



# **Deliverable D7.4**

# HoloRuminant Data Management Plan (update at M18)

Due date: M18

Actual submission date: M19

Start date of the project: October 1<sup>st</sup>, 2021 Duration: 60 months

Workpackage: WP7

Workpackage leader: INRAE Deliverable leader: INRAE

Partners involved: INRAE, QUB, UG, LUKE

Version: V1

Dissemination Level				
PU Public	✓			
CI Classified, as referred to Commission Decision 2001/844/EC				
<b>CO</b> Confidential, only for members of the consortium (including the Commission Services)				





# Table of contents

1.	Summary	3
2.	Introduction	3
3.	Methods and Datasets identified	3
4.	Conclusions	5
5.	Annexes	6





## 1. Summary

This document describes the second version of the HoloRuminant Data Management Plan (DMP). The first part of this report will describe the datasets categories identified in the initial workplan. The second part of this report will indicate the rules implemented by Holoruminant to comply with the Open Research Data policy in the context of EU-funded projects. This is living document that is going to be updated in the coming months and years with a finer description of datasets and of their metadata.

The objectives of this Data Management Plan are to:

- i. Describe the nature of data generated and collected by the project,
- ii. define the data sharing strategy based on the FAIR principles (data should be findable, accessible, interoperable and re-usable),
- iii. Identify the data management tools (data archiving and preservation),
- iv. Ensure that data management policy is in line with the project consortium agreement.

Work package	WP7		
Work package leader	Diego Morgavi	INRAE	
Author(s)	Chris Creevey/Diego Morgavi	QUB/INRAE	
Additional Authors	Dorte Becher	UG	
	Jarkko Niemi	LUKE	
Version	V1		

# 2. Introduction

HoloRuminant partners participate in the Pilot on Open Research Data in Horizon 2020, which aims to improve and maximise access to and re-use of research data generated by actions. However, participation in the Pilot is flexible in the sense that it does not mean that all research data needs to be open. Thereafter, participants have formulated this Data Management Plan (DMP), which addresses the relevant aspects of making data FAIR – findable, accessible, interoperable and re-usable, including what data the project will use or generate, whether and how it will be made accessible for verification and re-use, and how it will be curated and preserved. Through this DMP, HoloRuminant can define certain datasets to remain closed according to the principle "as open as possible, as closed as necessary".

# 3. Methods and Datasets identified

We adopted <u>ARGOS</u> as the tool to create the DMP and we are using the templates and Guidelines on Data Management in Horizon 2020. The datasets identified from the HoloRuminant work plan and covered by this first version of the DMP are DNA Sequencing data, Proteomic data, and Social Sciences data. Additional datasets might be identified during the course of the project and would therefore be included in an updated version.

The DMP in ARGOS can be seen in Annex 1 below. The ARGOS' template provides all the information necessary about the research and how the data complies with FAIR principles. We emphasize that data used





in this project will be managed using secure storage rules set up by each responsible partner producing or storing the data. All data will be included in a backup plan (local and remote backups on servers secured by access rights). The project partners will comply with GDPR as for the secure storage and protection of personal data (see D8.2 "POPD - Requirement No. 2").

Partners will follow the Ethics requirements governing the collection, storage and security of personal data for complying with GDPR for the dataset concerned. Animal studies in Partners' research facilities will follow the European Convention of the Protection of Vertebrate Animals used for Experimental and Scientific Purposes, directive 2010/63/EU and the 3R principles (Replacement, Reduction and Refinement). Details on the planned animal studies and certifications are described in D8.3 "A - Requirement No. 3".

The HoloRuminant project is not on the scope of the Nagoya Protocol (which is an international agreement which aims at sharing the benefits arising from the utilization of genetic resources in a fair and equitable way). The assessment was done after an initial consultation of available information on the <u>Nagoya Protocol</u> <u>documentation</u> website, contact with partners and ABS-Focal Points from concerned countries.

The HoloRuminant project will follow as much as possible the FAIR data principle, making data Findable, Accessible, Interoperable and Re-usable. To achieve this goal, is has been agreed within the consortium that the datasets that will be shared in open access will be associated with a version number and keywords (see Annex 2 "Rules for naming files").

