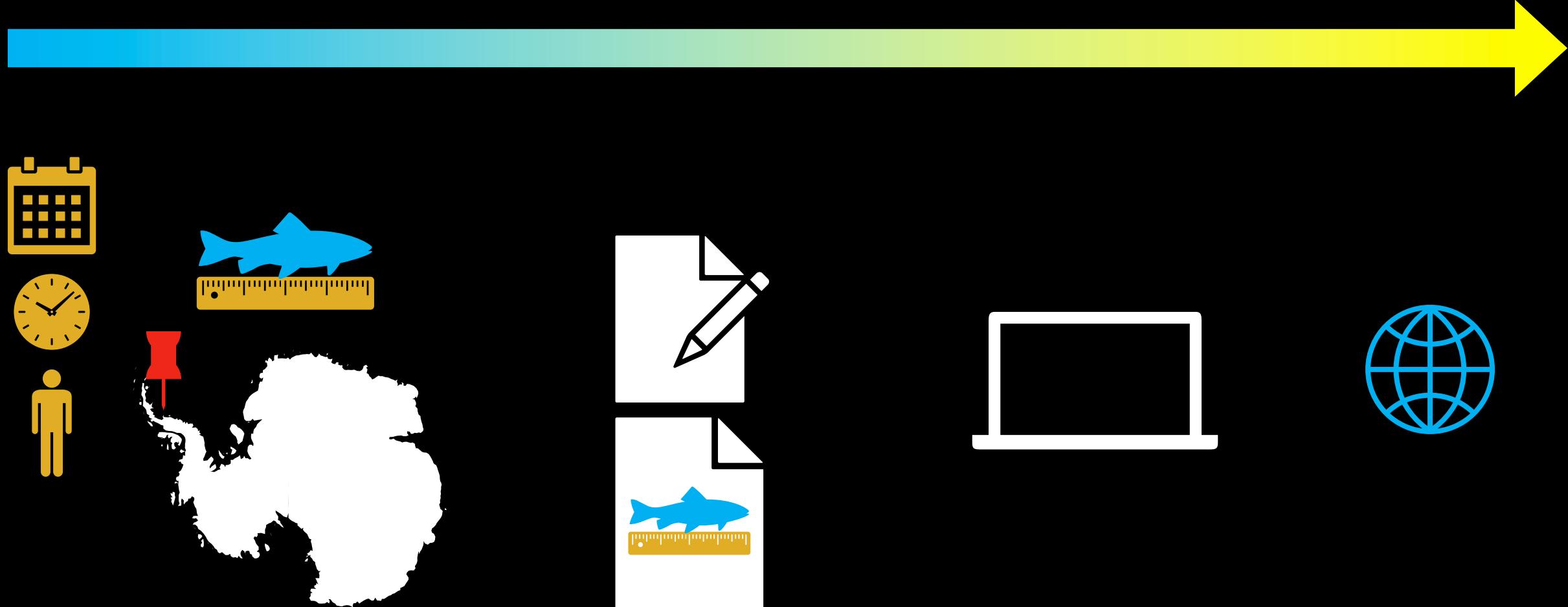
- Discover more based on the links that we provide
- Schedule a session to discuss directly
  - <https://doodle.com/meetme/qc/XxYYpJwmbG>



- Physical observation to digital data
- Example: sampling at sea
- Example: genomics data



# Physical observation to digital data



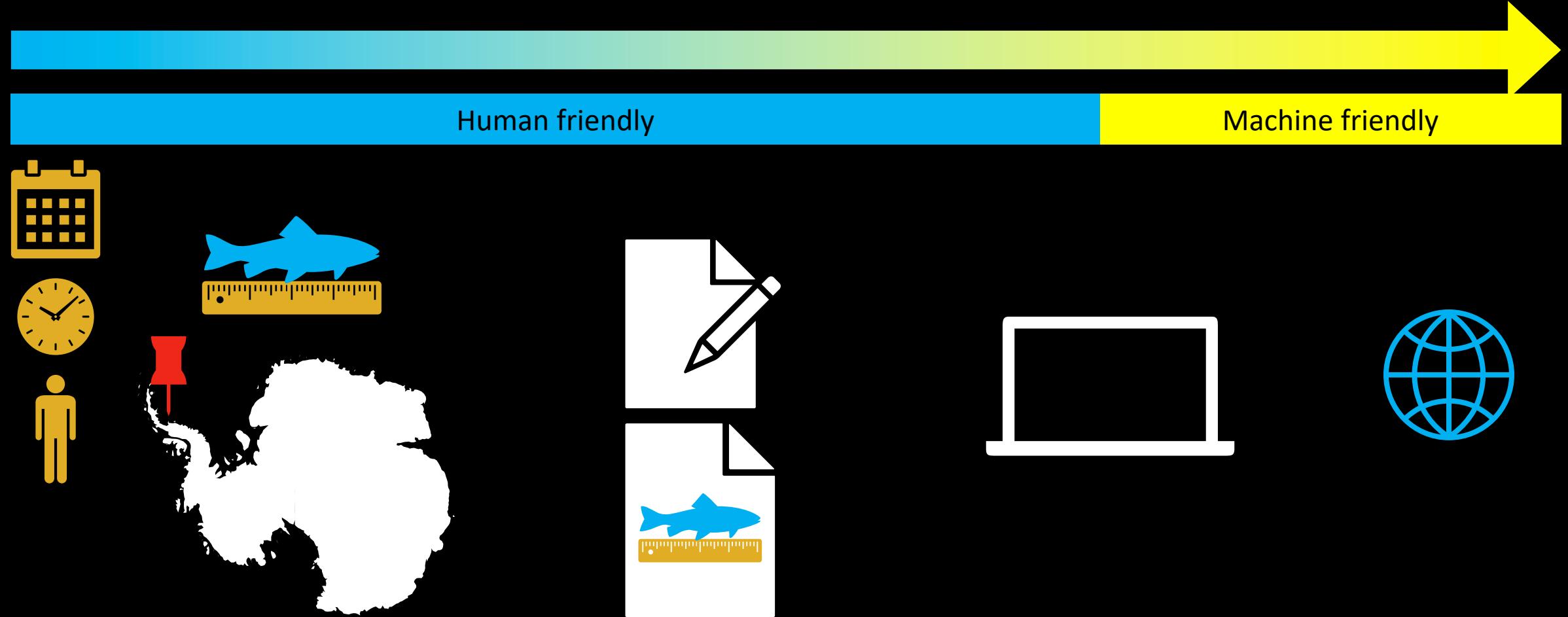
Physical observation

Physical records

Digital record

Distributed data

# Physical observation to digital data



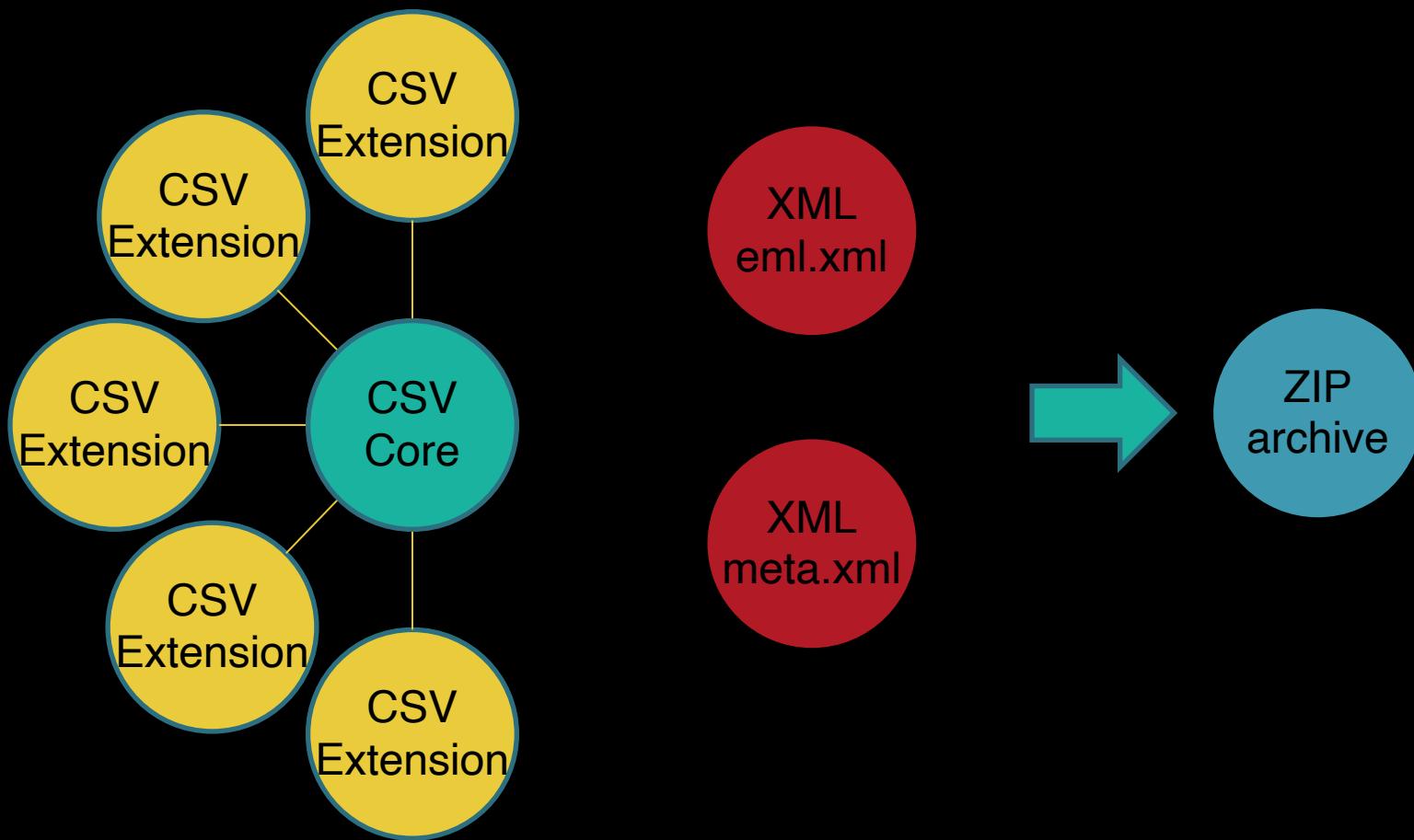
Physical observation

Physical records

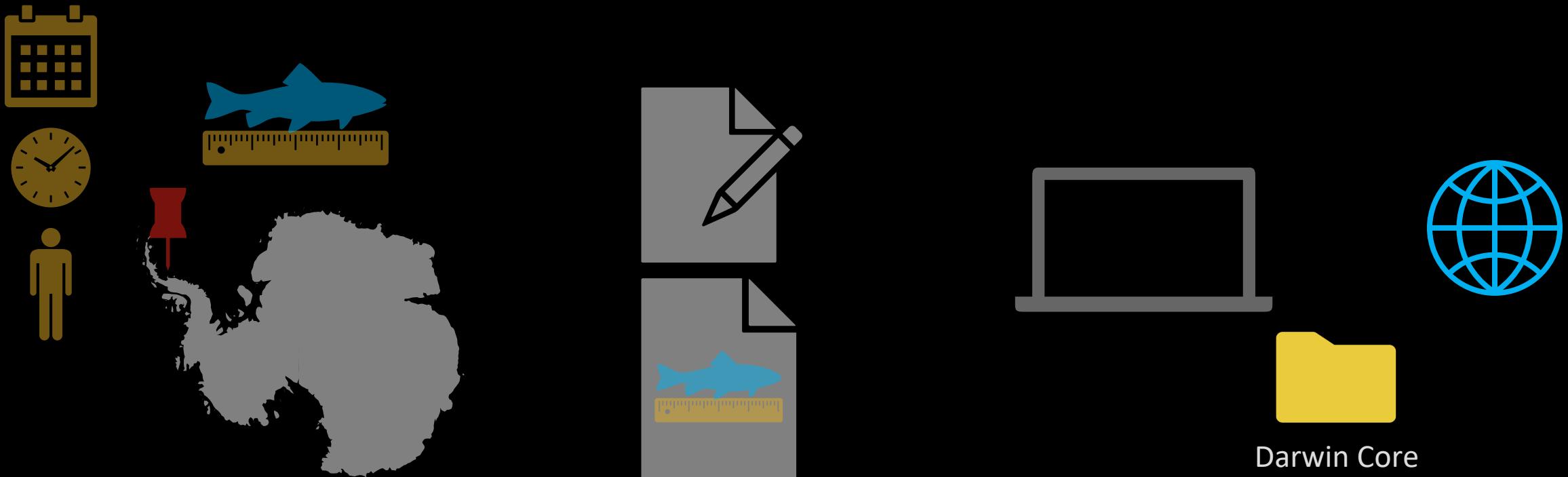
Digital record

Distributed data

# Darwin Core standard serves to distribute biodiversity data



# Darwin Core standard serves to distribute biodiversity data



Physical observation

Physical records

Digital record

Distributed data



---

## Darwin Core quick reference guide

<https://dwc.tdwg.org/terms/>

## Darwin Core extensions

<https://tools.gbif.org/dwca-validator/extensions.do>

## Darwin core validator

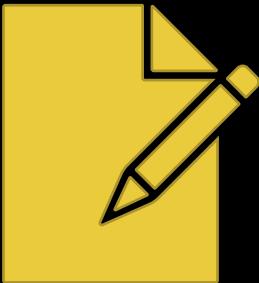
<https://www.gbif.org/tools/data-validator/about>

