# Recording Sample Metadata for the European Reference Genome Atlas Project*

## Sample Manifest Standard Operating Procedure

**Version: 1**

**Published Date:**

**Authors: Mara Lawniczak, Rob Davey, Mark Blaxter, Jennifer A. Leonard, Olga Vinnere Petterson, Seanna McTaggert, Ann McCartney, Astrid Böhne, and the ERGA Consortium**

Correct, ethical, and comprehensive recording of sample metadata is critical to the long-term utility of the work we do in the ERGA project: these metadata will link the genome sequences to their origins and remaining voucher material, and weave our work into the rich fabric of understanding of European eukaryotic biodiversity. Please read this Standard Operating Procedure (SOP) in full before completing the Sample Manifest as it contains detailed guidance on how to record metadata. Also contained is generic guidance on how to process specimens. Help for taxon-specific SOPs is available from each taxonomic working group (see contacts listed below) to provide guidance on sample processing and regulatory compliance. Guidance on sample submission to sequencing partners should be seeked from the sequencing facility (see list of sequencing partners below). Guidance on submission of vouchering and biobanking material should be seeked from the respective collection facility (see list of partners below). Suggestions to this manifest should be send to the sampling committee SSP [samples@erga-biodiversity.eu](mailto:samples@erga-biodiversity.eu)

*This sampling manifest builds on the work of the Darwin Tree of Life DToL sampling committee, in particular of Mara Lawniczak and Robert Davey.

**The ERGA consortium thanks DToL for allowing ERGA to adopt their sampling metadata manifest to ERGA's needs and for document sharing. All changes to this document only apply to ERGA not to DToL.**

**Preamble:** To be able to register a sample/specimen and its metadata for ERGA, the submitting person (most often identical to the species ambassador) must subscribe to the ERGA membership form confirming the ERGA code of conduct as well as confirm adhering to ERGA's ethical code of conduct for sampling.

**Purpose and responsibility**: ERGA aims to generate high quality genome sequences from samples and to embed these sequences into best practices in scientific research and the landscape of biodiversity science. To do this we must adhere to correct, legal and ethical physical handling of the specimens, and correct collation of rich metadata describing the specimens. This SOP contains specific instructions for filling in the metadata manifest. ERGA will not accession and process samples without complete associated metadata, have not been sampled in compliance with legal rules applying to each specimen and have not been sampled according to an ethical code of conduct. The legal responsibility for acquiring samples remains with the species' ambassador(s). By submitting the sampling manifest and providing information on compliance with sampling permits, the ambassador guarantees that the sample in question can be legally transferred to the sequencing center indicated, can be vouchered at the indicated collection(s), and has been sampled in compliance with all applicable rules. The responsibility for the oversight of all legal compliances remains with the species ambassador. Where necessary and applicable, material transfer agreements can be issued and signed by sequencing and collection facilities and the species ambassador; the oversight of this process lies with the species' ambassador(s).

**Additional SOPs:** Additional related SOPs are forthcoming on how to prepare samples for different taxonomic groups, which helps to assure delivery of high-quality samples that are more likely to be transformed into high quality genomes. As of now, help for actual sampling and vouchering can be requested from taxon chaperons of the ERGA SSP committee taxonomic focus groups as listed below.

**Future plans for this SOP:** This SOP will be reviewed on a quarterly basis by the SSP committee to incorporate feedback from the community. Metadata are currently collected manually by the species ambassador using a defined spreadsheet, referred to as the ERGA SAMPLE MANIFEST V1. This document will allow integration into the data management and brokering platform system COPO (http://copo-project.org. COPO allows for dry runs of metadata upload to validate compliance to format requirements. COPO will link to a database that tracks all samples and their associated metadata as they progress from collection to genome assembly. Finally, the sequencing data will be archived in the ENA (https://www.ebi.ac.uk/ena/browser) for all sequenced samples with the information provided in the metadata.

**Raising issues**: Elements of this SOP are subject to discussion, development and change. We expect that there will be questions to answer and lessons learned to share. Please raise specific issues by emailing the Samples Committee at **samples@erga-biodiversity.eu** or reach out on Keybase, Team erga.listserv, channel #Committee_SSP

## *Document History*

| *Major Version* | *Date* | *Changes* | *Contributors* |
|---|---|---|---|
| **1.0** | 2021-09-17 | first version | Mara Lawniczak, Jeena Rajan, Robert P Davey, Seanna McTaggart, Mark Blaxter, Alice Minotto, Felix Shaw, DToL SamplesWG, ERGA SSP committee, ERGA ELSI committee, Ann McCartney, Jennifer A Leonard, Olga Vinnere Petterson, Astrid Böhne |
| | | | |

# Completing the Sample Manifest: Overview

## *Scope of this document*

**Specific guidance on preparing samples** is not covered by this SOP. Please contact the chaperon for the specific taxonomic group you are working on.

**Submission of samples** is also not covered by this SOP. Please contact the involved sequencing partner to obtain necessary information .

## *The importance of "SPECIMEN_ID"*

The ERGA SPECIMEN_ID must reflect the genetic identity of the individual, serving to link the various samples, images, vouchers, DNA barcodes, etc. that derive from one individual organism together. The ERGA SPECIMEN_ID also allows the laboratory team to resample the same individual specimen (thus the same haplotypes) if needed, e.g., in the case of requiring more DNA to create a library. For example, ten different individual specimens each in their own tube would have ten distinct ERGA SPECIMEN_IDs, even if they are all from the same species or clone or culture. However, a single specimen split across ten tubes would result in each of those ten tubes having the same ERGA SPECIMEN_ID. This unique ERGA SPECIMEN_ID has three critical functions: identifying the species ambassador that holds responsibility for the specimen, tracking an individual sample's status and declaring the genetic uniqueness of the specimen.

Each ERGA specimen must be linked to a standardized, unique ID that begins with the prefix **ERGA_** followed by **SPECIES AMBASSADOR INITIALS** (up to 10 letters out of A-Z, if this is not possible please reach out to **samples@erga-biodiversity.eu**), followed by **underscore** followed by the l**ast four digits of the SPECIES AMBASSADOR's ORCID ID**, **underscore**, and running numbers (e.g. if you as ambassador register more than one sample make sure to use 01, 02…). SPECIMEN_IDs must be unique to an individual (e.g., ERGA_XY_1234_01 cannot be used again after it has been assigned to a specimen). SPECIMEN_IDs must follow the format described above.
Furthermore, Genome Acquisition Labs GALs will attribute identifiers to each SPECIMEN in the format following the table below for their internal sample tracking. These identifiers need to be provided from the Sequencing partners to the ambassadors.