To exchange data within a given WP and between WPs (inputs/outputs), all partners will use the same standards. To achieve these goals, INRAE has developed the <u>Animal Trait Ontology of Livestock</u> (ATOL). The ATOL is an ontology of characteristics defining phenotypes of livestock in their environment. The <u>EOL</u> (Environment Ontology for Livestock) ontology describes environmental conditions of livestock farms.

Data will then be made accessible by using as much as possible public repositories and domain-specific databases as described in the ARGOS document.

Bioinformatic tools developed and assembled during the project that cannot be deposited in domain-specific repositories will be stored at INRAE, that will be in charge of providing open access to user as well as the quality control, curation and preservation of the resources for at least 10 years.

All acquired data will be stored securely and curated by the academic partners' institutes, e.g. Data INRAE for INRAE (https://data.inra.fr/dataverse/root). Final results will be stored in partners' results repositories but also in publicly available standard data repositories (Zenodo, HAL, etc.) for at least ten years following the end of the project. In addition, partners will be encouraged to publish research data as supporting material with their publications to facilitate preservation of data for future re-use by other projects or research initiatives, through data paper journals (e.g. Data in Brief) providing open-access publication of datasets with DOI.

Project partners will preserve raw data and will be strongly encouraged to attribute DOI (Digital Object Identifier) to specific datasets for increased visibility. As for the bioinformatic tools mentioned above, data not adapted to domain-specific repositories (e.g. EMBL-ENI, NCBI) will be deposited in an open repository such as the INRAE <u>data repository</u> which allows identification of data sets with DOI, partners internal repositories and/or other publicly available standard data repositories (e.g. Zenodo, HAL, etc.).





Following the FAIR guidelines: Making data openly accessible, interoperable and reusable; the HoloRuminant partners will be strongly encouraged to make their data publicly accessible unless impossible for legal constraints.

To increase reusability, the consortium will standardise file names (Annex 2) and will notably use existing standards and ontologies to name and describe variables as mentioned above. The data in open access will be interoperable by using at maximum the ATOL, EOL and AHOL ontologies whenever possible. All parameters used will be in English language, and all units will follow international standard unit.

# 4. Conclusions

The Data Management Plan (DMP) of HoloRuminant is in accordance with EU Regulation and policies of research organisations which are partners of the project. Consequently, HoloRuminant partners will rely on this DMP to follow the open access rules for research data. Partners will follow this detailed DMP for making HoloRuminant data Findable, Accessible, Interoperable and Reusable (according to the FAIR principles) following guidelines published by Wilkinson et al. (2016; <u>The FAIR Guiding Principles for scientific data</u> management and stewardship) and the living document maintained by the <u>GO FAIR Initiative</u>

New data will be made available in a FAIR and Open manner, complying with EC Open Data policies and policies of research organisations which are partners of the project.

This first version of the DMP will be updated regularly. This DMP will be a "living" document outlining how the research data collected and generated (including computationally) will be handled during and after the project.





#### 5. Annexes

# Annex 1

DPM Created from Argos

# Data Management Plan Information

## 1. HoloRuminant DMP

Data Management Plan for the H2020 project HoloRuminant. The purpose is to describe the different datasets produced in the project and, according to FAIR principles, how data is formatted, shared, re-used, preserved and archived.

#### Funder

European Commission | | EC

#### Grant

Understanding microbiomes of the ruminant holobiont

#### Organisations

University of Alberta, Faculty of Science, Queen's University Belfast, Ben Gurion University of the Negev, THE AGRICULTURAL RESEARCH ORGANISATION OF ISRAEL - THE VOLCANI CENTRE, CSIRO Australia, Aarhus Universitet, IRTA, TEAGASC, Agencia Estatal Consejo Superior de Investigaciones Científicas, Natural Resources Institute Finland (LUKE), Wageningen University & Research, Institut national de recherche pour l'agriculture, l'alimentation et l'environnement (INRAE), NMBU - Norwegian University of Life Sciences, Agri Food and Biosciences Institute, National University of Ireland, Galway, Greifswald University, Ghent University, Scotlands Rural College (SRUC) SRUC Research Portal, AgResearch Limited, THE REGENTS OF THE UNIVERSITY OF CALIFOR