## Darwin core instructions

<https://github.com/gbif/ipt/wiki/howToPublish#instructions>

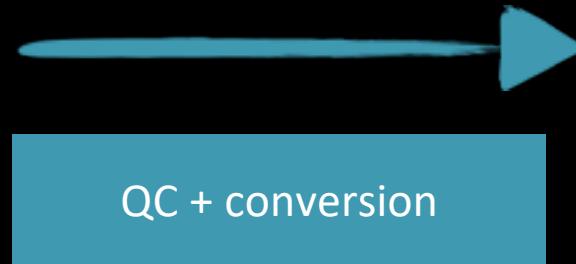
<https://obis.org/manual/ipt/>

# Objectives



field templates

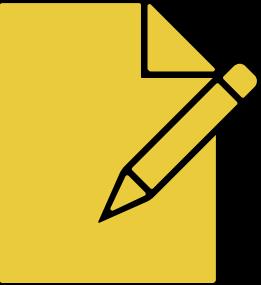
- Design to help scientists record and aggregate data
- Customized according to needs (1-1 feedback sessions)
- Human readable and writeable



Data templates

- Design for data exchange across the web
- Allows integrating data from different field expeditions
- Data standard
- Machine readable

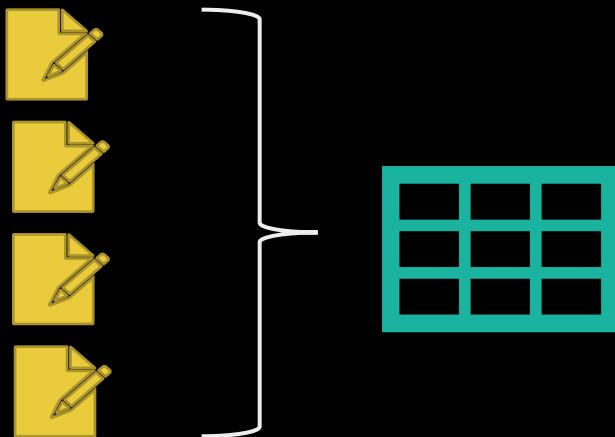
# 2 types of templates



## Field templates

- Record data in the **field**
  - Per sampling station/samples

# 2 types of templates



# Data templates

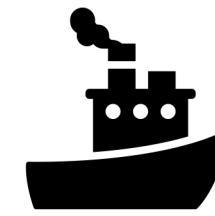
- Aggregate data from multiple field templates
  - Multiple stations/samples

# On identifiers

- Don't change them
- If you can keep them logical
- Don't use strange characters
- Did I say don't change them?
- An example!!!



Example: sampling at sea



## Original identifiers

RV105\_001

RV105\_001\_01

RV105\_001\_02

RV105\_001\_03

RV105\_002

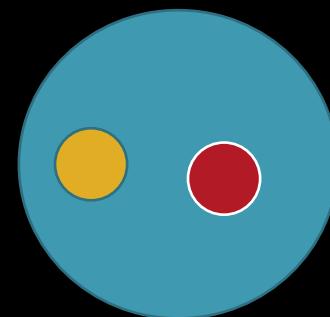
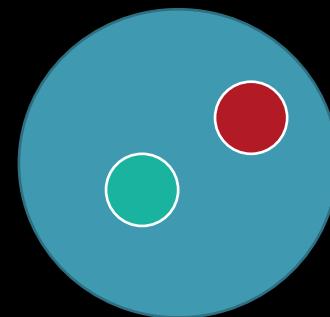
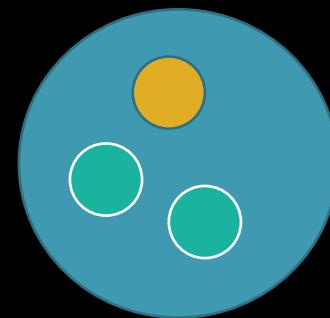
RV105\_002\_01

RV105\_002\_02

RV105\_003

RV105\_003\_01

RV105\_003\_02



## Identifiers people come up with

CTD\_01

Trawl\_01

Trawl\_02

Trawl\_03

Trawl\_01

Run1  
01\_1  
Run2  
Sample\_01\_01  
Sample\_02\_01  
01\_2

Sample\_1  
Sample\_2

Sample\_3

Trawl\_02

CTD\_02

## Original identifiers

RV105\_001

RV105\_001\_01

RV105\_001\_02

RV105\_001\_03

RV105\_002

RV105\_002\_01

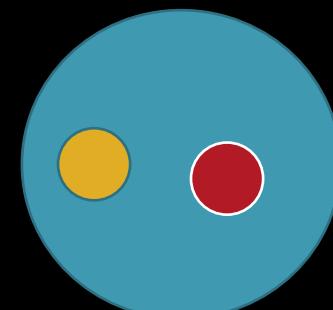
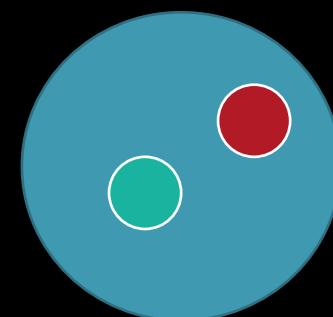
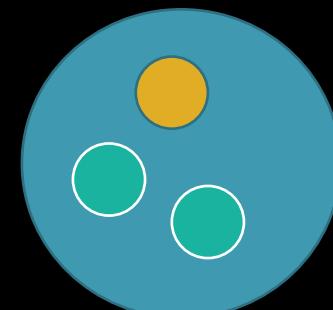
RV105\_002\_02

RV105\_003

RV105\_003\_01

RV105\_003\_02

## Identifiers people come up with



CTD\_01  
Trawl\_01  
Trawl\_02  
  
Trawl\_03  
Trawl\_01

Trawl\_02  
CTD\_02

Run1  
01\_1  
Run2  
Sample\_01\_01  
Sample\_02\_01  
01\_2

Sample\_1  
Sample\_2

Sample\_3

# Field Event template

Campaign:	Station:	Date:
	Event:	Time:
Recorded by:		Time zone:
Latitude:	Longitude:	Coordinate system:

Sampling Protocol:	
Gear:	
Sampling Effort:	

### **Station comments:**

For more information about the study, please contact Dr. [REDACTED] at [REDACTED].

## Sampling events

---

**Comments:**

For more information about the study, please contact Dr. John Smith at (555) 123-4567 or via email at [john.smith@researchinstitute.org](mailto:john.smith@researchinstitute.org).

# How are the templates organized?

The **field event** sheet for sampling events data

The diagram illustrates the organization of sampling event data. It features a central map of Antarctica with a red pushpin marking a 'Sampling station'. To the left of the map are three yellow icons: a calendar, a clock, and a person. Above the map, the text 'Sampling station' is written next to the pushpin. To the right of the map, there is a table titled 'Sampling events' with columns for event, measurement, unit, and remarks. Three lines point from the text 'Sampling event 1' and 'Sampling event 2' to the first two rows of the table. Another line points from the text 'Sub-event 1' and 'Sub-event 2' to the third and fourth rows of the table, indicating a hierarchical structure where multiple events can have sub-events.

Sampling events					
event	Measurement	Measureme	Measurement	Measuremen	Remarks
	Unit	Unit	Unit	Unit	

Campaign:		Station:		Date:	
		Event:		Time:	
Recorded by:				Time zone:	
Latitude:		Longitude:		Coordinate system:	
Sampling Protocol:					

# How are the templates organized?

The field **event** sheet for sampling events data



Station 1

RMT Net

Sampling event 2

RMT Net 1  
RMT Net 2

Sampling events		
event	start depth	end depth
	m	m
RMT Net 1	1000	500
RMT Net 2	500	200

Campaign: \_\_\_\_\_  
Recorded by: \_\_\_\_\_

**Station number:**

**Event ID:** \_\_\_\_\_  
**Date:** \_\_\_\_\_

## Summary

For more information about the study, please contact Dr. John Smith at (555) 123-4567 or via email at [john.smith@researchinstitute.org](mailto:john.smith@researchinstitute.org).

### Remarks

For more information about the study, please contact Dr. John Smith at (555) 123-4567 or via email at [john.smith@researchinstitute.org](mailto:john.smith@researchinstitute.org).

# Field occurrence template

# How are the templates organized?

The field **occurrence** sheet for sampling events data

Field **event** sheet  
**eventID: RMT Net**

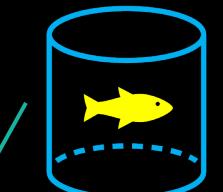
<u>Sampling events</u>		
event	start depth	end depth
	m	m
RMT Net 1	1000	500
RMT Net 2	500	200



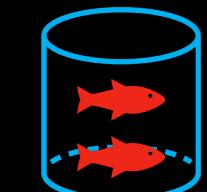
Station 1 — RMT Net



RMT Net 1  
RMT Net 2



Jar 01



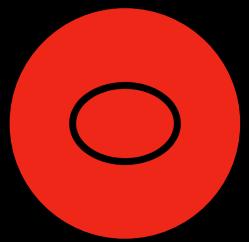
Jar 02

Field **occurrence** sheet  
for **RMT Net 1**

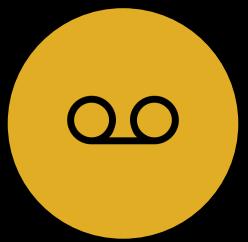
Sample ID	Species	Count
Jar 01		1
Jar 02		2



WHAT



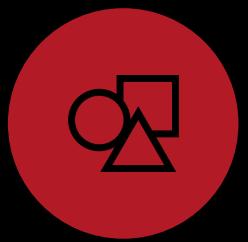
WHERE



WHEN



WHO



HOW

# Data occurrence template

# Occurrence Core: Key terms

**occurrenceID**

**basisOfRecord**

**scientificName**

**scientificNameID**

**occurrenceStatus**

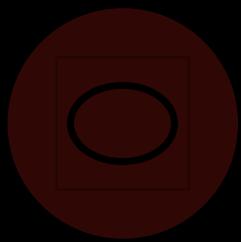
**decimalLongitude**

**decimalLatitude**

**eventDate**



WHAT



WHERE



WHEN

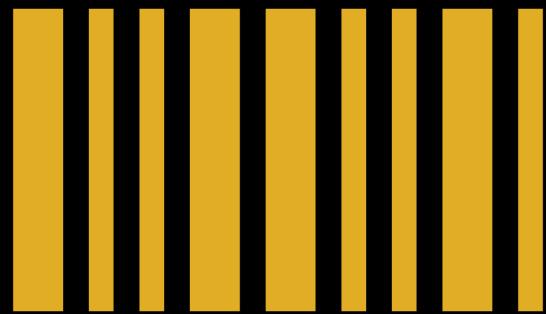


WHO



HOW

**occurrenceID**



# occurrenceID: example

occurrenceID	institutionCode	collectionCode	catalogNumber
<a href="http://arctos.database.museum/guid/MSB:Mamm:233627">http://arctos.database.museum /guid/MSB:Mamm:233627</a>	MSB	Mamm	233627
PS89_FF_000023	PS89	FF	000023

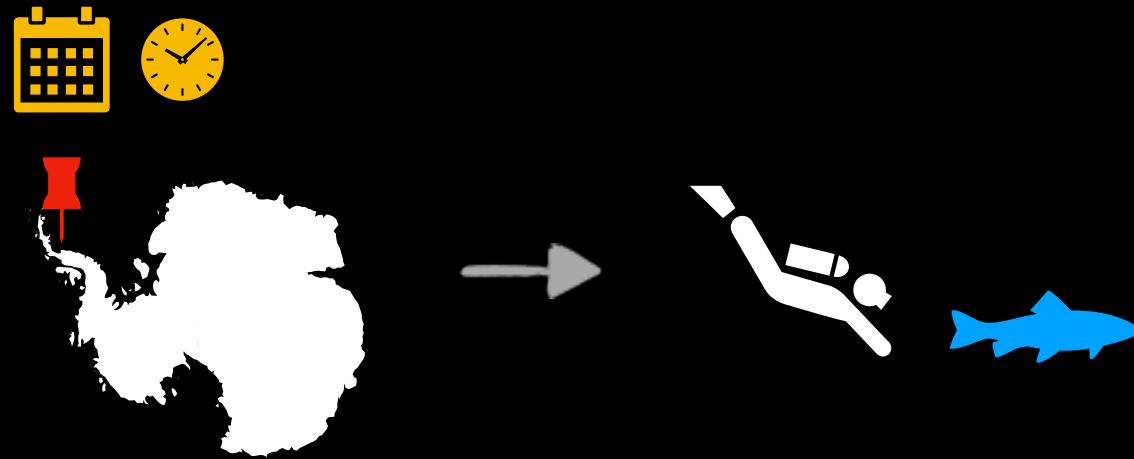
Mandatory term with standard labels

**basisOfRecord**

- HumanObservation
- MachineObservation
- PreservedSpecimen
- MaterialSample
- LivingSpecimen
- FossilSpecimen