## *Table: Genome Acquisition Labs Specimen Codes*          *GAL*

| GAL | Code Model | Number of digits | Contact Person Email address* |
|---|---|---|---|
| SANGER | SAN0000000 | 7 | Nancy Holroyd, neh@sanger.ac.uk |
| EARLHAM | EI_00000 | 5 | Seanna McTaggart, seanna.mctaggart@earlham.ac.uk |
| CNAG | | | Project Manager, projectmanager@cnag.crg.eu |
| SciLifeLab | | | Olga Vinnere Pettersson, olga.pettersson@igp.uu.se |
| WGGC West German Genome Centre | | | Peter Nürnberg, nuernberg@uni-koeln.de Antonella Succurro, a.succurro@uni-bonn.de |
| NCCT NGS Competence Center Tübingen | | | Nicolas Casadei, Nicolas.Casadei@med.uni-tuebingen.de |
| DresdenConcept | | | Sylke Winkler, winkler@mpi-cbg.de |
| FGC Zurich | | | Andrea Patrignani, andrea.patrignani@fgcz.ethz.ch |
| Genoscope | | | Pedro Oliveira, pcoutool@genoscope.cns.fr |
| GTF Lausanne | | | Julien Marquis, contactGTF@unil.ch |
| DNA sequencing and genomics laboratory, Helsinki Genomics Core Facility | | | Petri Auvinen, petri.auvinen@helsinki.fi |
| NGS Bern | | | |
| Norwegian Sequencing Centre (NSC) | | | Ave Tooming, ave.tooming-klunderud@ibv.uio.no |
| University of Bari | | | |
| University of | | | Claudio Ciofi, |

| | | | |
|---|---|---|---|
| Florence | | | claudio.ciofi@unifi.it |
| Neuromics Support Facility, UAntwerp, VIB | | | Mojca Strazisar, Mojca.Strazisar@uantwerpen.vib.be |
| GIGA-Genomics core facility University of Liège | | | Wouter Coppieters, wouter.coppieters@uliege.be |

\* As of September 2021

## *Other "_ID"s*

A sample can represent a set of specimens as well as multiple parts of the same specimen, and so the Genome Acquisition Labs GAL_SAMPLE_IDs and COLLECTOR_SAMPLE_IDs can refer to an individual organism or something else (e.g., a soil sample could be represented by the COLLECTOR_SAMPLE_ID and a specimen taken from within that collection of soil be assigned a SPECIMEN_ID). The COLLECTOR_SAMPLE_ID is the identifier assigned by the collector to the specimen or the sample, hence the use of the term SAMPLE rather than SPECIMEN in this metadata field. The same is true of the GAL_SAMPLE_ID. For example, if a collector collects a sample that could have mixed genotypes or species, this will have a single COLLECTOR_SAMPLE_ID, and will need to be split further into specimens, each of which is assigned a unique ERGA_SPECIMEN_ID.

It is permitted to have identical names for any or all of three categories (COLLECTOR_SAMPLE_ID, GAL_SAMPLE_ID and ERGA_SPECIMEN_ID). The SPECIMEN_ID is the only one that is required for sequencing to commence.

Management of COLLECTOR_SAMPLE_ID, GAL_SAMPLE_ID and their relationship to ERGA_SPECIMEN_ID is the responsibility of the species ambassador and the GAL.

## *Manifest Validation Process*

Choose whether you prefer to use the Sample Manifest from the google spreadsheet or another option (e.g., epicollect, ARCGIS). We recommend that you retain a copy in Excel (XLS/XLSX) or Google spreadsheet form so as not to lose the data validation given the likelihood that further edits will be required, even after upload to COPO

**Google spreadsheet:** The Google spreadsheet can be used by ***making a copy*** and using it as an online spreadsheet, or by downloading it and entering data locally. If you choose to do the latter, please download as an XLS/XLSX (Microsoft Excel format) file to ensure that the data validation fields are retained.

Please carefully read the guidance in this SOP for each field, and attempt to get your submitted manifests as close to the guidance as possible (you will be provided with a mock example to facilitate input). If your sample requires metadata fields or terms that are not present in the manifest, please contact samples@erga-biodiversity.eu to discuss and define new fields or terms.

Once you have completed entering all metadata, the initial check will confirm that each TAXON_ID maps to the correct species name. If mismatches are found, this will require the submitter to examine the mismatches and determine the nature of the problem. Please read the guidance on TAXON_ID below carefully as you should be able to ensure that each TAXON_ID precisely and accurately matches a species name in advance of submitting your manifest. There are too many possibilities to enumerate them all here, but three of the most common issues are a misspelling in the SCIENTIFIC_NAME or the TAXON_ID fields, a species for which no TaxonID is available in the NCBI TaxonomyDB, or a change in the taxonomy not reflected in NCBI TaxonomyDB. These will need to be addressed before the manifest can be validated. More information on how to fix these issues is below in the discussion of the TAXON_ID field.

Once you have ensured that your manifest is ready for validation, follow the guidance of the GAL and collections (if applicable at this stage) involved for sample submission. If any other issues with the information provided within the sample manifest are identified (e.g., missing mandatory entries, duplicate rows, incorrect date formats), the sample manifest will be returned to you to resolve these issues; within COPO, this will be an iterative process pointing you towards mal formatted or missing information.

Once this process is complete and every sample has a TAXON_ID together with complete metadata, the manifest is considered to be "validated". However, prior to samples being accepted at the sequencing institute, DNA barcoding data may be required. Manifests can be validated and held until barcoding results are back and relevant fields (e.g., SCIENTIFIC_NAME, PUBLIC_NAME, TAXON_ID) can be updated. The process for "updating" a validated manifest will be developed over the coming months. COPO has a version history.

For samples, if there is any possibility of species misidentification (SPECIMEN_RISK = Y), samples will only be accepted for genome sequencing after DNA barcoding data is returned and samples are confirmed as the species they were declared . At this stage, each sample will be allocated a "PUBLIC_NAME" that reflects both the species and the SPECIMEN_ID (i.e., the genetic identity of the sample).

When data are submitted to ENA for release (as part of BioSample, raw data and assembly submissions), the submissions will include all of the fields below indicated by **ENA_submission**. If the field name is in **turquoise**, then an entry for each specimen is mandatory for that field, even if only to declare why the information is missing. For all other fields, we strongly encourage data entry, but it is not mandatory if it has not been collected.