#### Researchers

Diego Morgavi (orcid:0000-0002-3883-0937), Jarkko Niemi (orcid:0000-0002-9545-3509), Chris Creevey (orcid:0000-0001-7183-1555), Dörte Becher (orcid:0000-0002-9630-5735)





# Datasets

#### 2. Title: Social Sciences data for HoloRuminant project

#### Template: Horizon 2020

The dataset will include qualitative and quantitative information on innovations which are studied by the HoloRuminant project. This dataset is collected to facilitate stakeholder engagement and to consult stakeholders and consumers about these innovations. The dataset will be used to determine their expectations and views concerning these interventions.

#### **Dataset Description**

#### 1.1 Data Summary

1.1.1 What is the purpose of the data collection/generation and its relation to the objectives of the project?

To obtain information, To share information, To keep on record, To make informed decisions, To improve a product

Comment: The dataset is collected to facilitate stakeholder engagement, to consult stakeholders and consumers about their view concerning innovations developed by HoloRuminant, and to help to improve and adopt the innovations.

#### 1.1.2 What types of data will the project generate/collect?

• observational (e.g., sensor data, data from surveys), simulation (e.g., climate modeling data), reference or canonical (e.g., static, peer-reviewed data sets, likely published or curated, such as gene sequence databanks or chemical structures)

- Surveys data, focus groups, review data, modelled data
- 1.1.3 What formats of data will the project generate/collect?
  - Text files, Numerical, Multimedia, Models
  - Scripts/summaries/recordings of focus group discussions, survey data in tables, model results

# 1.1.4 What is the origin of the data?

Primary data

1.1.5 What is the expected size of the data? **GB (gigabyte)** 

*Comment: Megabytes; Audio-visual recordings' total size may be gigabytes (depending on the format)* 

#### 1.1.6 To whom might it be useful ('data utility')?

- Researchers, Research communities, Decision makers, Industry
- Reuse and analysis by researchers. The use of confidential data are restricted..



#### 2.1 Reused Data

2.1.1 Will you re-use any existing data and how? No

3.1.1 Making data findable, including provisions for metadata

3.1.1.1 Will you use metadata to describe the data?

• Yes

#### • CERIF (Common European Research Information Format)

- 3.1.1.3 Will your metadata use standardised vocabularies?
  - Yes
  - AGROVOC Multilingual Agricultural Thesaurus

3.1.1.5 Will you make the metadata available free-of-charge? Yes

3.1.1.6 Will your metadata be harvestable? Yes

#### Comment: When permitted by data confidentiality and personal data protection reasons.

3.1.1.7 Will you use naming conventions for your data? Yes

3.1.1.9 Will you provide clear version numbers for your data? Yes

#### Comment: Running numbers to identify versions

3.1.1.10 Will you provide persistent identifiers for your data? Yes

# 3.1.1.11 Persistent identifiers **URN**

3.1.1.12 Will you provide searchable metadata for your data? No

3.1.1.15 Will you use standardised formats for your data?

• Yes

#### • Microsoft Excel for Windows

3.1.1.18 Are the file formats you will use open? Yes

3.1.1.20 Do supported open-source tools exist for accessing the data? Yes

3.1.1.21 Please describe if data require proprietary tools to access the data? Open access programs to read and access the data exist.





3.1.1.22 Will you provide metadata describing the quality of the data? Yes

#### Comment: Brief data description will be provided.

3.1.2 Making data openly accessible

3.1.2.1 Are there ethical or legal issues that can impact sharing the data? Yes

Comment: Some data may not be shared because of data protection and informed consent reasons.