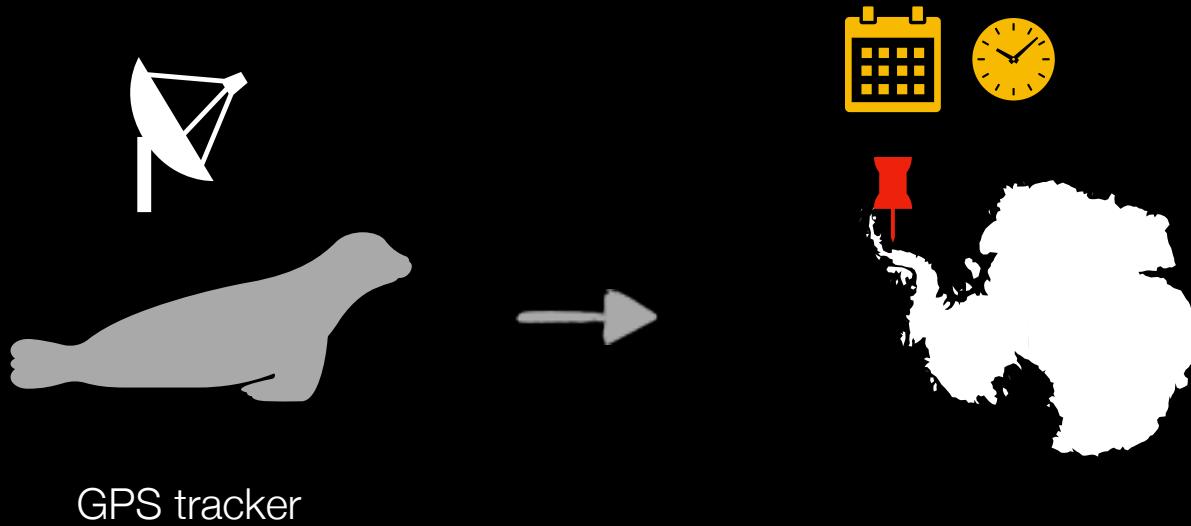
# Basis of record

HumanObservation



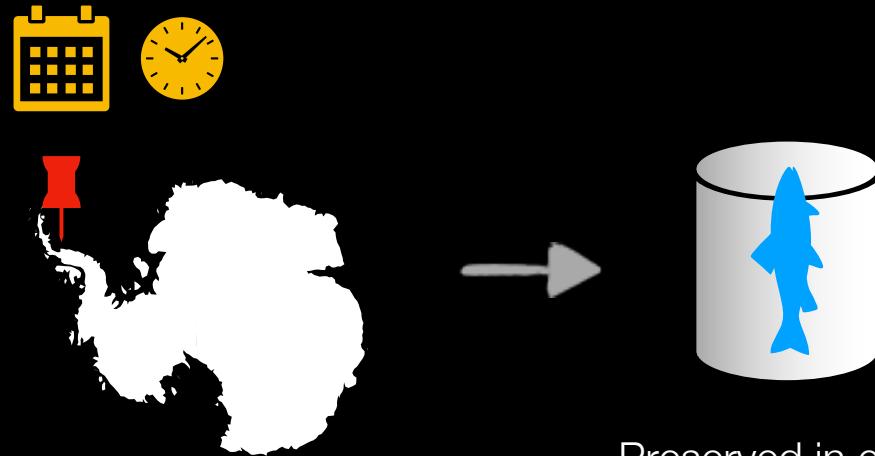
# Basis of record

MachineObservation



# Basis of record

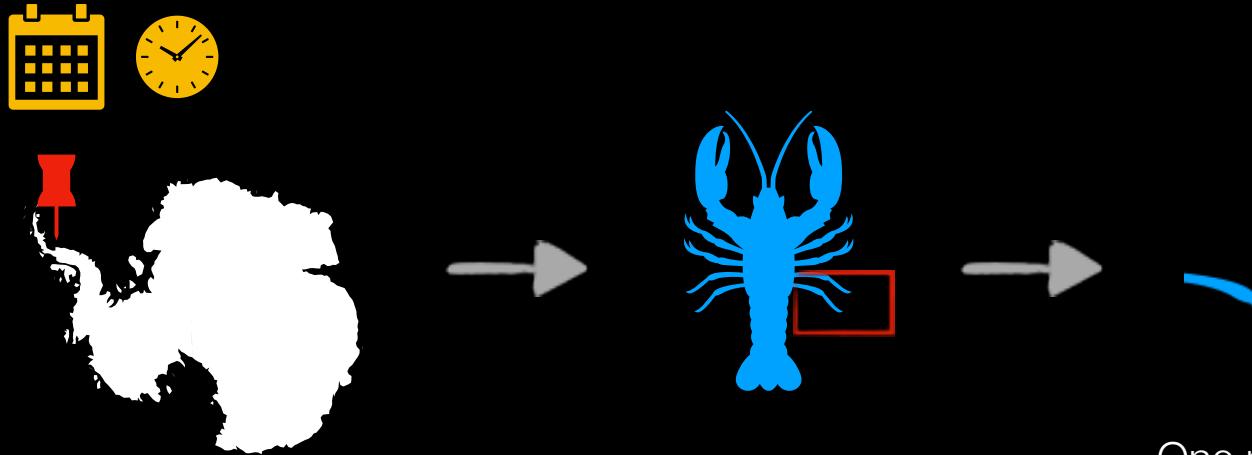
PreservedSpecimen



Preserved in ethanol

# Basis of record

MaterialSample



One pleopod

# Basis of record

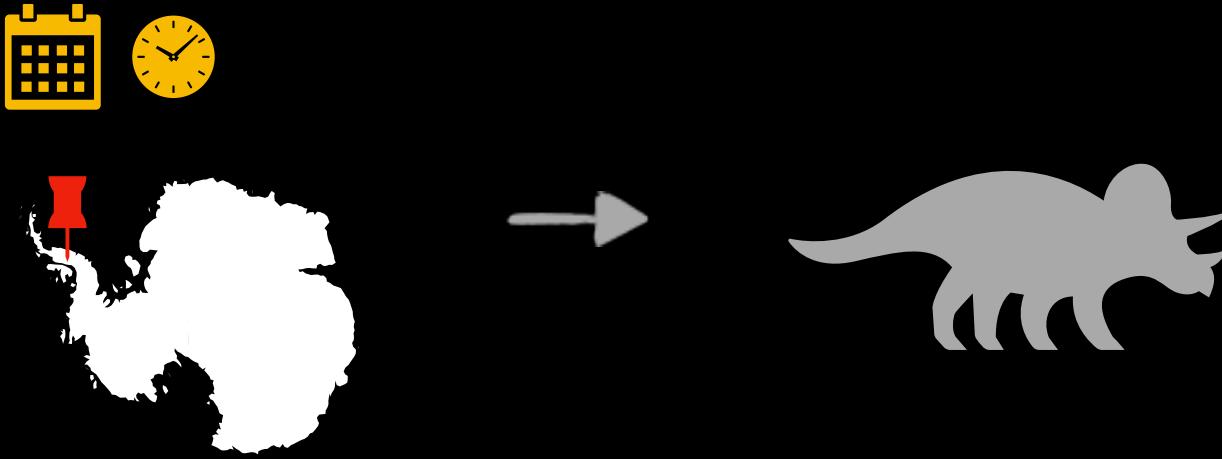
Living Specimen



A living plant in a botanical garden

# Basis of record

FossilSpecimen



# BasisOfRecord: Provide context in combination of other fields

**basisOfRecord**

- HumanObservation

***type***

- StillImage

***associatedMedia***

- Link to media (doi)

This record is an observation made by human based on still image from the associated media

# IndividuCount

- Not A required term for OBIS/GBIF
- Highly recommended by us

# Filling the occurrence data template

Field **event** sheet  
**eventID:** RMT Net

Sampling events		
event	start depth	end depth
	m	m
RMT Net 1	1000	500
RMT Net 2	500	200



Station 1



**occurrenceID**

**Species**

**individualCount**

**basisOfRecord**

ANT1\_Stn1\_RMT-Net\_1\_Jar01



1

PreservedSpecimen

ANT1\_Stn1\_RMT-Net\_1\_Jar02

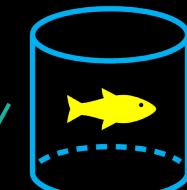


2

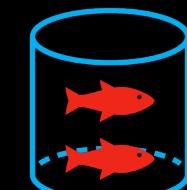
PreservedSpecimen

Field **occurrence** sheet  
**eventID:** RMT Net 1

Sample ID	Species	Count
Jar 01		1
Jar 02		2



Jar 01



Jar 02

**scientificName**

**scientificNameID**

## WoRMS

- Register of Antarctic  
(marine) Species

## Global Names Index (GNI)

## Integrated Taxonomic Information System (ITIS)

scientificName, scientificNameID

# Taxonomy

Uncertain identifications

Reduce taxonomy to **first common higher classification level**

*cf.*

- *confer*
- resembles, looks like

*aff.*

- *affinis*
- related to but not identical to

Species | *Electrona antarctica*

Species | *Electrona aff. antarctica*

Genus

*Electrona*

# Taxonomy

Uncertain identifications

Reduce taxonomy to **first common higher classification level**



scientificName	genus	identificationQualifier	taxonRank
Acanthonchus	Acanthonchus	cf. duplicatus	genus

# Examples

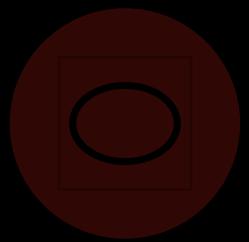
scientificName	scientificName Authorship	Identification Qualifier	taxonrank	scientificNameID
<u>Electrona</u>	Goode & Bean, 1896	aff. risso	genus	urn:lsid:marinespecies.org:taxname:125821)
<u>Electrona antarctica</u>	Günther, 1878		species	(urn:lsid:marinespecies.org:taxname:217697)
<u>Electrona</u>	Goode & Bean, 1896	sp.	genus	urn:lsid:marinespecies.org:taxname:125821)
<u>Electrona</u>	Goode & Bean, 1896	spp.	genus	urn:lsid:marinespecies.org:taxname:125821)

*identifiedBy*  
*dateIdentified*

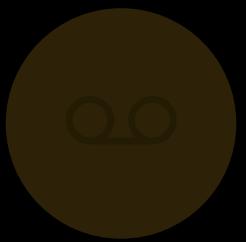
Acknowledge the personnel who carry out different tasks:  
e.g. Taxonomist who identified the species  
(<https://bionomia.net/>)



WHAT



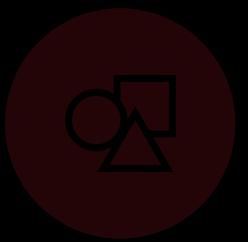
WHERE



WHEN



WHO



HOW

# Filling the occurrence data template

Field **event** sheet  
**eventID:** RMT Net

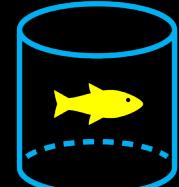
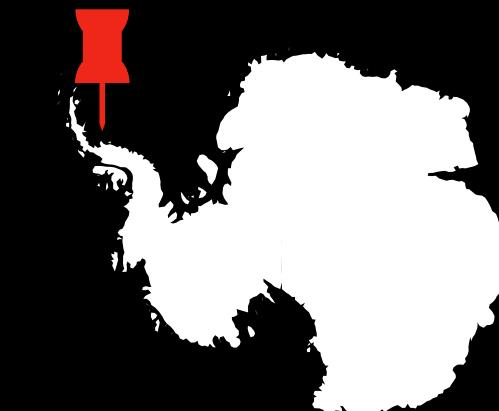
Sampling events		
event	start depth	end depth
	m	m
RMT Net 1	1000	500
RMT Net 2	500	200



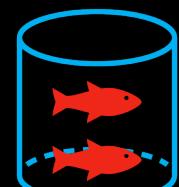
Station 1



— RMT Net



Jar 01



Jar 02

Field **occurrence** sheet  
**eventID:** RMT Net 1

Sample ID	Species	Count
Jar 01		1
Jar 02		2

occurrenceID	ScientificName	individualCount	basisOfRecord	scientificNameID	identifiedBy	IdentifiedDate
ANT1_Stn1_RMT-Net_1_Jar01		1	PreservedSpecimen	125821	Anton Van de Putte	2008-11-05
ANT1_Stn1_RMT-Net_1_Jar02		2	PreservedSpecimen	217697	Anton Van de Putte	2008-11-05

**occurrenceStatus**

Present/absent

# Filling the occurrence data template

Field **event** sheet  
**eventID:** RMT Net

Sampling events		
event	start depth	end depth
	m	m
RMT Net 1	1000	500
RMT Net 2	500	200