## *Vouchers of Specimen or Sample*

Every submitted specimen has to be accompanied by voucher material. This material should be accessioned by a registered collection for permanent storage. Physical voucher material may be separated on collection, and be submitted directly to the designated collection organisation, or material remaining after processing may be returned to the designated collection from the sequencing centre if the collection permits. In cases where the entire specimen is consumed by processing, we request that digital images are recorded prior to destructive sampling and submitted in lieu of physical samples. We regard it as good practice to record digital images of all specimens and samples destined for ERGA processing, whether or not physical vouchers are retained, as this provides a close-to-life record of the organism sampled (see below).

## *Photographs of Specimen or Sample*

Every submitted specimen should be accompanied by a photograph with explicit labelling as described below. Please upload your image, as of now for each SPECIMEN_ID, COPO will create a directory and all necessary subdirectories in ERGA's reserved and secured space on a B2Drop instance on Nextcloud in which all documents associated with the registered specimen will be stored. Images need to be dropped to the subdirectory called "IMAGES".

In preparation for linking images to metadata, please name images using the following format: SPECIMEN_ID-X.Y where X is a numerical identifier for the number of photographs you have taken of the same individual, and Y is the file format, e.g., NHMUK014110995-1.png and NHMUK014110995-2.png for two photos of the same specimen provided. When uploading photographs, please use PNG or JPG format.

File names must exactly match the SPECIMEN_ID in order to match photographs to samples automatically.

# Detailed instructions for filling in the Sample Manifest

I.    The manifest has several tabs. Please only fill in the **Metadata Entry** tab. If you discover a missing attribute in the drop-down menus, new attributes can be suggested by raising a request to the SSP committee at **samples@erga-biodiversity.eu**. Please only do this if absolutely required (i.e., no available term is a good proxy, and the absence of the attribute likely to affect many samples).

II.    **Information must be entered for all fields below with** ==turquoise bold names== [in the Google spreadsheet version of the manifest, these fields are represented by cells with purple fill. The fill will go white when an entry has been made to help you identify where mandatory fields still require data.] For all mandatory fields with ==turquoise bold names==, even if information is unavailable, they must be populated with the appropriate term describing why this information is missing. The acceptable missing value terms are:

        **NOT_APPLICABLE** = information is inappropriate to report. This can also indicate that the standard itself fails to model or represent the information appropriately.

        **NOT_COLLECTED** = information was not given because it has not been collected.

        **NOT_PROVIDED** = information of an expected format was not given but a value may be given at the later stage (this may be a particularly useful missing information term for VOUCHER_ID)

Fields that are named in **BOLD** without color do not require an entry describing why the information is missing because we expect that many samples will not have information for these fields (e.g., most samples will not have DEPTH information). However, if you have collected the information related to these terms, please do enter it.

Many terms will have the data released publicly as part of the ENA record. For every field for which this is true, you will find "**ENA_submission**" next to the name of the term.

III.    **All dates** in the manifest must be formatted consistently as **YYYY-MM-DD** (ISO8601).

IV.    In fields that are "free text", we ask that you use only the core alphanumeric characters, plus full stop ".", hyphen "-", underscore "_" and spaces (summarised in coding parlance as "`-_.a-zA-Z0-9`"). Please avoid "|" (the vertical pipe symbol) except where we indicate it should be used to separate elements in a list. Please **do not** use "special characters" (such as other punctuation and "logical" marks: "`#"';:?!@*()[]{}/\,=+`", etc.).

## *Column by column instructions for the Metadata Entry tab*.

A. <mark>TUBE_OR_WELL_ID</mark>: This field should record the FluidX barcode for each tube in a rack (or each well in a plate, where relevant) if available, else the position of the well if submitted in a well plate or the label on the submitted tubes. If barcodes are entered, use a barcode scanner in advance of preparing samples to reduce errors – do not enter barcodes manually.

B. <mark>SAMPLE_COORDINATOR:</mark> (ENA_submission) Also known as the ERGA ambassador, Enter the name of the person or people who is responsible for the genome project of the sample using all CAPITALS, and separate names with "|" (vertical pipe symbol), e.g., "CAROLUS LINNAEUS | JEAN_BAPTISTE LAMARCK".

- We note that storage of names with affiliations in a database brings the system under the aegis of the GDPR regulations, and we must ask all involved to agree to their data being stored in COPO and to those data being propagated to secondary databases (including ENA and the final collections of record).

C. <mark>SAMPLE_COORDINATOR_AFFILIATION</mark>: (ENA_submission) Free text field to supply the university, institution, or society that is responsible for the genome project of the sample. This is typically the society or institution of the person(s) specified in the SAMPLE_COORDINATOR field. If multiple people are specified in COLLECTED_BY, ensure that their institutional affiliations are also separated by a vertical pipe symbol. Position in the list of affiliations should match the person in the same position in the list of names (e.g., PERSON A | PERSON X | PERSON C will have their affiliations as: (INSTITUTE A | INSTITUTE X | INSTITUTE C). If multiple people are listed but all from the same affiliation, no need to repeat the affiliation.

D. <mark>SAMPLE_COORDINATOR_ORCID_ID</mark>: (ENA_submission) Enter the 16 digits ORCID ID of the person or people who is responsible for the genome project of the sample.

E. <mark>SPECIMEN_ID</mark>: (ENA_submission) This is a unique identifier that refers to the genetic identity of the supplied material. It is assumed that the ERGA SPECIMEN_ID refers to a singular genetic individual. If the same individual specimen is split into several samples submitted in separate tubes, the ERGA SPECIMEN_ID for these samples would be the same. If multiple individuals of a species are sampled (e.g., from the same population), they must be placed in multiple, individual tubes, each with a unique ERGA SPECIMEN_ID. If sampling from organisms where distinguishing genetic individuals is difficult (e.g., mat-forming species like mosses or bryozoans), tease out individual units as far as is possible (e.g., single strands from a moss mat), and place each in a separate specimen tube with a unique ERGA SPECIMEN_ID. Each ERGA specimen must be linked to a standardized, unique ID that begins with the prefix ERGA_ followed by SPECIES AMBASSADOR INITIALS (up to 10 letters out of A-Z, if this is not possible please reach out to **samples@erga-biodiversity.eu**) _underscore and SPECIES AMBASSADORS AFFILIATION COUNTRY (alpha 2 ISO3166 letter codes, in case your place of affiliation has no ISO3166 letter code

please reach out to **samples@erga-biodiversity.eu** XXXX to integrate an appropriate 2 letter code for you) and _underscore running numbers (e.g. if you as ambassador register more than one sample make sure to use 01, 02…). SPECIMEN_IDs must be unique to an individual (e.g., ERGA_XY_DE_01 cannot be used again after it has been assigned to a specimen). SPECIMEN_IDs must follow the format described above.