3.1.2.2 Will your data be openly accessible? some

*Comment: upon publication of associated papers, if permitted by data confidentiality and personal data protection reasons.* 

3.1.2.3 Please provide the URL of the data which can be made available **not available** 

- 3.1.2.4 How will the data be made available?
  - University repository
  - RADAR Luke

3.1.2.6 Is the storage sufficiently secure for the data and does the storage provide backup and recovery procedures?

secure with backup and recovery

3.1.2.7 Are there any methods or tools required to access the data? Yes

3.1.2.8 Please provide information about the method(s) needed to access the data **Data requests: kirjaamo@luke.fi** 

3.1.2.9 Please provide information about the tools needed to access the data. Data requests: kirjaamo@luke.fi

3.1.2.10 Will you also make auxiliary data that may be of interest to researchers available? after publication

Comment: Auxiliary data will be provided in the publications

#### 3.1.3 Making data interoperable

3.1.3.1 Will you use a controlled vocabulary for your data? Yes

#### 3.1.4 Increase data reuse

3.1.4.1 When do you plan to make your data available for reuse? after article publication





# 3.1.4.4 What internationally recognised licence will you use for your data? Creative Commons Attribution 4.0 International

3.1.4.5 Do you have documented procedures for quality assurance of your data? Yes

Comment: To be provided in the data collection protocols developed during the project.

3.1.4.7 Describe the data quality assurance processes Data conform to format specification, Other

Comment: Quality checks by data collecting persons.

3.1.4.8 Will you provide any support for data reuse? Yes

3.1.4.9 How long do you intend to support data reuse? Up to 5 years

#### 4.1 Allocation of resources

4.1.1 How will the cost of making your data findable, accessible, interoperable and reusable be covered? Use of institution infrastructure

4.1.2 Will you identify a data manager to manage your data, if not who will be responsible for the management of your data?

Yes

4.1.3 Identify the people or roles that will be responsible for the management of the project data Jarkko Niemi (orcid:0000-0002-9545-3509)

4.1.4 How do you intend to ensure data reuse after your project finishes? Institutional archive

#### 5.1 Data Security

5.1.1 What do you plan to do with research data of limited use? Kept on secure, managed storage for limited time

#### 6.1 Ethical aspects

6.1.1 Are there any ethical or legal issues that can have an impact on data sharing? Yes

Comment: Protection of personal data, consent given by those who provide the data for the project

6.1.2 What are the methods used for processing sensitive/personal data?

Anonymising data where necessary, Privacy constraints and applicable ethical norms, Data accompanied by informed consent statements, Privacy policies, National laws

#### 7.1 Other

7.1.1 Do you make use of other procedures for data management? No





#### 3. Title: Proteomics dataset

#### Template: Horizon 2020

#### Obtaining proteomic information from microbes/microbiomes associated to ruminants.

#### **Dataset Description**

#### 1.1 Data Summary

1.1.1 What is the purpose of the data collection/generation and its relation to the objectives of the project?

To obtain information, To share information, To keep on record, To make informed decisions, To combine with other data

Comment: (Meta)proteome data generation, storage and analysis and combinations with other data for the purpose of understanding the ruminant microbiome

1.1.2 What types of data will the project generate/collect?

• sample or specimen data, experimental (e.g., gene sequencing data)

• Protein data

1.1.3 What formats of data will the project generate/collect? Text files, Numerical, Discipline specific formats, Instrument specific formats

1.1.4 What is the origin of the data? Primary data

1.1.5 What is the expected size of the data? **TB (terabyte)** 

1.1.6 To whom might it be useful ('data utility')?

- Researchers, Research communities
- Reuse and re-analysis by researchers in the community.

#### 2.1 Reused Data

2.1.1 Will you re-use any existing data and how?

No

- 3.1.1 Making data findable, including provisions for metadata
  - 3.1.1.1 Will you use metadata to describe the data?
    - Yes
    - MIBBI (Minimum Information for Biological and Biomedical Investigations)
  - 3.1.1.3 Will your metadata use standardised vocabularies?
    - Yes
    - Couldn't find it? Insert it manually



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement N° 101000213.