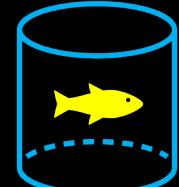
Station 1



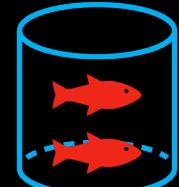
RMT Net



RMT Net 1  
RMT Net 2



Jar 01



Jar 02

Field **occurrence** sheet  
**eventID:** RMT Net 1

Sample ID	Species	Count
Jar 01		1
Jar 02		2

occurrenceID	ScientificName	individualCount	basisOfRecord	scientificNameID	Occurrence status
ANT1_Stn1_RMT-Net_1_Jar01		1	PreservedSpecimen	125821	present
ANT1_Stn1_RMT-Net_1_Jar01		0	HumanObservation	217697	absent
ANT1_Stn1_RMT-Net_1_Jar02		2	PreservedSpecimen	217697	present
ANT1_Stn1_RMT-Net_1_Jar02		0	HumanObservation	125821	absent

## Recommended information

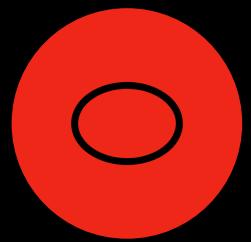


WHAT

- *organismQuantity*
- *organismQuantityType*
- *sex*
- *lifeStage*
- *behavior*
- *occurrenceRemarks*



WHAT



WHERE



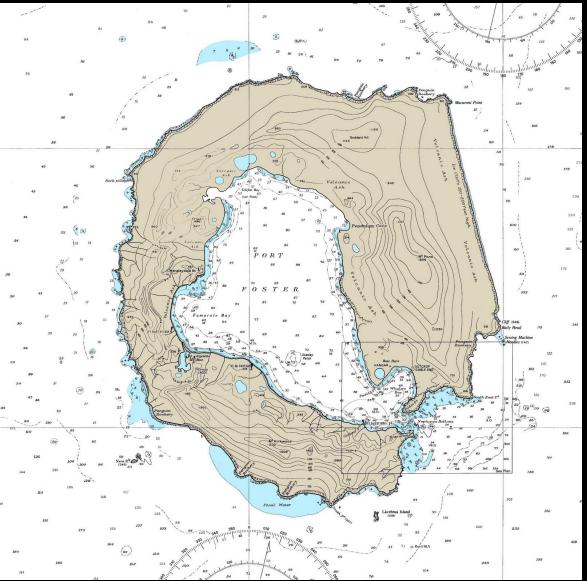
WHEN



WHO



HOW



# Geographic coordinates

**decimalLongitude**

**decimalLatitude**

# Geographic coordinates

## Coordinates precision

### WHAT THE NUMBER OF DIGITS IN YOUR COORDINATES MEANS

#### LAT/LON PRECISION

#### MEANING

28°N, 80°W	YOU'RE PROBABLY DOING SOMETHING SPACE-RELATED
28.5°N, 80.6°W	YOU'RE POINTING OUT A SPECIFIC CITY
28.52°N, 80.68°W	YOU'RE POINTING OUT A NEIGHBORHOOD
28.523°N, 80.683°W	YOU'RE POINTING OUT A SPECIFIC SUBURBAN CUL-DE-SAC
28.5234°N, 80.6830°W	YOU'RE POINTING TO A PARTICULAR CORNER OF A HOUSE
28.52345°N, 80.68309°W	YOU'RE POINTING TO A SPECIFIC PERSON IN A ROOM, BUT SINCE YOU DIDN'T INCLUDE DATUM INFORMATION, WE CAN'T TELL WHO
28.5234571°N, 80.6830941°W	YOU'RE POINTING TO WALDO ON A PAGE
28.523457182°N, 80.683094159°W	"HEY, CHECK OUT THIS SPECIFIC SAND GRAIN!"
28.523457182818284°N, 80.683094159265358°W	EITHER YOU'RE HANDING OUT RAW FLOATING POINT VARIABLES, OR YOU'VE BUILT A DATABASE TO TRACK INDIVIDUAL ATOMS. IN EITHER CASE, PLEASE STOP.

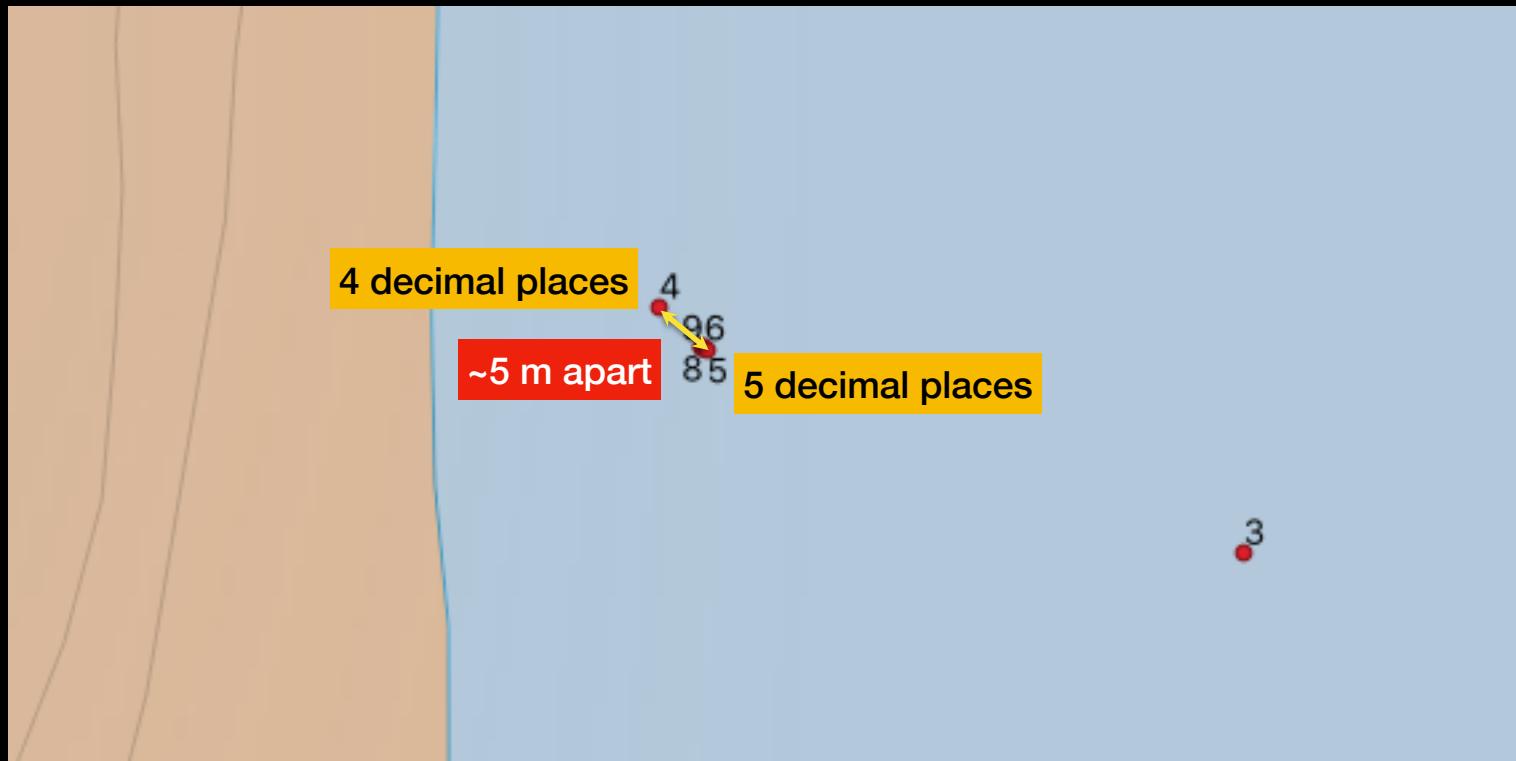
# Geographic coordinates

## Coordinates precision

Coordinates with many (>5) decimal places maybe resulted from data conversion  
(e.g. converting degrees, minutes, seconds to decimal degrees)

# Geographic coordinates

## Coordinates precision



# Geographic coordinates

## Coordinates precision



# Geographic coordinates

## Coordinates precision



# Geographic coordinates

## Coordinates precision



# Geographic coordinates

Coordinates precision

What about sensitive data?

# Geographic coordinates

Coordinates precision

Which precision is appropriate?

# Geographic coordinates

Coordinates precision

How sensitive is your species?

To what extend your data needs to be generalized

# Geographic coordinates

## Coordinates precision

Chapman AD (2020) Current Best Practices for Generalizing Sensitive Species Occurrence Data (Community review draft). Version 2. Copenhagen: GBIF Secretariat.  
<https://doi.org/10.15468/doc-5jp4-5g10>.

# Geographic coordinates

Longitude, latitude switched

# Geographic coordinates

Longitude, latitude switched

Often due to human errors when recording  
coordinates

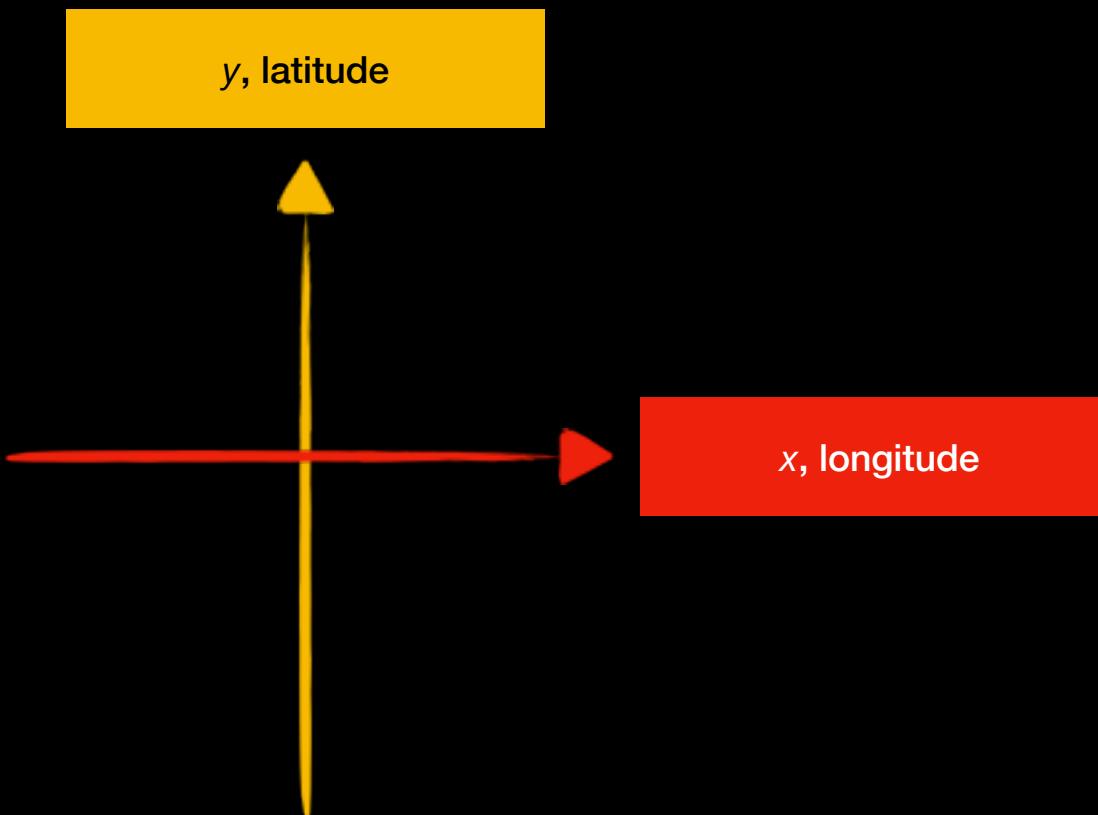
## verbatimCoordinates

58°28'30"S, 62°7'0,10"W
58°28'42,00"S, 62°7'15,00"W
62°08.0958"S, 058°24,1625"W
62°10.4680"S, 058°25.2137"W
62°09.7174"S, 058°21,5886"W
62°11.2075"S, 058°18,9951W
62°11.9764"S, 058°22.5800W
62°12.1942"S, 058°23.4483"W
62°53.6407"S, 58°26.8232"W