F. ORDER_OR_GROUP: The taxonomic Order into which the Family is placed or (if this is not defined) the monophyletic group to which the Family or Genus belongs. This should correspond to the taxonomy as represented in the NCBI Taxonomy Database. If you or your taxonomist have a disagreement with the taxonomy represented on NCBI Taxonomy Database, please raise this with the NCBI TaxonomyDB curators as described below.

G. FAMILY: The taxonomic Family into which the Genus is placed. This should correspond to the taxonomy as represented in the NCBI Taxonomy Database. If you or your taxonomist have a disagreement with the taxonomy represented on NCBI Taxonomy Database, please raise this with the NCBI TaxonomyDB curators as described below.

H. GENUS: The taxonomic Genus to which the Species belongs. This should correspond to the taxonomy as represented in the NCBI Taxonomy Database, and with the generic component of the scientific name given below. If you or your taxonomist have a disagreement with the taxonomy represented on NCBI Taxonomy Database, please raise this with the NCBI TaxonomyDB curators as described below.

I. TAXON_ID: (ENA_submission) A valid NCBI TAXON_ID to the species level is mandatory in order to submit data to public repositories. The species name in the manifest must be identical to that listed in the "current name" box in the Taxonomy Browser for that species. If this is not the case, you must write to ena-dtol@ebi.ac.uk to request the change.

> If there is another taxon database for your group, e.g., EukRef, please fill in the NCBI TAXON_ID, and then use the TAXON_REMARKS field to specify the taxon database and the ID/accession/URL.

> > ■ TAXON_IDs can be looked up based on the species at the following links: https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi

> > or https://www.ncbi.nlm.nih.gov/Taxonomy/TaxIdentifier/tax_identifier.cgi.

> > ■ If no TAXON_ID exists, or a credible TAXON_ID exists that likely is a synonym of the species name the collector or submitter would use (through differential usage, error or lack of currency of the NCBI taxonomy), please ask for assistance by writing to ena-dtol@ebi.ac.uk, providing the full name, authority, and publication for the chosen name where possible. If required (e.g.,

newly described species, species missing from taxonomy browser), a new TAXON_ID should be available within 14 days. In the case of conflict, the sample submitter will be contacted and may be required to provide further information. Please note that the final species name on submission of the data to INSDC will be the one associated with the TAXON_ID in NCBI Taxonomy.

■ When a sample is provided that requires DNA barcoding before a species ID is possible, please provide the lowest taxonomic rank identification as possible (ORDER_OR_GROUP, FAMILY, GENUS) and leave SCIENTIFIC_NAME blank. You may care to place comments on what the specimen is likely to be in TAXON_REMARKS.

J. **SCIENTIFIC_NAME**: (**ENA_submission**) The latin binomial/combined genus and species name with a space in between.

■ See TAXON_ID above if you or the taxonomic expert have substantive issues with the species name present for the taxon in the NCBI TaxonomyDB.

K. **TAXON_REMARKS**: Free text to summarise any known issues with the mapping of TAXON_ID to SCIENTIFIC_NAME or add other taxon database identifiers here e.g., EukRef. Here you can also comment on STRAIN availability, if the specimen is a representative of a living and accessible strain/colony/culture. If there are no issues, leave this field **blank**.

L. **INFRASPECIFIC_EPITHET**: Where the sample is from a formally named infraspecific taxon, give the infraspecific name here, with prefixes in the following format: ssp. (for subspecies), var. (for variety), cv. (for cultivar), br. (for breed). Entries in this field should reflect organisms that can be found living outside of laboratories (see next attribute for lab strains). If there is no epithet here, leave this field **blank**.

M. **CULTURE_OR_STRAIN_ID:** (**ENA_submission**) Please give the reference ID from the source culture collection, such that the culture accession can be found in the collection's database. This is only relevant if the sequenced material is derived from a living, culturable, named laboratory strain (e.g., *Anopheles coluzzii* N'Gousso strain). This field should not be used to record a variant or type that has been collected anew from the wild: such information should be placed in **OTHER_INFORMATION**. Leave this field **blank** if it is not relevant.

N. **COMMON_NAME**: Vernacular name, if the species has one. If multiple names are required, separate names with a | (vertical pipe) character. If you are unsure of or the species has no vernacular name leave this field **blank**.

O. **LIFESTAGE**: (**ENA_submission**) The life stage of the specimen from which the sample was derived. This field has a controlled vocabulary: use the drop-down menu or look at the available terms on the second tab to complete. Please note that there are currently curated attributes for animals, for plants/fungi/macroalgae, and for some protists.

- If these do not fit your taxa, please contact samples@erga-biodiversity.eu. Please enter **NOT_PROVIDED** if your proposal for a lifestage term has not yet been accepted.

P. <mark>SEX</mark>: (<mark>ENA_submission</mark>) The sex of the specimen from which the sample was derived. This field has a controlled vocabulary: use the drop-down menu. If the sex of the organism is not known, use **NOT_COLLECTED**. The sex may be determined at a later date using the genome sequence data, but this will be captured in a different field, so this field should refer solely to the sex as determined by morphological examination of the specimen or strong inference (e.g., the species is from a clade that is always hermaphroditic/monoecious).

Q. <mark>ORGANISM_PART</mark>: (<mark>ENA_submission</mark>) A description of the exact tissue(s) in the tube or well. Accurate information here is important for downstream analyses on the symbiome, chromosomal diminution, RNAseq, etc. There is a tab in the Sample Manifest that defines the terms that can be used for ORGANISM_PART. This tab lists definitions for the full tissue, but pieces of that tissue are acceptable (e.g., LUNG is defined as 'the lung of a vertebrate', but the whole lung is not expected and a small piece of lung is expected).

- Please combine tissues by entering multiple terms from the ontology using the | (vertical pipe) symbol (e.g., for head + abdomen of an insect enter "HEAD | ABDOMEN"). When using multiple body parts, there will be a data validation error that arises, but these can be ignored as long as the spelling and capitalization of the terms is identical to the provided list.

- If the tissue or organism part you are providing is not present in the drop- down menu, please choose the best generic category (these start with **) and add the name of the tissue that you have put into the tube in the "OTHER_INFORMATION" free text field. Please also email the ERGA SSP committee at samples@erga-biodiversity.eu to request the necessary additions. We will update attributes quarterly.

- If the sample is shipped as a DNA or RNA extract, select the tissue from which this was extracted and add further information in the OTHER_INFORMATION field regarding quality, quantity, etc. Note that any shipment of DNA should be discussed in advance with the involved GAL.