#### 3.1.1.4 Please provide URL/Description of used vocabularies

• https://sicpa-web.cati.inrae.fr/ontologies/visualisation/ontologie/

#### • Animal Trait Ontology for Livestock

3.1.1.5 Will you make the metadata available free-of-charge? Yes

3.1.1.6 Will your metadata be harvestable? Yes

3.1.1.7 Will you use naming conventions for your data? Yes

3.1.1.9 Will you provide clear version numbers for your data? Yes

3.1.1.10 Will you provide persistent identifiers for your data? Yes

3.1.1.11 Persistent identifiers **DOI** 

3.1.1.12 Will you provide searchable metadata for your data? No

3.1.1.15 Will you use standardised formats for your data?

• Yes

• Couldn't find it? Insert it manually

3.1.1.16 Provide information about used standardised formats

- https://en.wikipedia.org/wiki/Mass\_spectrometry\_data\_format
- e.g. MZxml, .raw

3.1.1.18 Are the file formats you will use open? Yes

Comment: partially, as raw-files in specific formats of manufacturers of mass spectrometers will be

accessible as well

3.1.1.20 Do supported open-source tools exist for accessing the data? Yes

Comment: e.g. Trans-Proteomic Pipeline

3.1.1.22 Will you provide metadata describing the quality of the data? Yes

Comment: Quality data are included in the data





#### 3.1.2 Making data openly accessible

3.1.2.1 Are there ethical or legal issues that can impact sharing the data? No

3.1.2.2 Will your data be openly accessible? all

Comment: Comment: upon publication of associated papers

3.1.2.4 How will the data be made available?

- Repository of Archive
- Couldn't find it? Insert it manually

3.1.2.5 Please provide URL/Name of used data repositories https://www.ebi.ac.uk/pride/archive/

3.1.2.6 Is the storage sufficiently secure for the data and does the storage provide backup and recovery procedures?

secure with backup and recovery

3.1.2.7 Are there any methods or tools required to access the data?

• No

• All data is freely available over the web using multiple different access methods

3.1.2.10 Will you also make auxiliary data that may be of interest to researchers available? after publication

Comment: Auxiliary data will be provided with publications describing the data.

#### 3.1.3 Making data interoperable

3.1.3.1 Will you use a controlled vocabulary for your data? Yes

3.1.4 Increase data reuse

3.1.4.1 When do you plan to make your data available for reuse? after article publication

3.1.4.5 Do you have documented procedures for quality assurance of your data? Yes

Comment: For proteome data, quality will be assessed by the discipline standard approaches.

3.1.4.7 Describe the data quality assurance processes Use of tools for automatic checks

3.1.4.8 Will you provide any support for data reuse? No





#### 4.1 Allocation of resources

4.1.1 How will the cost of making your data findable, accessible, interoperable and reusable be covered? Use of institution infrastructure

4.1.2 Will you identify a data manager to manage your data, if not who will be responsible for the management of your data?

Yes

4.1.3 Identify the people or roles that will be responsible for the management of the project data

a.

- Dörte Becher (orcid:0000-0002-9630-5735)
- Data management of metaproteome data.

b.

- Diego Morgavi (orcid:0000-0002-3883-0937)
- Project coordinator

4.1.4 How do you intend to ensure data reuse after your project finishes? Institutional archive

Comment: all data are stored with back up for at least 20 years

#### 5.1 Data Security

5.1.1 What do you plan to do with research data of limited use? Kept on secure, managed storage for limited time

6.1 Ethical aspects

6.1.1 Are there any ethical or legal issues that can have an impact on data sharing? No

6.1.2 What are the methods used for processing sensitive/personal data? National laws

#### 7.1 Other

7.1.1 Do you make use of other procedures for data management? No

4. Title: DNA Sequencing data Template: Horizon 2020 Metagenomic sequencing





#### **Dataset Description**

#### 1.1 Data Summary

1.1.1 What is the purpose of the data collection/generation and its relation to the objectives of the project?

To obtain information, To share information, To keep on record, To make informed decisions, To combine with other data

*Comment: DNA sequencing generation, storage and analysis and combinations with other data for the purpose of understanding the ruminant microbiome* 

#### 1.1.2 What types of data will the project generate/collect?

- experimental (e.g., gene sequencing data)
- Gene sequencing data
- 1.1.3 What formats of data will the project generate/collect?
  - Text files, Numerical, Discipline specific formats, Instrument specific formats
  - Fastq sequencing data files, with associated summaries etc.