# Geographic coordinates

Longitude, latitude switched

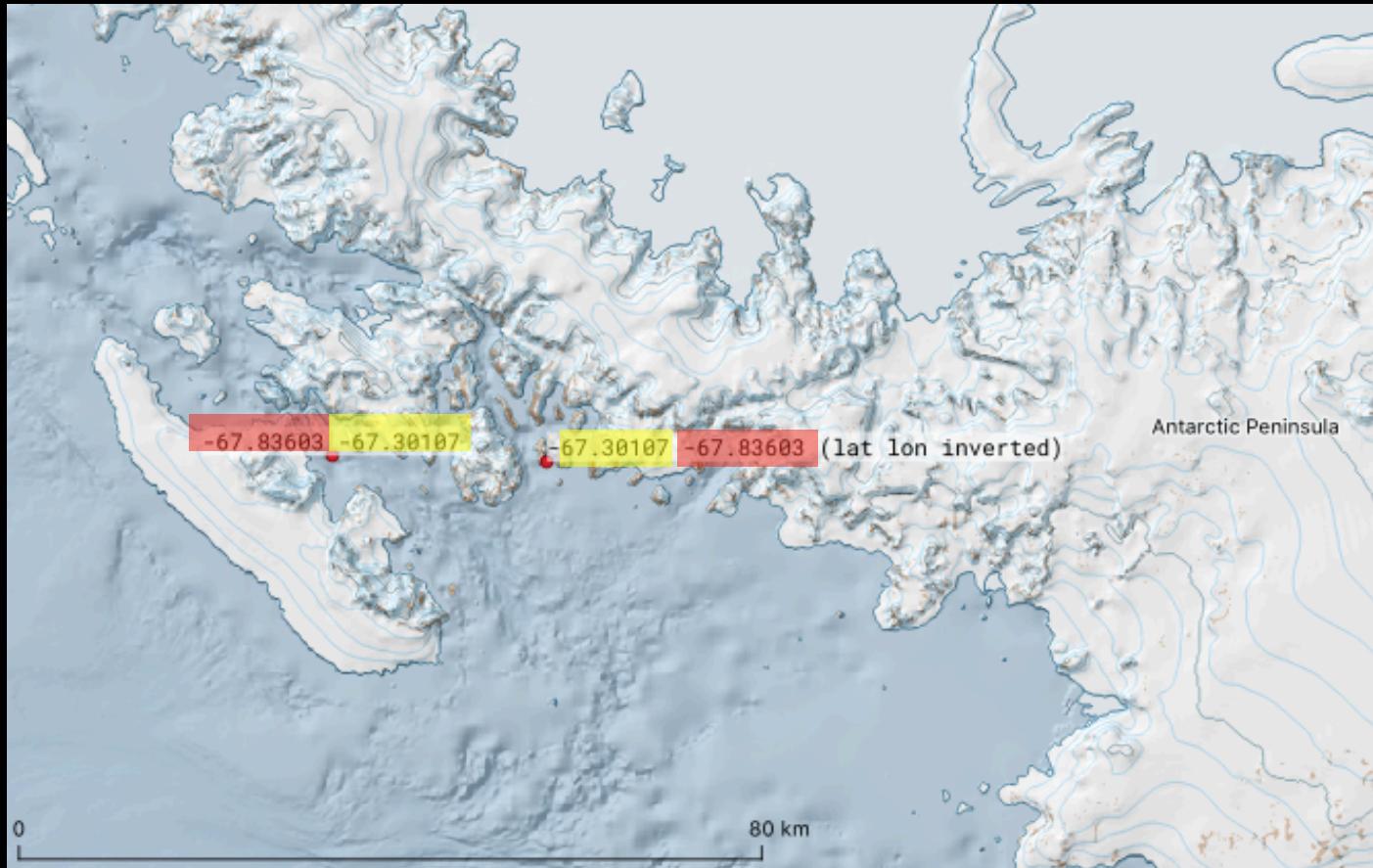
Sometimes due to confusion



# Geographic coordinates

Consequence of switched longitude and latitude

Longitude  
Latitude



Marine species - Both points are in water, which is correct?

# Atomize columns

Best to record latitude and longitude in **2 separate columns**

verbatimLongitude	verbatimLatitude
-67° 50' 9.71"	-67° 18' 3.85"

# Atomize columns

To ease downstream batch operations

verbatimLongitude	verbatimLatitude
-67° 50' 9.71"	-67° 18' 3.85"



Function can be applied to  
whole column

decimalLongitude	decimalLatitude
-67.83603	-67.30107

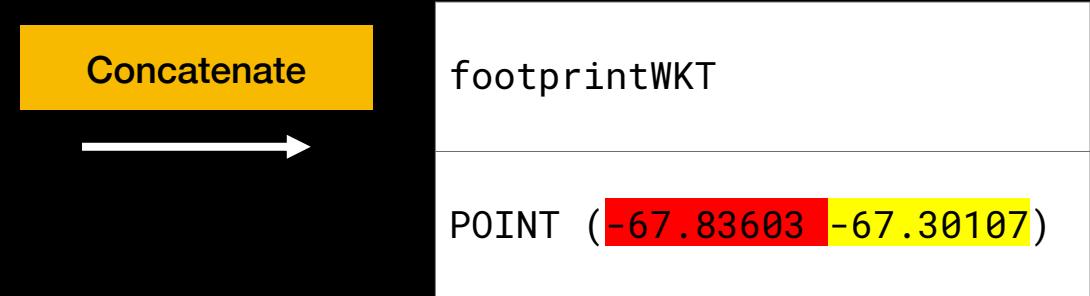
# Atomize columns

Function can be applied to whole column

verbatimLongitude	verbatimLatitude
-67° 50' 9.71"	-67° 18' 3.85"



decimalLongitude	decimalLatitude
-67.83603	-67.30107



# Geographic coordinates

Consistent format

# Geographic coordinates

Be consistent

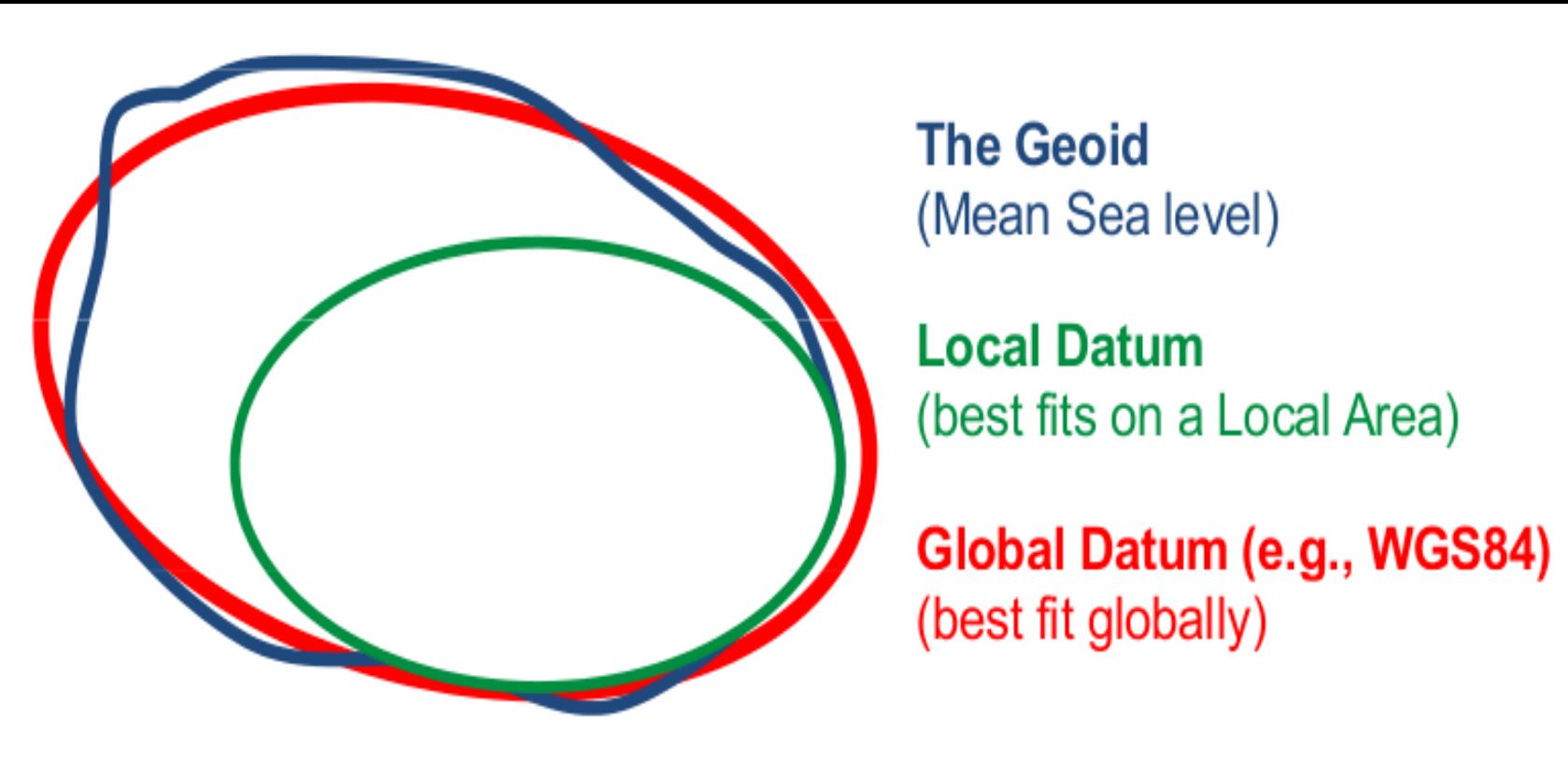
- Type of decimal point
- Order of longitude, latitude
- Degree minute second vs degree decimal minute
- Padding with 0

verbatimCoordinates
58°28'30"S, 62°7'0,10"S
58°28'42,00"W, 62°7'15,00"S
62°08.0958'S, 058°24,1625'W
62°10.4680'S, 058°25.2137'W
62°09.7174'S, 058°21,5886'W
62°11.2075'S, 058°18,9951W
62°11.9764'S, 058°22.5800W
62°12.1942'S, 058°23.4483'W
62°53.6407'S, 58°26.8232'W

# Geodetic datum

# Geodetic datum

Most of the datasets we received lack this information



# Geodetic datum

decimalLongitude	decimalLatitude	geodeticDatum
-67.83603	-67.30107	EPSG:4326

Use EPSG code of the  
spatial reference system

EPSG code for World Geodetic System 1984 (WGS 84) is EPSG:4326

# Georeferencing best practice

Georeferencing Best Practice from GBIF

<http://mb.gbif.org/documents/doc-georeferencing-best-practices/en/>

# Filling in the occurrence data template

Field **event sheet**  
**eventID:** RMT Net

Sampling events		
event	start depth	end depth
	m	m
RMT Net 1	1000	500
RMT Net 2	500	200

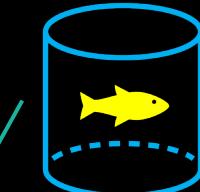


Station 1

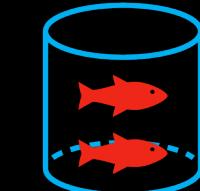


RMT Net

RMT Net 1  
RMT Net 2



Jar 01



Jar 02

Field **occurrence sheet**  
**eventID:** RMT Net 1

Sample ID	Species	Count
Jar 01		1
Jar 02		2

occurrenceID	scientificName	individualCount	basisOfRecord	decimalLongitude	decimalLatitude	geographicDatum
ANT1_Stn1_RMT-Net_1_Jar01		1	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326
ANT1_Stn1_RMT-Net_1_Jar02		2	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326

eventDate

## PUBLIC SERVICE ANNOUNCEMENT:

OUR DIFFERENT WAYS OF WRITING DATES AS NUMBERS CAN LEAD TO ONLINE CONFUSION. THAT'S WHY IN 1988 ISO SET A GLOBAL STANDARD NUMERIC DATE FORMAT.

THIS IS **THE** CORRECT WAY TO WRITE NUMERIC DATES:

2013-02-27

THE FOLLOWING FORMATS ARE THEREFORE DISCOURAGED:

02/27/2013 02/27/13 27/02/2013 27/02/13

20130227 2013.02.27 27.02.13 27-02-13

27.2.13 2013. II. 27. 2 $\frac{1}{2}$ -13 2013.158904109

MMXIII-II-XXVII MMXIII  $\frac{LVII}{CCCLXV}$  1330300800

$((3+3)\times(111+1)-1)\times3/3-1/3^3$  2013 Mississipi

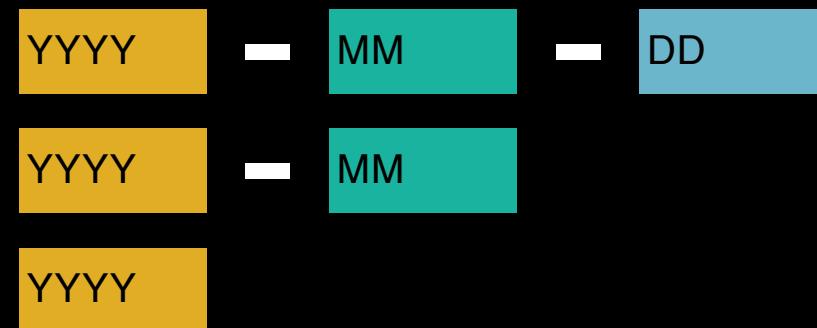
10/11011/1101 02/27/20/13  $\begin{matrix} 2 & 3 & 1 & 4 \\ \hline 5 & 6 & 7 & 8 \end{matrix}$



# Date

ISO 8601 standard

Fixed number of digits padded with leading zero



# Filling in the occurrence data template

Field **event sheet**  
**eventID:** RMT Net

Sampling events		
event	start depth	end depth
	m	m
RMT Net 1	1000	500
RMT Net 2	500	200



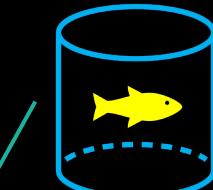
Station 1



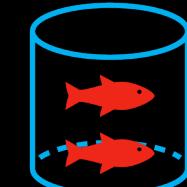
RMT Net



RMT Net 1  
RMT Net 2



Jar 01



Jar 02

Field **occurrence sheet**  
**eventID:** RMT Net 1

Sample ID	Species	Count
Jar 01		1
Jar 02		2

occurrenceID	scientificName	individualCount	basisOfRecord	decimalLongitude	decimalLatitude	geodeticDatum	eventDate	year	month	day
ANT1_Stn1_RMT-Net_1_Jar01		1	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326	2020-11-02	2020	11	02
ANT1_Stn1_RMT-Net_1_Jar02		2	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326	2020-11-02	2020	11	02

# Excel saves what you see, not what you typed

A screenshot of Microsoft Excel version 16.42 (2020) showing a blank worksheet titled "Book1". The ribbon menu is visible at the top, showing tabs for Home, Insert, Draw, Page Layout, Formulas, Data, Review, View, Tell me, Share, and Comments. The Home tab is selected. The formula bar shows "D15" and includes buttons for Paste, Share, and Find & Select. The main area shows a grid from A1 to N11.