R. **SYMBIONT**:   This is to indicate whether the sample contains a known endo- or ectosymbiont (i.e., you have metadata for it and a species-level and ENA-submittable TAXON ID). Select "TARGET" if only the "host" metadata is known OR if it is a symbiont-only culture. Select "SYMBIONT" if you have a known symbiont in the sample and you have metadata (including, critically, a species-level identification supported by a valid taxon ID) for the symbiont. If you need to select "SYMBIONT" you will then need to copy and paste your "TARGET" row and amend the following fields to reflect the symbiont data:

■ORDER_OR_GROUP, FAMILY, GENUS, TAXON_ID, SCIENTIFIC_NAME, TAXON_REMARKS, INFRASPECIFIC_EPITHET, CULTURE_OR_STRAIN_ID, COMMON_NAME, LIFESTAGE, SEX, ORGANISM_PART

If there is no explicit information on potential symbionts, this field should be left **blank**.

S. **RELATIONSHIP:** (==ENA_submission==) This is a free text field to permit declaration of any known parental, child, or sibling relationship between the specimen and any other specimens that are submitted for the ERGA project, OR to declare if the specimen is a "barcode exemplar" for another specimen.

- If there are known genetic relationships between submitted specimens, please concisely state the relationship: "Full sibling to SPECIMEN_ID1", "Mother to SPECIMEN_ID2", "Maternal half sibling to SPECIMEN_ID1, SPECIMEN_ID2, and SPECIMEN_ID3", or "Trio child of SPECIMEN_ID1 and SPECIMEN_ID2". If knowledge of the relationships is not confident but suspected, do not add anything here and instead add this information to the "OTHER_INFORMATION" field (e.g., "suspected full or half sibling to SPECIMEN_ID2").

- If the specimen is acting as a barcoding exemplar or a proxy voucher for another specimen because the entire organism must be used for reference genome sequencing and it is not possible to take a sample for DNA barcoding or biobanking or vouchering (e.g., midges from the same swarm where one is submitted for sequencing and 5 are submitted individually for DNA barcoding), then add "barcode/voucher/biobank exemplar for SPECIMEN_IDx" and insert the SPECIMEN_ID for the specimen that is going for reference genome sequencing, potentially without its own DNA barcoding.

- If there is no relationship to note, this field can be left **blank**.

T. ==GAL==: (==ENA_submission==) Use the drop-down menu to select the Genome Accession Lab (GAL) responsible for this sample.

U. ==GAL_SAMPLE_ID==: (==ENA_submission==) This is the unique name assigned to the sample by the GAL. This will include an abbreviation for the GAL and a simple shorthand identifier. This is a free text field, but please do not use spaces or special characters, e.g., #, !, ^, *, etc. It is fine for the GAL_SAMPLE_ID to be the same as the COLLECTOR_SAMPLE_ID and the ERGA_SPECIMEN_ID if warranted. Each GAL maintains its own register of GAL_SPECIMEN_IDs for the project. Please ensure that you do not use IDs that have already been used, and that you stick to the format required by the GAL. Please see table on page 5 for details.

V. ==COLLECTOR_SAMPLE_ID==: This is the unique name assigned to the sample by the

COLLECTOR or COLLECTOR_AFFILIATION. This is a free text field, but please **do not use spaces or special characters**, other than hyphens and underscores ("-" and "_") i.e do not use #, !, ^, *, etc.

- In some cases, you will be splitting a single specimen across multiple tubes (see SPECIMEN_ID), and you will want to consider what kind of information you want in your unique sample names for this. For example, if the specimen is a butterfly with SPECIMEN_ID = Ox000005, and you put the head in one tube and the thorax in another, your COLLECTOR_SAMPLE_IDs might reflect this with one tube called Ox000005-h and the other called Ox000005-t. Likewise, the COLLECTOR_SAMPLE_ID may be the name given to a collection consisting of a 'clump' from a mat-forming species, which may then be subdivided into different specimen tubes, each given a unique SPECIMEN_ID.

W. <mark>COLLECTED_BY</mark>: (<mark>ENA_submission</mark>) Enter the name of the person or people who collected the sample using all CAPITALS, and separate names with "|" (vertical pipe symbol), e.g., "CAROLUS LINNAEUS | JEAN_BAPTISTE LAMARCK".

- We note that storage of names with affiliations in a database brings the ERGS system under the aegis of the GDPR regulations, and we must ask species ambassadors, GALs, and collaborators to agree to their data being stored in COPO and to those data being propagated to secondary databases (including ENA and the final collections of record). The species ambassador is asked to seek agreement from all involved collaborators before uploading the metadata sheet into COPO.

X. <mark>COLLECTOR_AFFILIATION</mark>: (<mark>ENA_submission</mark>) Free text field to supply the university, institution, or society that is responsible for the collected specimen. This is typically the society or institution of the person(s) specified in the COLLECTED_BY field. If multiple people are specified in COLLECTED_BY, ensure that their institutional affiliations are also separated by a vertical pipe symbol. Position in the list of affiliations should match the person in the same position in the list of names (e.g., PERSON A | PERSON X | PERSON C will have their affiliations as: (INSTITUTE A | INSTITUTE X | INSTITUTE C). If multiple people are listed but all from the same affiliation, no need to repeat the affiliation.

Y. <mark>COLLECTOR_ORCID_ID:</mark> Enter the 16 digits ORCID ID of the person or people who is responsible for the collection of the sample.

Z. <mark>DATE_OF_COLLECTION</mark>: (<mark>ENA_submission</mark>) The date the sample was collected, with year, month, and day specified (**YYYY-MM-DD**).

AA. <mark>COLLECTION_LOCATION</mark>: (<mark>ENA_submission</mark>) General description of the location. This should start with the geographical origin of the sample country as defined by the country or sea in agreement with ISNDC country list  (look up accepted country names     here

https://www.ebi.ac.uk/ena/browser/view/ERC000053), but also include more specific locations (e.g., "Barton's Pond") ranging from least to most specific and separated by | character, e.g., "United Kingdom | East Anglia | Norfolk | Norwich | University of East Anglia | UEA Broad". It is important to give the name of the site here if possible.

| | | |
|---|---|---|
| a | If the specimen is from a zoo, botanic garden, culture collection and has known origin elsewhere, please note this information in OTHER_INFORMATION and only include here information about the location of the specimen at the time from which a sample was taken | |
| (e.g., | "London Zoo", "Millennium Seed Bank", etc). | |

AB. **ORIGINAL_COLLECTION_DATE**: (ENA_submission) If the specimen is from a zoo, botanic garden, culture collection and has a known date of collection **from a known origin elsewhere** (e.g., the wild), please record the date here in as much detail as possible, with year, month and day specified (**YYYY-MM-DD**). YYYY-MM and YYYY is acceptable where further detail is not known. This information is important for regulatory compliance checks. Leave this field **blank** if it is not applicable.