1.1.4 What is the origin of the data? Primary data

1.1.5 What is the expected size of the data? **TB (terabyte)** 

#### 1.1.6 To whom might it be useful ('data utility')?

- Researchers, Research communities
- Reuse and re-analysis by researchers in the community.

#### 2.1 Reused Data

2.1.1 Will you re-use any existing data and how?

#### No

- 3.1.1 Making data findable, including provisions for metadata
  - 3.1.1.1 Will you use metadata to describe the data?
    - Yes
    - MIBBI (Minimum Information for Biological and Biomedical Investigations)
  - 3.1.1.3 Will your metadata use standardised vocabularies?
    - Yes
    - Couldn't find it? Insert it manually

#### 3.1.1.4 Please provide URL/Description of used vocabularies

- https://sicpa-web.cati.inrae.fr/ontologies/visualisation/ontologie/
- Animal Trait Ontology for Livestock



### HoloRuminant - H2020 n°101000213



3.1.1.5 Will you make the metadata available free-of-charge? Yes

3.1.1.6 Will your metadata be harvestable? Yes

3.1.1.7 Will you use naming conventions for your data?

Yes

3.1.1.9 Will you provide clear version numbers for your data? Yes

3.1.1.10 Will you provide persistent identifiers for your data? Yes

3.1.1.11 Persistent identifiers

a.

DOI

b.

- Other
- NCBI Bioproject ID + sample ID

3.1.1.12 Will you provide searchable metadata for your data? No

3.1.1.15 Will you use standardised formats for your data?

- Yes
- Couldn't find it? Insert it manually

3.1.1.16 Provide information about used standardised formats

- https://en.wikipedia.org/wiki/FASTQ\_format
- FASTQ

3.1.1.18 Are the file formats you will use open? Yes

3.1.1.20 Do supported open-source tools exist for accessing the data? Yes

3.1.1.22 Will you provide metadata describing the quality of the data? Yes

Comment: Quality data is included in the data format provided (fastq)

#### 3.1.2 Making data openly accessible

3.1.2.1 Are there ethical or legal issues that can impact sharing the data? No



#### 3.1.2.2 Will your data be openly accessible?

all

#### Comment: upon publication of associated papers

- 3.1.2.4 How will the data be made available?
  - Domain-specific database
  - NCBI

3.1.2.6 Is the storage sufficiently secure for the data and does the storage provide backup and recovery procedures?

#### secure with backup and recovery

3.1.2.7 Are there any methods or tools required to access the data?

- No
- All data is freely available over the web using multiple different access methods (HTTP/FTP/etc)

3.1.2.10 Will you also make auxiliary data that may be of interest to researchers available? after publication

Comment: Auxiliary data consisting of animal phenotypic data will be provided with publications

#### describing the data.

#### 3.1.3 Making data interoperable

3.1.3.1 Will you use a controlled vocabulary for your data? Yes

#### 3.1.4 Increase data reuse

3.1.4.1 When do you plan to make your data available for reuse? after article publication

3.1.4.5 Do you have documented procedures for quality assurance of your data? Yes

Comment: For sequencing data, quality will be assessed by the discipline standard approaches such as

FASTQC etc....

3.1.4.6 Please provide URL with the documented procedures https://www.bioinformatics.babraham.ac.uk/projects/fastqc/

3.1.4.7 Describe the data quality assurance processes Use of tools for automatic checks, Data conform to format specification

3.1.4.8 Will you provide any support for data reuse? No

#### 4.1 Allocation of resources

4.1.1 How will the cost of making your data findable, accessible, interoperable and reusable be covered? Use of institution infrastructure





Comment: INRAE facilities for storage and sharing of data within the project prior to deposition in public

repositories.

4.1.2 Will you identify a data manager to manage your data, if not who will be responsible for the management of your data?