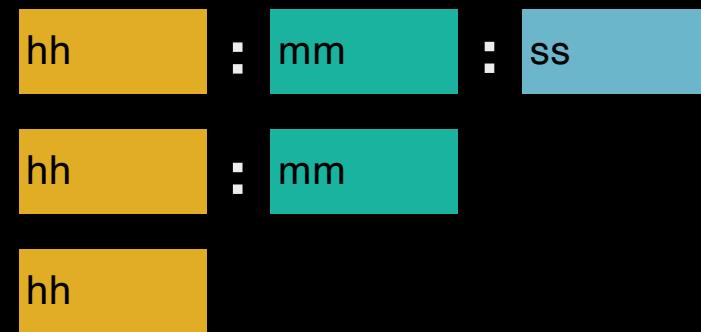
	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1														
2														
3														
4														
5														
6														
7														
8														
9														
10														
11														

Microsoft Excel version 16.42 (2020)

# Time

ISO 8601 standard

Fixed number of digits padded with leading zero



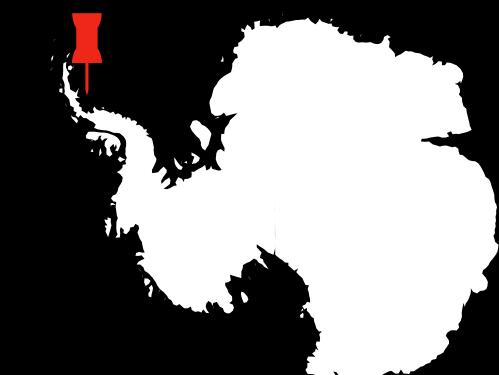
# Filling in the data template

Field **event sheet**  
**eventID:** RMT Net

Sampling events		
event	start depth	end depth
	m	m
RMT Net 1	1000	500
RMT Net 2	500	200

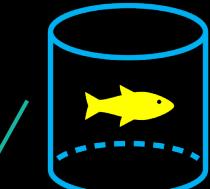


Station 1

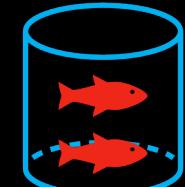


RMT Net

RMT Net 1  
RMT Net 2



Jar 01



Jar 02

Field **occurrence sheet**  
**eventID:** RMT Net 1

Sample ID	Species	Count
Jar 01		1
Jar 02		2

occurrenceID	scientificName	individualCount	basisOfRecord	decimalLongitude	decimalLatitude	geographicDatum	eventDate	year	month	day	eventTime
ANT1_Stn1_RMT-Net_1_Jar01		1	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326	2020-11-02	2020	11	02	08:32:00
ANT1_Stn1_RMT-Net_1_Jar02		2	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326	2020-11-02	2020	11	02	08:32:00

# Filling in the occurrence data template

Field **event sheet**  
**eventID:** RMT Net

Sampling events		
event	start depth	end depth
	m	m
RMT Net 1	1000	500
RMT Net 2	500	200



Station 1

— RMT Net

RMT Net 1  
RMT Net 2

Field **occurrence sheet**  
**eventID:** RMT Net 1

Sample ID	Species	Count
Jar 01		1
Jar 02		2

occurrenceID	scientificName	individualCount	basisOfRecord	decimalLongitude	decimalLatitude	geodeticDatum	eventDate	year	month	day	eventTime
ANT1_Stn1_RM T-Net_1_Jar01		1	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326	2020-11-02	2020	11	02	08:32:00
ANT1_Stn1_RM T-Net_1_Jar02		2	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326	2020-11-02	2020	11	02	08:32:00

# Time zone

Time zone is very often neglected in the datasets we received

# Time zone

ISO 8601 standard

z

UTC

$\pm$ hh:mm

Offset from UTC

# Time zone

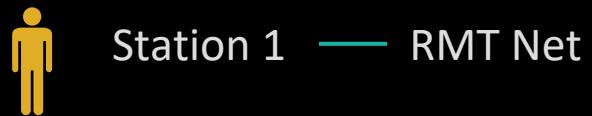
ISO 8601 standard

Example	Description
14:07-06:00	2:07 pm in the time zone 6 hours earlier than UTC
08:40:21Z	8:40:21am UTC

# Filling in the occurrence data template

Field event sheet  
eventID: RMT Net

Sampling events		
event	start depth	end depth
	m	m
RMT Net 1	1000	500
RMT Net 2	500	200



Field occurrence sheet  
eventID: RMT Net 1

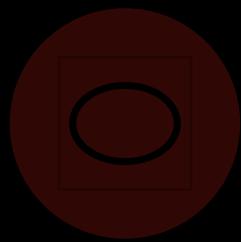
Sample ID	Species	Count
Jar 01		1
Jar 02		2

occurrenceID	scientificName	individualCount	basisOfRecord	decimalLongitude	decimalLatitude	geodeticDatum	eventDate	year	month	day	eventTime
ANT1_Stn1_RMT-Net_1_Jar01		1	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326	2020-11-02	2020	11	02	08:32:00
ANT1_Stn1_RMT-Net_1_Jar02		2	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326	2020-11-02	2020	11	02	08:32:00

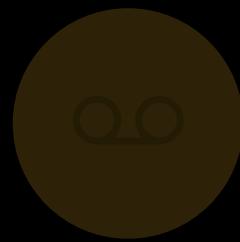
occurrenceID	timeZone
ANT1_Stn1_RMT-Net_1_Jar01	Z
ANT1_Stn1_RMT-Net_1_Jar02	Z



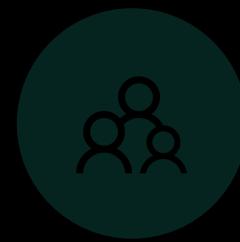
WHAT



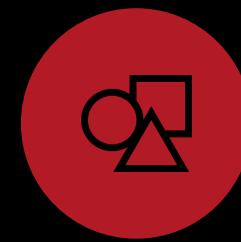
WHERE



WHEN



WHO



HOW

# Filling in the occurrence data template

Field event sheet  
eventID: RMT Net

Sampling events		
event	start depth	end depth
	m	m
RMT Net 1	1000	500
RMT Net 2	500	200



Field occurrence sheet  
eventID: RMT Net 1

Sample ID	Species	Count
Jar 01		1
Jar 02		2

occurrenceID	scientificName	individualCount	basisOfRecord	decimalLongitude	decimalLatitude	geographicDatum	eventDate	year	month	day	eventTime
ANT1_Stn1_RMT-Net_1_Jar01		1	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326	2020-11-02	2020	11	02	08:32:00
ANT1_Stn1_RMT-Net_1_Jar02		2	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326	2020-11-02	2020	11	02	08:32:00

occurrenceID	timeZone	identifiedBy	eventID	samplingProtocol	minimumDepthInMeters	maximumDepthInMeters
ANT1_Stn1_RMT-Net_1_Jar01	Z	John Doe	ANT1_Stn1_RMT-Net_1	RMT Net	500	1000
ANT1_Stn1_RMT-Net_1_Jar02	Z	John Doe	ANT1_Stn1_RMT-Net_1	RMT Net	500	1000

# Filling in the occurrence data template

Field event sheet  
eventID: RMT Net

Sampling events		
event	start depth	end depth
	m	m
RMT Net 1	1000	500
RMT Net 2	500	200



Station 1

— RMT Net

RMT Net 1  
RMT Net 2

Field occurrence sheet  
eventID: RMT Net 1

Sample ID	Species	Count
Jar 01		1
Jar 02		2

occurrenceID	scientificName	individualCount	basisOfRecord	decimalLongitude	decimalLatitude	geodeticDatum	eventDate	year	month	day	eventTime
ANT1_Stn1_RMT-Net_1_Jar01		1	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326	2020-11-02	2020	11	02	08:32:00
ANT1_Stn1_RMT-Net_1_Jar02		2	PreservedSpecimen	-67.83603	-67.30107	EPSG:4326	2020-11-02	2020	11	02	08:32:00

occurrenceID	timeZone	identifiedBy	eventID	minimumDepthInMeters	maximumDepthInMeters	measurement
						unit
ANT1_Stn1_RMT-Net_1_Jar01	Z	John Doe	ANT1_Stn1_RMT-Net_1	500	1000	Value 1
ANT1_Stn1_RMT-Net_1_Jar02	Z	John Doe	ANT1_Stn1_RMT-Net_1	500	1000	Value 2

# Epipelagic mesozooplankton distribution and abundance in Southern Ocean Atlantic sector and the North Atlantic and Arctic 1996-2013

Published by British Antarctic Survey

Peter Ward • Geraint Tarling • Rachael Shreeve • Petra ten Hoopen

Event data

samplingProtocol	samplingEffort	sampleSizeValue	sampleSizeUnit
<a href="https://www.bodc.ac.uk/resources/inventories/cruise_inventory/report/5386/">https://www.bodc.ac.uk/resources/inventories/cruise_inventory/report/5386/</a>	average of 30 min net deployment	1	square metre

The screenshot shows the BODC website interface. At the top, there's a logo for the National Oceanography Centre and the text "British Oceanographic Data Centre BODC". Below the header, there's a navigation bar with links for HOME, SEARCH THE DATA, SUBMIT YOUR DATA, PROJECTS, RESOURCES, ABOUT, and a search icon. The main content area has a dark background with a blue gradient overlay. A large, bold heading "Resources" is visible. Below it, a breadcrumb navigation shows "Resources > Inventories > Cruise inventory > Report > 5996". A prominent title "RRS James Clark Ross JR20001217 (JR57)" is displayed in large white font. Underneath, there's a section titled "Cruise summary report" with a link "Search the Cruise Inventory". At the bottom, there's a "Cruise Info." section containing details about the ship name (RRS James Clark Ross), cruise identifier (JR20001217), cruise period (2000-12-17 – 2001-01-12), and status (Completed).

**Methodology**

**Study extent**

Data were gathered on a series of oceanographic cruises aboard the RRS James Clark Ross during expeditions to the Atlantic Sector of the Southern Ocean and the North Atlantic/Arctic.

**Sampling**

Samples were retrieved by either a motion compensation Bongo net or a mini Bongo net, and then preserved with Borax buffered 10% formalin for analysis back at the home laboratory. Taxa were identified through examination by light microscopy. The abundance of each taxa within a sample was determined through examining a known fraction of the sample and then making an inverse multiplication of that fraction. Known fractions were principally achieved through the use of a Folsom splitter. Volumetric abundances (individuals m<sup>-2</sup>) were determined through dividing sample abundance by the volume of water sampled by the respective net. Net volume was mainly derived through multiplying the net opening diameter by the maximum depth of sampling. This calculation assumes 100% net sampling efficiency and that specimens were only captured during the ascent phase of the net deployment.

**Method steps**

All species were identified according to the taxonomic guides available at time of analysis and the user must be aware that some species names may have since been updated.