AC. **ORIGINAL_GEOGRAPHIC_LOCATION**: (ENA_submission) If the specimen is from a zoo, botanic garden, culture collection and has a **known origin elsewhere**, please record the general description of the original location here. This should start with the country (United Kingdom, or look up other accepted country names here https://www.ebi.ac.uk/ena/browser/view/ERC000053), but also include more specific locations (e.g., "Barton's Pond") ranging from least to most specific and separated by | character, e.g., "United Kingdom | East Anglia | Norfolk | Norwich | University of East Anglia | UEA Broad". It is important to give the name of the site here if possible. This information is important for regulatory compliance checks. Leave this field **blank** if it is not applicable.

AD. DECIMAL_LATITUDE: (ENA_submission) In decimal degrees, between -90 and 90. We advise that locations are specified to a minimum of 3 decimal places (https://en.wikipedia.org/wiki/Decimal_degrees).

AE. DECIMAL_LONGITUDE: (ENA_submission) In decimal degrees, between -180 and 180. We advise that locations are specified to a minimum of 3 decimal places (https://en.wikipedia.org/wiki/Decimal_degrees).

AF. **GRID_REFERENCE**: Information to geolocate the sample area, preferably with a map or standardised geolocation reference, e.g., OS GRID REF: SP45998 08751. https://osmaps.ordnancesurvey.co.uk/ is useful to map lat-long to grid references. This field is optional and can be left **blank**.

AG. HABITAT: (ENA_submission) Any comments about the location, habitat or substrate, e.g., *damp mossy ground in moderate shade.* If substrate is living and there is a chance that it is included in the sample, add this to the SYMBIONT category. We recommend using terms from the ENVO ontology. If the specimen is from a zoo or botanic garden, you can add its original habitat to

"OTHER_INFORMATION" but here, please only capture its habitat at the time of collection (e.g., "reptile cage at London Zoo").

AH. **DEPTH**: (ENA_submission) Depth below (water body) surface, supplied in metres. This is not the absolute depth of the water body. Do not supply the unit, e.g., use 200 for 200 m below sea level, 100-200 for 100-200 m range below sea level, etc. Leave this field **blank** if the depth was not recorded or it is not an applicable field.

AI. **ELEVATION:** (ENA_submission) Altitude above sea level, supplied in metres. Do not supply the unit, e.g., use 200 for 200 m above sea level, 100- 200 for 100-200 m range above sea level, etc. Please supply elevation of water surface for inland water bodies. Leave this field **blank** if the elevation was not recorded or it is not an applicable field. For specimens isolated below the earth surface in sediment or soil, please provide a negative value in metres corresponding to sampling depth.

AJ. **TIME_OF_COLLECTION**: Time of day of sample collection in 24-hour clock format, with hours and minutes separated by colon e.g., 13:35, 04:53, etc. This should be in GMT/UTC. This field may be particularly relevant for RNAseq but it is not mandatory. Leave this field **blank** if the time was not recorded.

AK. **DESCRIPTION_OF_COLLECTION_METHOD**: A detailed as possible description of the sample collection methods, e.g., "*caught with fibre net within densely wooded area, and immediately placed into the collection container*".

AL. **DIFFICULT_OR_HIGH_PRIORITY_SAMPLE**: Drop down menu to flag species/samples that are difficult to collect (rare) or high priority to push through sequencing for any reason.

AM. **IDENTIFIED_BY**: (ENA_submission) Enter the name of the person or people who identified the sample to species level. Use ALL CAPs, and separate names with | (vertical pipe symbol), e.g., "CAROLUS LINNAEUS | JEAN-BAPTISTE LAMARCK".

We note that storage of names with affiliations in a database brings the ERGA system under the aegis of the GPDR regulations, and we must ask species ambassadors, GALs, and collaborators to agree to their data being stored in COPO and to those data being propagated to secondary databases (including ENA and the final collections of record). The species ambassador is asked to seek agreement from all involved collaborators before uploading the metadata sheet into COPO.

AN. **IDENTIFIER_AFFILIATION**: (ENA_submission) Free text field to supply the university, institution, or society that is responsible for the collected specimen. This is typically the society or institution of the person(s) specified in the IDENTIFIED_BY field. If multiple people are specified in IDENTIFIED_BY, ensure that their institutional affiliations are also separated by a vertical pipe symbol. Position in the list of affiliations should match the person in the same position in the list of names (e.g., "Person A | Person X | Person C" will have their affiliations as: "Institute A | Institute X | Institute C". If multiple people are listed but all from the same affiliation, no need to repeat the

affiliation.

AO. <mark>IDENTIFIED_HOW</mark>: Indicate what method(s) were used to identify the specimen to the nominal species (e.g., morphology, ITS barcoding). This is free text and should include reference to an authoritative key if possible. If the identification is by a taxon expert, note that here and ensure the name of that person is in the IDENTIFIED_BY column.

AP. <mark>SPECIMEN_ID_RISK</mark>: Y/N field to indicate if there is any risk that the SPECIMEN_ID provided does not reflect a single genetic entity OR the species names it has been submitted under. Examples of this include 1) a clump of tissue or cells that could comprise multiple individuals; 2) a species that is part of a species complex or group where it can be difficult to be certain of species identity. Please make every effort to ensure this field is N if possible (e.g., by taking single strands of clumpy organisms that are most likely to reflect a single genetic entity or ensuring molecular barcode data support the species name provided).

AQ. <mark>PRESERVED_BY</mark>: Name of person that carried out the preservation, supplied in CAPITALS. Multiple preserver names should be separated by a | character.

> We note that storage of names with affiliations in a database brings the ERGA system under the aegis of the GPDR regulations, and we must ask species ambassadors, GALs and collaborators to agree to their data being stored in COPO and to those data being propagated to secondary databases (including ENA and the final collections of records). The species' ambassador is asked to seek agreement from all involved collaborators before uploading the metadata sheet into COPO.

AR. <mark>PRESERVER_AFFILIATION</mark>: Free text field to supply the university, institution, or society that is responsible for the collected specimen. This is typically the society or institution of the person(s) specified in the PRESERVED_BY field. If multiple people are specified in PRESERVED_BY, ensure that their institutional affiliations are also separated by a vertical pipe symbol. Position in the list of affiliations should match the person in the same position in the list of names (e.g., Person A | Person X | Person C will have their affiliations as: (Institute A | Institute X | Institute C). If multiple people are listed but all from the same affiliation, there is no need to repeat the affiliation.

AS. <mark>PRESERVATION_APPROACH</mark>: e.g., snap frozen, dry ice, ethanol/dry ice slurry, in RNALater, lyophilised, air dried, etc.