Yes

4.1.3 Identify the people or roles that will be responsible for the management of the project data

- a.
  - Chris Creevey (orcid:0000-0001-7183-1555)
  - Data management of metagenomic sequencing data.
- b.
- Diego Morgavi (orcid:0000-0002-3883-0937)
- Project coordinator

4.1.4 How do you intend to ensure data reuse after your project finishes? Data Center Archive Storage

Comment: NCBI genbank/SRA/etc.

#### 5.1 Data Security

5.1.1 What do you plan to do with research data of limited use?

- Delete at end of project
- Intermediate data files generated as part of the analyses will be kept until publication and deleted

at the end of the project.

#### 6.1 Ethical aspects

6.1.1 Are there any ethical or legal issues that can have an impact on data sharing?

No

6.1.2 What are the methods used for processing sensitive/personal data? National laws

Comment: GPDR rules will apply if required (there is no plan for the use or collection of

sensitive/personal data)

7.1 Other

7.1.1 Do you make use of other procedures for data management? No





# Annex 2

# **Rules for naming files for HoloRuminant**

Common to all members of the consortium and for files included in the DMP

#### File names will be:

- Descriptive and provide just enough contextual information.
- Short, 40-50 character Operating systems have different limits to the number of characters.
- Unique.

Do not create case-sensitive files, e.g. Data.csv and DATA.csv.

- Will use ISO 8601 standard for date (YYYYMMDD).
- If using a sequential numbering system, use leading zeros to make sure files sort in sequential order. Example: 001, 002, ...010, 011 ... 100, 101 ...
- Use versioning to indicate the most current version of a file. Example: filename\_v02.xxx
  - Use one leading zero for version (see point above)
- Avoid special characters, such as: ~ ! @ # \$ % ^ & \* () `; : <> ? . , [] { } ' " |
- Do not use spaces as some software will not recognize file names with spaces.
- Use underscores for separating sections
  - Use dashes, first letter capitalized for other information
- Avoid empty words like: the, of, and, or, ...
- The name of the file must inform on the type of data contained in the document according to a standard list of abbreviations that are reported in the README.txt (see below).
  - $\circ~$  PR for Protocol;
  - RD for Raw Data;
  - $\,\circ\,$  AD for Analyzed Data;
  - LIST for lists (e.g. sample list);
  - $\,\circ\,$  SCRIPT script or program ;
  - $\circ~$  Notes for Study Notes

File with information about the data, e.g. sick animal number, antibiotic treatment, etc.

• FILE NAMES SHOULD BE MACHINE READABLE, HUMAN READABLE, AND PLAY WELL WITH DEFAULT ORDERING

# File Name Convention for HoloRuminant

[project-name]\_[WP#]\_[partner-abbreviation]\_[data-type]\_[subject-or-keywords\_or\_study-ID]\_[YYYYMMDD]\_[version].[ext]

➔ Example

HR\_1\_QUB\_RD\_metaG-calves\_20221202\_v01.fastq



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement N° 101000213.

limit.



- The naming conventions is documented in a README.txt file so that all partners can follow this standard (example in Appendix 1).
- Datasets will have an associated metadata that will have a persistent identifier that can be a Handle or DOI in a dataverse (<u>https://data.inrae.fr/</u>) (example in Appendix 2 from <u>datainrae guide</u>).

#### References

- 1. SOP AQ1204 V1 from UMR Herbivores and references therein
  - a. DataPartage, https://www6.inra.fr/datapartage/Gerer/Nommer-et-organiser-ses-fichiersde-donnees
  - b. <u>http://datalib.edina.ac.uk/mantra/organisingdata/</u>
  - c. <u>http://libraries.mit.edu/data-management/store/organize/</u>
  - d. <u>https://www.jisc.ac.uk/guides/managing-information/good-file-name</u>
  - e. K. McNeill, 2014. Research Data Management: 101 The Lifecycle of a Dataset. MIT library
  - f. K. McNeill, H. Bailey, 2014. Research Data Management: File Organization. MIT library
  - g. UK data service, 2019, https://www.ukdataservice.ac.uk/manage-data/format/organising
- 2. Harvard Biomedical Data Management
  - a. https://osf.io/dpu45/