Project metadata

<https://www.gbif.org/dataset/402c81b6-8d69-4650-8826-92f91d2728a7>

# Summary

- Use original identifiers
- Basis of record
- Record the most certain taxonomy
- Occurrence status: present/absent
- Atomize columns for latitude and longitude – record what you see, transform later
- Take precision into account
- Geodetic datum
- Atomize columns for date
- Time and time zone
- Sampling protocol and environmental measurement

# Citations and resources

- Cartoons by XKCD more at <https://m.xkcd.com/>
- [www.biodiversity.aq](http://www.biodiversity.aq)
- [www.obis.org](http://www.obis.org)
- [www.gbif.org](http://www.gbif.org)
- <https://scar.github.io/EGABIcourse19/metadata-darwin-core.html>
- <https://www.scar.org/science/egabi/>
- Georeferencing Best Practice from GBIF  
<http://mb.gbif.org/documents/doc-georeferencing-best-practices/en/>

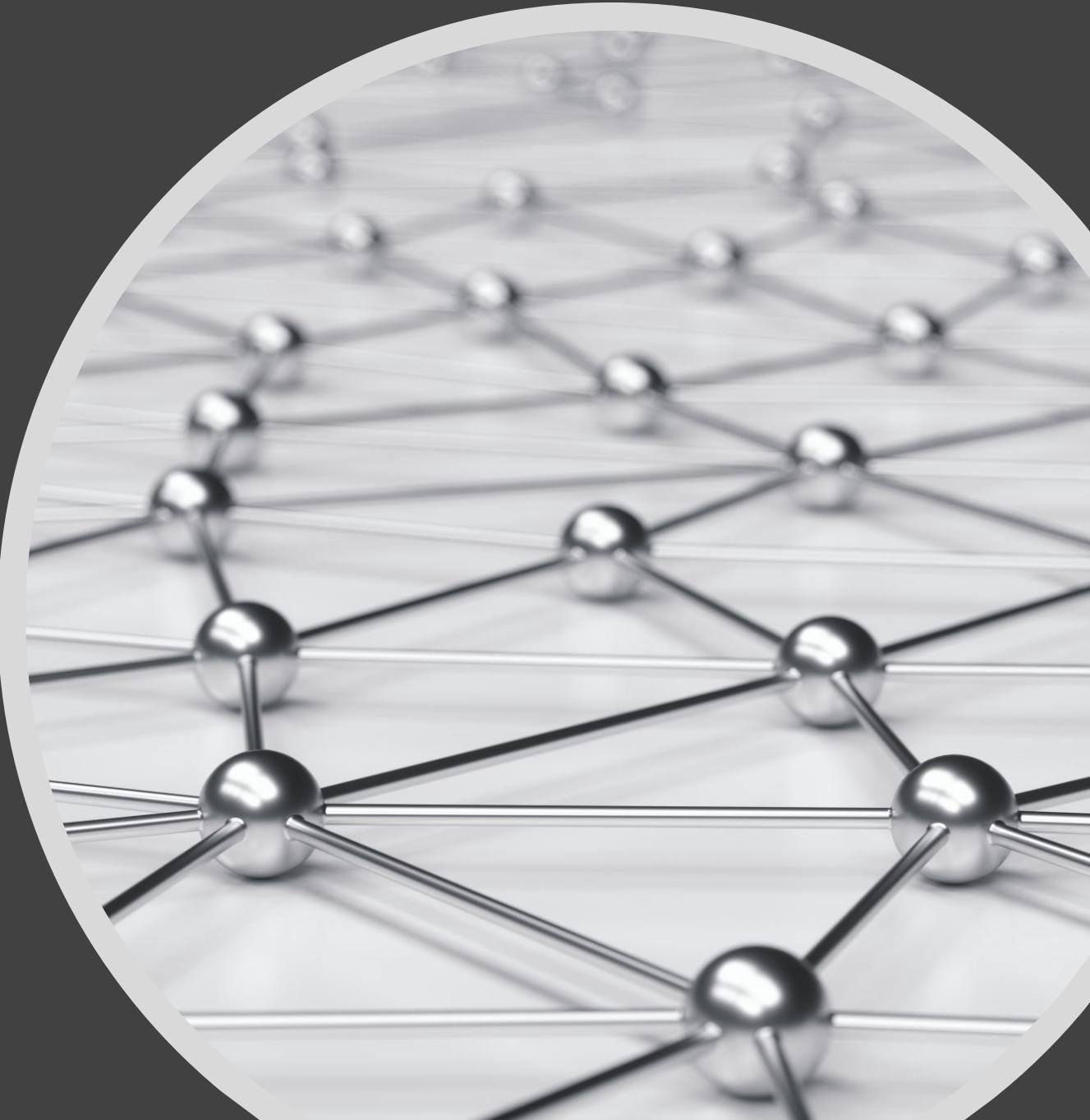
# Acknowledgement

- Scientists who contributed to templates design:  
Huw Griffiths, Anton Van de Putte, Fokje Schaafsma, Hauk Flores
- Early career scientists who provided valuable input:  
Robyn Samuel, Raissa Meyer, Louraine Salabao, Jasmine Lee,  
Kimberlee Bradley, Svenja Hafter



Thank you!

Additional  
resources



# Templates

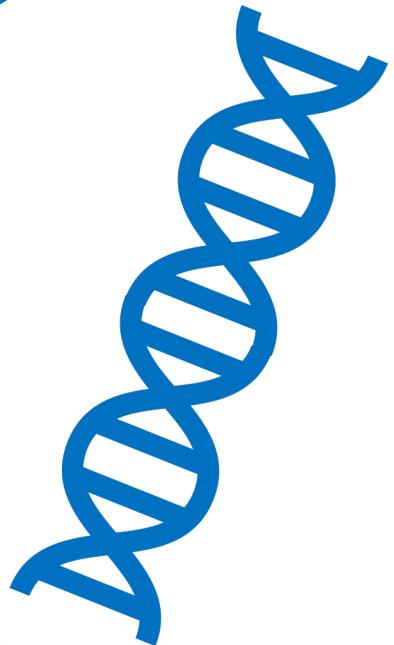
- <https://sios-svalbard.org/cgi-bin/darwinsheet/index.cgi?setup=darwin>

# Taxonomic data

- World Register of Marine Species (WoRMS) initial focus on marine species but extending to terrestrial and species traits. Also Register of Antarctic species
- Integrated Taxonomic Information System( ITIS) authoritative taxonomic information on plants, animals, fungi, and microbes of North America and the world.
- Species 2000 is an autonomous federation of taxonomic database custodians, involving taxonomists throughout the world
- Catalogue of Life is the most comprehensive and authoritative global index of species currently available. It consists of a single integrated species checklist and taxonomic hierarchy.

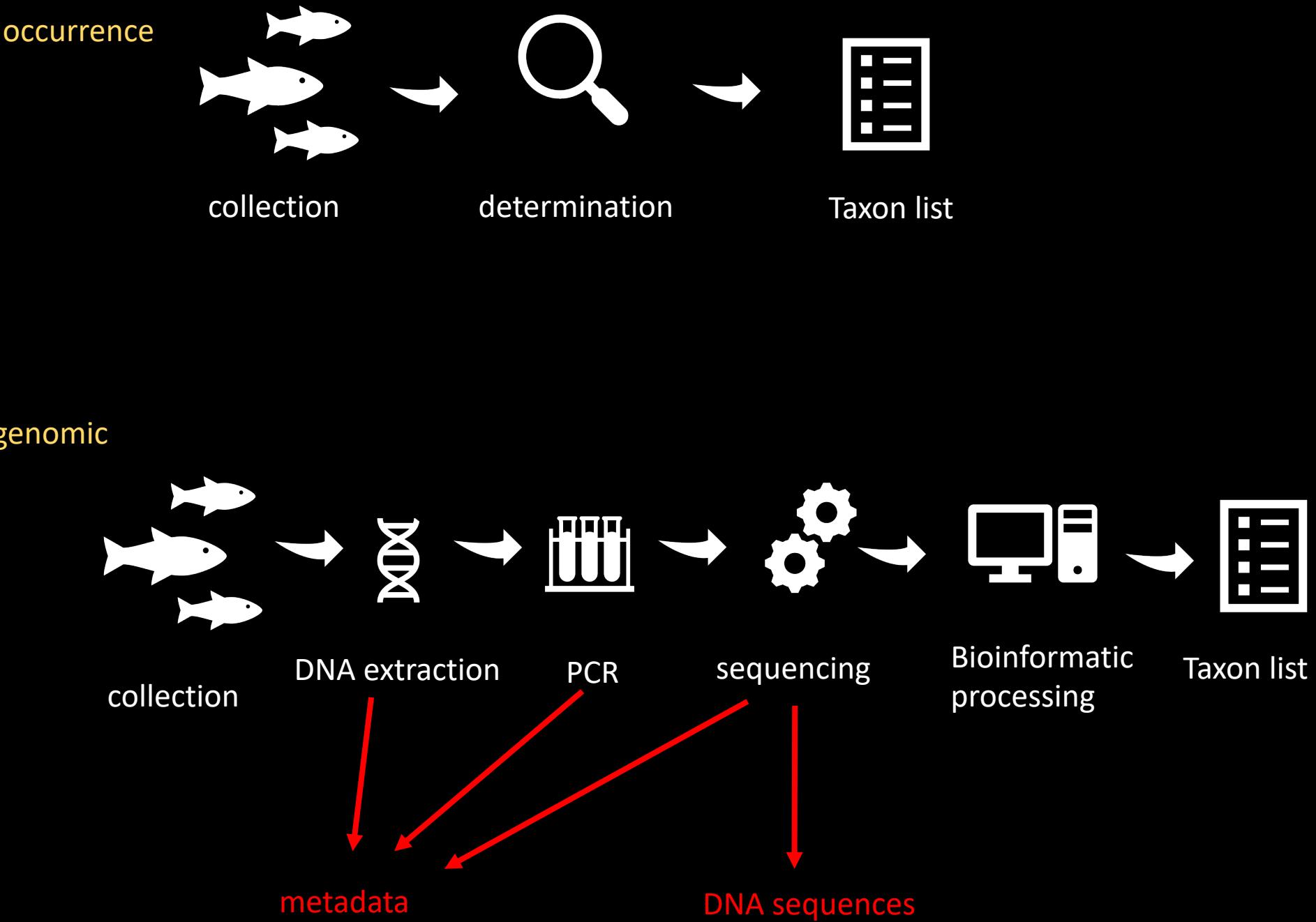
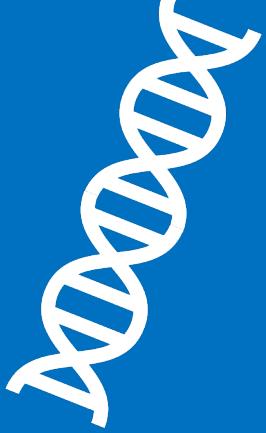
# Genomic data

- The [Genomic Standards Consortium \(GSC\)](#) aims to make genomic data discoverable. The GSC enables genomic data integration, discovery and comparison through international community-driven standards
- [Global Genome Biodiversity Network \(GGBN\)](#) network of well-managed collections of genomic tissue samples from across the Tree of Life, benefiting society through biodiversity research, development and conservation. This network will foster collaborations among repositories of molecular biodiversity in order to ensure quality standards, improve best practices, secure interoperability, and harmonize exchange of material in accordance with national and international legislation and conventions.

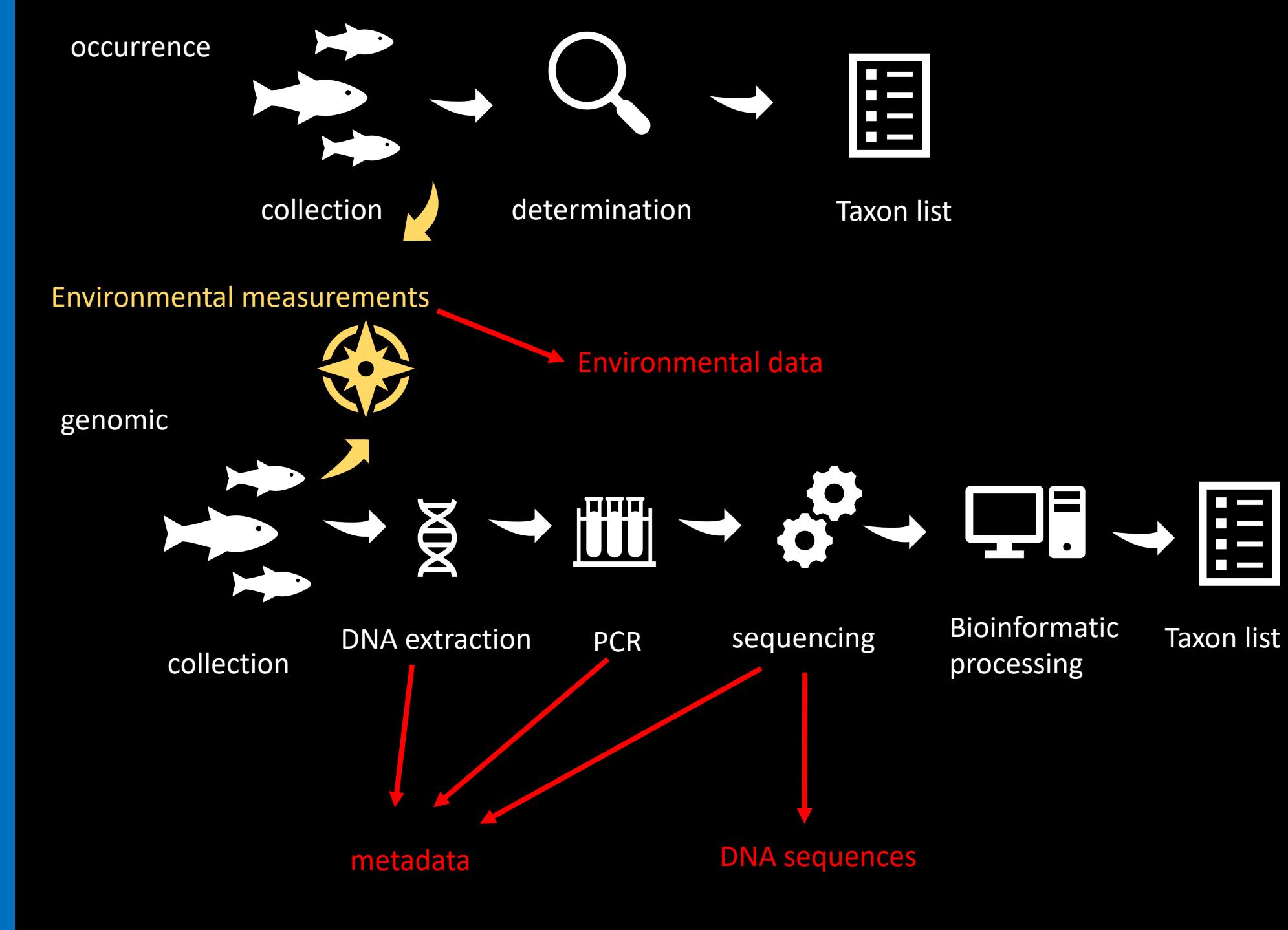
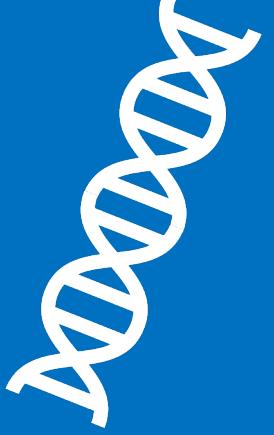