AT. **PRESERVATIVE_SOLUTION**: Suspension liquid used to preserve the sample, e.g., RNALater, RLT Buffer, DESS. If no preservative was used, this field should be left **blank**.
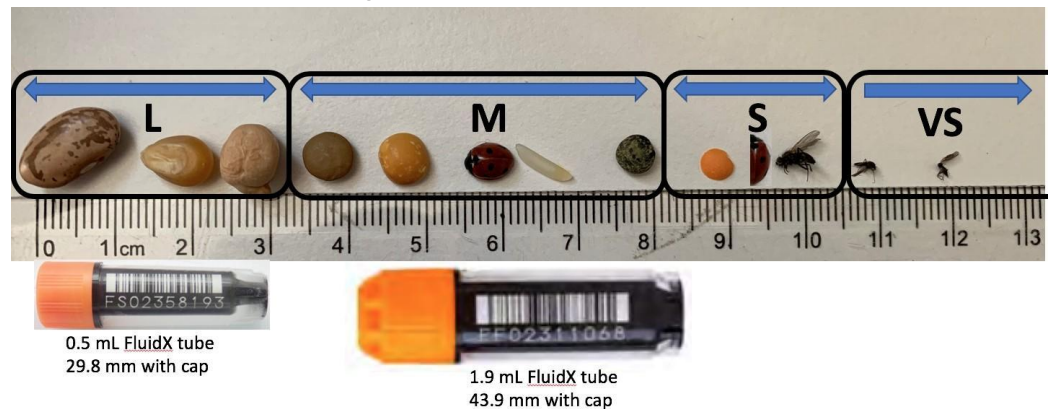
AU. <mark>TIME_ELAPSED_FROM_COLLECTION_TO_PRESERVATION</mark>: some organisms may be held living in collection for a period of time for starvation or other factors. This entry should be specified in hours, but no unit, e.g., 0.5 for half an hour, 3 for 3 hours,

etc.

AV. <mark>DATE_OF_PRESERVATION</mark>: Date on which the species was preserved. Please use **YYYY-MM-DD** format.

AW. <mark>SIZE_OF_TISSUE_IN_TUBE</mark>: Select from drop down menu how large is the sample in the tube. We aim for one lentil-sized piece per tube but sometimes adding more or less tissue than this will be necessary. Please note the approximate size of the piece or pellet: use the following shorthand:

- "VS" for very small

- "S" for small (~red lentil sized)

- "M" for medium (~yellow lentil/ladybird sized/5mm)

- "L" for large (>5mm, chickpea/bean sized)

- If the specimen is a single cell, use "SINGLE_CELL"

- Aim for single lentil sized (S or M) pieces in tubes whenever possible. If the sample is L, then wherever possible process this into multiple tubes of S or M sized pieces . See visual guidance below.

- If the sample has been shipped as extracted DNA please enter "NOT_APPLICABLE". Note that we expect that all samples will be extracted at Sanger.



0.5 mL FluidX tube
29.8 mm with cap

1.9 mL FluidX tube
43.9 mm with cap

Guidance for "Size_of_tissue_in_tube"
L = popcorn kernel or dried chickpea sized and larger
M = green, yellow lentil sized, whole ladybird size
S = red lentil, half a ladybird size
VS = smaller than half a red lentil
SC = single cell

AX. <mark>TISSUE_REMOVED_FOR_BARCODING</mark>: Select from drop down menu "**Y**" or "**N**". Instructions for appropriate Molecular Barcoding SOPs has to be arranged by the species ambassador with the Barcoding partner , noting that barcoding requires materials in specific tube or plate types.

AY. <mark>TUBE_OR_WELL_ID_FOR_BARCODING</mark>: This is either the well number on a plate (there are 96 wells per tissue plate) OR the barcode/unique identifier on the tube containing the tissue sample.

AZ. <mark>TISSUE_FOR_BARCODING</mark>: Please select from drop down menu what part of the organism was dissected for DNA barcoding (e.g. leg, soft-body tissue etc.). Muscle tissue is ideal for barcoding. This list is a repeat of the attributes available for "ORGANISM_PART" with one addition of "DNA_EXTRACT"

BA. <mark>BARCODE_PLATE_PRESERVATIVE</mark>: Typically, animal samples will be submerged in 70% ethanol, plant tissue will be preserved in silica gel, and fungal tissue will be frozen or lyophilized. Record the volume, concentration, and type of preservative/method of preservation used here.

BB. <mark>TISSUE_REMOVED_FOR_BIOBANKING</mark>: Select from drop down menu "**Y**" or "**N**". Instructions for appropriate Biobanking SOPs has to be arranged by the species ambassador with the Biobanking partner, noting that biobanking requires materials in specific tube or plate types.

BC. <mark>TISSUE_VOUCHER_ID_FOR_BIOBANKING</mark>: Accession number of frozen, biobanked material from the sequenced specimen. This ID should be prefixed by the name of the collection (e.g., ATCC:12345) and refers to a frozen, physical voucher of the specimen that is accessioned and curated into a collection accessible over GGBN (https://www.ggbn.org/ggbn_portal/) or the collection's webportal. This field can be updated in COPO at a later date if accession numbers are not available at the time of sample preparation (e.g., moth bodies sent to Sanger, moth wings to be accessioned and curated at Oxford Museum of Natural History and given a museum accession or acquisition number = Voucher_ID). In such cases please use **NOT_PROVIDED** as a placeholder, allowing for update at a later time.

> In some cases, voucher material will need to be made from a specimen that is different than the one being submitted for sequencing (e.g., a midge is too small to provide both a voucher and a specimen for sequencing, so another midge from the same swarm may provide a para-genomotype voucher). When this is the case, it should be noted.

BD. <mark>TISSUE_FOR_BIOBANKING</mark>: Pease select from drop down menu what part of the organism was dissected for biobanking (e.g. leg, soft-body tissue etc.). This list is a repeat of the attributes available for "ORGANISM_PART".

BE. <mark>DNA_REMOVED_FOR_BIOBANKING</mark>: Select from drop down menu "**Y**" (yes) or "**N**" (no).

BF. <mark>DNA_VOUCHER_ID_FOR_BIOBANKING</mark>: Accession number of DNA biobanked from the sequenced specimen. This ID should be prefixed by the name of the collection (e.g. ATCC:12345) and refers to a frozen sample of DNA of the specimen that is accessioned and curated into a collection accessible over GGBN (https://www.ggbn.org/ggbn_portal/) or the collection's webportal. This field can be

updated in COPO at a later date if accession numbers are not available at the time of sample preparation. In such cases please use **NOT_PROVIDED** as a placeholder, allowing for update at a later time.