Example:  
genomics data

# Genomic data is complex



# Genomic data is complex





# Metadata and environmental data

- Every action that can affect the data
  - Sample collection, culturing or experimental settings
  - DNA extraction and sequencing protocols
  - essential to interpret the sequence data
- Associated Environmental measurements
  - Date + geographic coordinates
  - Any other measurement taken at the time
    - E.g. pH, ions, conductivity, weather conditions,...
- Put bluntly: sequence data is useless without



# MIxS standard

- Developed by: Genomics Standards Consortium (GSC)
- Minimum Information on Any Sequence (MIxS)
  - a list of standardized terms
    - = a minimum set to describe sequences
    - + any other relevant information
  - Allows standardization of environmental measurements
  - Accepted (required) by ENA and NCBI

List of the MIxS terms → <https://gensc.org/mixs/>



# MIxS example

Example: 2 samples on 1 sampling day:

Sample location 1: Gerlache strait seawater for 16S rRNA



Sample location 2:  
Anvers Island to look for  
soil metazoans (COI)



# MIxS example



	MIxS_Term_1	MIxS_Term_2	MIxS_Term_3				
1		value					
2							

Rows = sample = corresponds to an event

Columns = variable = standard MIxS terms

# MIxS example



	sample_name	project_name					
	2020_11_Gerlache_day1s1	webinaR2020					
	2020_11_Anvers_day1_s2	webinaR2020					

Cf. DarwinCore eventID

Cf. DarwinCore parentEventID

- Sample\_name:
  - choose something logical (e.g. place\_date\_letter\_number)
  - Only underscore or capitals, no special characters!

# MIxS example

sample_name	project_name	lat_lon	decimalLatitude	decimalLongitude			
2020_11_Gerlache_day1s1	webinar2020	-67.43 123.09	-67.43	123.09			
2020_11_Anvers_day1_s2	webinar2020	-69.43 123.59	-69.43	123.59			

- lat\_lon
  - Space separated
  - Decimal degrees
  - Write lat and lon in separate fields (e.g. “decimalLatitude”, “decimalLongitude”)





# MIxS: required terms

sample_name	project_name	decimalLatitude	decimalLongitude	collection_date	year	month	day
2020_11_Gerlache_day1s1	webinar2020	-67.43	123.09	2020-11-06	2020	11	06
2020_11_Anvers_day1_s2	webinar2020	-69.43	123.59	2020-11-06	2020	11	06

- collection\_date
  - Use ISO format (YYYY-MM-DD)
  - Best also include year, month, day



# MIxS: required terms

sample_name	project_name	geo_loc_name	seq_meth	env_biome	env_feature		
2020_11_Gerlache_day1s1	webinar2020	Gerlache Stait	Illumina MiSeq 300bp-PE	Ocean ENVO:01000048	Marine pelagic ENVO_01000044		
2020_11_Anvers_day1_s2	webinar2020	Anvers Island	Illumina MiSeq 300bp-PE	Polar ENVO:01000339	Glacial feature ENVO_00000131		

- Free text fields: try to standardize if possible
  - Use ontologies (e.g. ENVO: <http://bioportal.bioontology.org/ontologies/ENVO>)



# MIxS: environmental data

- The environmental package
  - Different packages of terms specific to one environment
  - E.g. soil, water, host\_associated,...
  - To provide some structure in long list of terms
  - Focused on associated measurements
    - soil: depth, pH,...
    - water: cond, turbidity, part\_org\_nitro, phosphate,...
  - Pick and mix, as long as you use as many standardized terms as possible



# MIxS: environmental data

sample_name	depth	part_org_nitro					
units	meter	micromole per liter					
2020_11_Gerlaче_day1s1	100	0.4					
2020_11_Anvers_day1_s2	0	NA					

- Include units, and write in full
- Also include any information that has no appropriate MIxS term
- New terms can be requested/proposed to GCS through their GitHub issues
- Only use “NA” where there is no information (not: “ND”, “NULL”, “0”,...)



# MIxS: sequencing terms

- Best also include
  - pcr\_primers
  - target\_gene; target\_subfragment
  - subspecf\_gen\_lin
    - = target taxonomic group
  - samp\_collect\_device and samp\_mat\_process
  - nucl\_acid\_amp; nucl\_acid\_ext; pcr\_cond,...



# DNA and RNA sequences

- FASTA format (=text file)

```
Header sequence 1 <--> NC_010533.1:3027-3709 Cryptopygus antarcticus mitochondrion  
sequence 1 <--> ATCCCCACATGAGCTTCTTAGGGTTCCAAAACGCAGCTTCTCCTCTAGAGCAATT  
Header sequence 2 <--> NC_010533.1:3845-4009 Cryptopygus antarcticus mitochondrion  
sequence 2 <--> ATCCCTCAAATAGCCCCATTAAGATGATTAATTCTATTTTATTTCGATTTATTTCT  
TTTAGCTAAAATTTTTCTAAAATAATCTTCATTTACTACACACCCGCTAATAA
```

- IUPAC notation of bases
  - Known bases: A, T, G, C
  - Ambiguous bases: W (=A or T), S (=C or G), R, Y, K, M, B, D, H, V, N



# DNA and RNA sequences

- FASTQ variant (=FASTA + quality score)

The diagram shows a sequence record in FASTQ variant format:

```
@NC_010533.1:3027-3709 Cryptopygus antarcticus mitochondrion
ATCCCCACATGAGCTTCTTAGGGTTCCAAAACGCAAGCTTCTCCTCTTAGAGCAATT
+
eeeffff>>efdccccdd"[[[_BA^YBBBBBBBTTTBYY^^^[_["dddcccd...[][[[]'dd c
```

Annotations with red arrows and text:

- Header sequence 1: Points to the header line starting with '@'.
- Start with @: Points to the '@' symbol.
- sequence 1: Points to the first sequence line.
- Header sequence 1 repeated or blank: Points to the header line starting with '+'.
- Start with +: Points to the '+' symbol.
- Quality score sequence 1: Points to the quality score line.

A red circle highlights the '@' symbol in the header. Another red circle highlights the '+' symbol in the header. A third red circle highlights a single character in the quality score line.

Per base Q-score,  
single ASCII character  
e.g. B = 0.6309 chance of being correct



# MIxS useful terms

sample_name	pcr_primers	target_gene	subspecf_gen_lin	nucl_acid_ext			
2020_11_Gerlache_day1s1	ATTGTGAAGT AAGGTGTCW	16S rRNA	Bacteria Archaea	<a href="https://dx.doi.org/10.17504/protocols.io.bdwsi7ee">dx.doi.org/10.17504/protocols.io.bdwsi7ee</a>			
2020_11_Anvers_day1_s2	CCTGHAATCG GTTGGAGCA	COI	Metazoans	<a href="https://dx.doi.org/10.17504/protocols.io.bdwsi7ee">dx.doi.org/10.17504/protocols.io.bdwsi7ee</a>			

- Don't use “, “ as separator (cf. CSV file)
- For protocols: use DOIs or URLs to the publications



# Using MIxS

- Possible to upload alongside sequences to NCBI or ENA
- Can be added to a GBIF DarwinCore archive
  - Using event core (=list of sequence samples)
  - Add MIxS file as extended Measurement or Fact (eMoF) extension
  - See <https://docs.gbif-uat.org/publishing-dna-derived-data/1.0/en/#environmental-dna-as-a-source-for-dna-derived-occurrence-data>
- List the dataset on POLA3R
  - Specific data repository of any Antarctic sequence data
  - internal database (searchable metadata)
  - link to sequences on INSDC (not in MIxS)
  - [www.biodiversity.aq/pola3r](http://www.biodiversity.aq/pola3r)





# Thank You



# way to deal with metadata

1. In the publication
  - Not physically linked with the sequence data
  - Often too concise to be fully reproducible
  - No standardization
  - + High visibility in the scientific community
2. Register to a biodiversity repository (e.g. GBIF or POLA3R)
  - + Better linkage to the sequence data
  - + May include environmental metadata
  - + Possibility to be very extensive
  - + Standardized format (EML=Ecological Metadata Language and/or MIxS)
3. Attach to the sequence data on ENA or NCBI
  - But difficult to find or download
  - + Physically linked with the sequence data
  - + Standardized format (MIxS)
  - + Ideal for PRE-sampling

# Register metadata a biodiversity repository

<https://ipt.biodiversity.aq>

**INTEGRATED PUBLISHING TOOLKIT<sup>(IPT)</sup>**  
free and open access to biodiversity data

Logged in as msweetlove@naturalsciences.be [ACCOUNT](#) [LOGOUT](#) [ENGLISH](#)

Resource Title [Polychaetes from the JR17003 Expedition](#)

## Basic Metadata

Please enter all the mandatory properties on the Basic Metadata page, and then continue entering metadata in the other pages that are applicable to your resource. The more metadata you provide, the greater the chance that your resource will be found, reused by other researchers, and cited.

Title\*

Polychaetes from the JR17003 Expedition

Publishing Organisation\*

Select an organisation

Update Frequency\*

Unknown

Data Licence\*

No licence selected

Type\*

Occurrence

Subtype

Select a subtype

Metadata Language\*

English

Data Language\*

English

## Section

[Basic Metadata](#)  
[Geographic Coverage](#)  
[Taxonomic Coverage](#)  
[Temporal Coverage](#)  
[Keywords](#)  
[Associated Parties](#)  
[Project Data](#)  
[Sampling Methods](#)  
[Citations](#)  
[Collection Data](#)  
[External links](#)  
[Additional Metadata](#)

Resource Title [Polychaetes from the JR17003 Expedition](#)

## Sampling Methods

Please enter metadata about the sampling methods used for the data represented by the resource.

Study Extent\*



Sampling Description\*



Quality Control



Step Description\*



[Add new method step](#)

## Section

[Basic Metadata](#)

[Geographic Coverage](#)

[Taxonomic Coverage](#)

[Temporal Coverage](#)

[Keywords](#)

[Associated Parties](#)

[Project Data](#)

[Sampling Methods](#)

[Citations](#)

[Collection Data](#)

[External links](#)

[Additional Metadata](#)

This XML file does not appear to have any style information associated with it. The document tree is shown below.

```
<?xml version="1.0" encoding="UTF-8"?>
<eml:eml xmlns:eml="eml://ecoinformatics.org/eml-2.1.1" xmlns:dc="http://purl.org/dc/terms/" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
xsi:schemaLocation="eml://ecoinformatics.org/eml-2.1.1 http://rs.gbif.org/schema/eml-gbif-profile/1.1/eml.xsd" packageId="https://ipt.biodiversity.aq/resource?
id=2503_antarctic_nemertea/v1.0" system="http://gbif.org" scope="system" xml:lang="eng">
  <dataset>
    <alternateIdentifier>
      https://ipt.biodiversity.aq/resource?r=2503_antarctic_nemertea
    </alternateIdentifier>
    <title xml:lang="eng">A collection of Antarctic Nemertea</title>
    <creator>
      <individualName>
        <givenName>Claude</givenName>
        <surName>De Boyer</surName>
      </individualName>
      <organizationName>Belgian Royal Institute of Natural Sciences</organizationName>
    </creator>
    <address>
      <city>Brussels</city>
      <country>BE</country>
    </address>
  </creator>
  <metadataProvider>
    <individualName>
      <givenName>Maxime</givenName>
      <surName>Sweetlove</surName>
    </individualName>
    <organizationName>Belgian Royal Institute of Natural Sciences</organizationName>
    <address>
      <city>Brussels</city>
      <country>BE</country>
    </address>
    <electronicMailAddress>msweetlove@naturalsciences.be</electronicMailAddress>
  </metadataProvider>
  <associatedParty>
    <individualName>
      <givenName>Claude</givenName>
      <surName>De Boyer</surName>
    </individualName>
    <organizationName>Belgian Royal Institute of Natural Sciences</organizationName>
    <address>
      <city>Brussels</city>
      <country>BE</country>
    </address>
    <role>user</role>
  </associatedParty>
  <pubDate>2020-09-10</pubDate>
  <language>en</language>
  <abstract>
    A collection of Antarctic Nemertea
  </abstract>
</dataset>
</eml:eml>
```



# Date range

ISO 8601 standard

Start date / End date

YYYY - MM - DD / YYYY - MM - DD

YYYY - MM / YYYY - MM

YYYY / YYYY

etc etc ...

# Date, time intervals

ISO 8601 standard

2007/2008

2007-11/12

2007-11-13/15

2007-03-01/2008-05-11

2007-12-14T13:30/15:30Z

2007-03-01T13:00:00Z/2008-05-11T15:30:00Z