BG. <mark>PURPOSE_OF_SPECIMEN</mark>: Please select appropriate from drop down menu.

- The majority of specimens will be for "REFERENCE_GENOME". All samples listed for REFERENCE_GENOME sequencing are assumed to also need DNA BARCODING and RNA-SEQUENCING, and the term "REFERENCE GENOME" encompasses all three things (reference genome, barcoding, rna-seq) wherever samples allow. Please use REFERENCE GENOME for all specimens / samples of a particular species unless they should be destined for an alternative use only.

- If a particular tissue is needed solely for RNAseq use "RNA-SEQUENCING"

- If the specimen is intended for population genetics or resequencing please use "SHORT_READ_SEQUENCING.

- The drop-down option for DNA_BARCODING_ONLY is reserved for those specimens submitted solely for DNA barcoding (e.g., when the sample is too small to provide material for both reference genome and barcoding and genome paratype / other specimens must be used as proxies, or when the specimen was identified to species level but died before being preserved, or is otherwise unsuitable for HMW DNA, but the material is valuable for barcoding).

- The drop-down option for PROXY_VOUCHERING_ONLY is reserved for those specimens submitted solely for physical specimen vouchering (collection, biobank) in case this is appropriate (e.g., sampled specimen at the same moment and time as a specimen sampled for genome sequencing that will be completely consumed during sequencing).

BH. <mark>HAZARD_GROUP</mark>: If the specimen needs to be processed in a containment level 1, 2, or 3 lab. Please note that any specimens above Hazard Group 1 must be discussed prior to shipping samples. To determine if the species is above HG1, please check both the HSE "Approved List of Biological Agents" and the SAPO list of animal pathogens. If the species is not listed on either of these lists, then it is HG1. Select from drop down menu.

BI. <mark>REGULATORY_COMPLIANCE</mark>: Please select from drop down menu Y (yes), NOT_APPLICABLE or N (not known). Note that ERGA will not be able to process further any samples where N is entered.

- Enter Y if you have affirmed that the necessary regulatory compliance documents have been obtained by the species ambassador and are available to the species ambassador and all involved partners including ther GAL. These documents need

to cover all regulatory compliance including sampling, vouchering, sample transfers, sequencing, and sequence deposition. These may include landowner permission, restricted area (SSSI, Nature Reserve, etc.) permission, BAP, CITES or other endangered species permission, ethical and Home Office Licencing for sampling for specified animals (vertebrates, cephalopods), phytosanitary permissions, veterinary pathogen sampling permissions etc.

- ■ If you have determined that no regulatory permissions or documents are required (for example where the sample is from a long-established culture) please enter NOT_APPLICABLE.

- ■ This is an important "per species" check that ensures that permissions were granted to collect and transfer the specimen for this research purpose. The sample provider should ensure this documentation is obtained, and that copies of the relevant paperwork are shared with the sequencing institution where necessary and as stipulated, for example, by regulations/approvals or licencing authorities..

BJ. **VOUCHER_ID**: (**ENA_submission**) Accession number of voucher material from the sequenced specimen. This ID should be prefixed by the name of the collection (e.g., ATCC:12345) and refers to the physical voucher of the specimen that is accessioned and curated into a collection. This field can be updated in COPO at a later date if accession numbers are not available at the time of sample preparation. In such cases please use **NOT_PROVIDED** as a placeholder, allowing for update at a later time.

- ■ In some cases, voucher material will need to be made from a specimen that is different than the one being submitted for sequencing ("proxy vouchering", e.g., a midge is too small to provide both a voucher and a specimen for sequencing, so another midge from the same swarm may provide a para-genomotype voucher). When this is the case, it should be noted and the vouchered specimen be entered with an independent ERGA_SPECIMEN_ID and relationship be indicated over column RELATIONSHIP

BK. **INDIGENOUS_RIGHTS_APPLICABLE:** Mandatory information upon if indigenious rights are applicable to the sample/the species the sample was derived from, select "**Y**" (yes) or "**N**" (no) from drop down menu.

BL. <mark>INDIGENOUS_RIGHTS_DEF</mark>:  Free text, Please state which rights are applicable if the previous column says yes, else NA.

> -> upload field pdf

BM. <mark>ASSOCIATED_TRADITIONAL_KNOWLEDGE_APPLICABLE:</mark>  Mandatory information upon if associated traditional knowledge exists for the sample/the species the sample was derived from, select "**Y**" (yes) or "**N**" (no) from drop down menu.

BN. <mark>ASSOCIATED_TRADITIONAL_KNOWLEDGE_LABEL:</mark>

> Choose label for corresponding icon from drop down menu according to icons used in https://localcontexts.org/labels/traditional-knowledge-labels, https://localcontexts.org/labels/biocultural-labels/ [localcontexts.org]

BO. <mark>ASSOCIATED_TRADITIONAL_KNOWLEDGE_CONTACT:</mark> Provide reference, free text allowed, could be linked to an ORCID

BP. <mark>ETHICS_PERMITS_MANDATORY:</mark>  Mandatory information upon if an ethics permit is needed to sample/sequence/voucher/biobank the sample/the species the sample was derived from, select "**Y**" (yes) or "**N**" (no) from drop down menu.

BQ. <mark>ETHICS_PERMITS_DEF:</mark>  Free text explaining permits, permit issuing entity and permit number. -> upload field pdf

BR. <mark>SAMPLING_PERMITS_MANDATORY:</mark>  Mandatory information upon if sampling permits are needed to sample/sequence/voucher/biobank the sample/the species the sample was derived from, select "**Y**" (yes) or "**N**" (no) from drop down menu.

BS. <mark>SAMPLING_PERMITS_DEF</mark>:  Free text explaining permits, permit issuing entity and permit number. -> upload field pdf

BT. <mark>NAGOYA_PERMITS_MANDATORY:</mark>  Mandatory information upon if a permit in compliance with the *Nagoya Protocol on Access to Genetic Resources and the Fair and Equitable Sharing of Benefits Arising from their Utilization to the Convention on Biological Diversity* is needed for the sample in question/the species the sample was derived from, Select "**Y**" (yes) or "**N**" (no) from drop down menu.

BU. <mark>NAGOYA_PERMITS_DEF:</mark>  Free text explaining permits, permit issuing entity and permit number. -> upload field pdf

BV. **OTHER_INFORMATION**: Free text field for further relevant information not captured by the other fields. This is a place also for partners to flag species that should be prioritized in the sequencing queue. If this species represents one of the two family representatives submitted for the project, please note this here. If there is nothing else to add here, this field should be left **blank**.