# **APPENDIX 1**

#### **README.txt** file

-----

README, File & Folder Schema Template V01 Last Updated: YYYYMMDD by XXXX

File naming schema:

File type: Protocol Filename schema: keywords]\_[YYYYMMDD]\_[version].[ext] Schema key: PR for protocol Example:

File type: Script or program Filename schema: keywords]\_[YYYYMMDD]\_[version].[ext] Schema key: SCRIPT for script or program Example:

File type: Lists Filename schema: keywords]\_[YYYYMMDD]\_[version].[ext] Schema key: LIST for lists (sample list for example) Example:

File type: Raw data Filename schema: keywords]\_[YYYYMMDD]\_[version].[ext] Schema key: RD for Raw data Example:

File type: Analysed data Filename schema: keywords]\_[YYYYMMDD]\_[version].[ext] Schema key: AD for analysed data Example:

File type: Research Notes Filename schema: keywords]\_[YYYYMMDD]\_[version].[ext] Schema key: Notes for research notes Example: HR\_[WP#]\_[partner-abbreviation]\_PR\_[subject-or-

HR\_[WP#]\_[partner-abbreviation]\_SCRIPT\_[subject-or-

HR\_[WP#]\_[partner-abbreviation]\_LIST\_[subject-or-

HR\_[WP#]\_[partner-abbreviation]\_RD\_[subject-or-

HR\_[WP#]\_[partner-abbreviation]\_AD\_[subject-or-

HR\_[WP#]\_[partner-abbreviation]\_Notes\_[subject-or-





File type: Filename schema: Schema key: Example:





# Appendix 2

Metadata controlled by a list of values - User Guide

#### Common Metadata Author / Identifier Scheme

- DAI
- DOI
- GND
- idHAL
- ISNI
- LCNA
- ORCID
- ResearcherID
- ScopusID
- VIAF

Note : ORCID recommended

#### Kind of Data :

- Audiovisual
- Collection
- Dataset
- Event
- Image
- Interactive Resource
- Model
- Physical Object
- Service
- Software
- Sound
- Text
- Workflow
- Other

#### **Data Origin**

- observational data
- experimental data
- survey data
- analysis data
- text corpus
- simulation data
- aggregate data
- audiovisual corpus
- computer code
- Other





#### Subject : based on classification OCDE

- Agricultural Sciences
- Arts and Humanities
- Astronomy and Astrophysics
- Business and Management
- Chemistry
- Computer and Information Science
- Earth and Environmental Sciences
- Engineering
- Law
- Mathematical Sciences
- Medicine, Health and Life Sciences
- Physics
- Social Sciences
- Other

#### Life cycle step: DDI Controlled Vocabulary for Lifecycle Event Type

- Study proposal
- Funding
- Study design
- Instrument design
- Questionnaire translation
- Questionnaire adaptation
- Interviewer training
- Ethics review
- Legal review
- Sampling
- Instrument pre-testing
- Pilot study
- Data collection
- Data collection reports
- Data processing
- Data processing: Coding
- Data processing: Classification
- Data processing: Transcriptions of interviews
- Data processing: Weighting
- Data processing: Aggregation
- Data processing: Composite measures
- Data processing: Derivation
- Data processing: Quality checks
- Data processing: Data integration
- Data processing: Disclosure limitation
- Data processing: Imputation
- Metadata production
- Metadata editing
- Metadata translation





- Final report
- Evaluation
- Original release
- Deposit
- Preservation package production
- Dissemination package production
- Dissemination package release
- Data analysis reports
- New version release
- Other

#### **Related Publication / ID Type**

- ark
- arXiv
- bibcode
- doi
- ean13
- eissn
- handle
- isbn
- issn
- istc
- lissn
- Isid
- pmid
- purl
- upc
- url
- urn

#### Contributor / based on DataCite 4.0

- Data collector
- Data curator
- Data manager
- Editor
- Funder
- Hosting institution
- Metadata author
- Project leader
- Project manager
- Project member
- Registration agency
- Registration authority
- Related person
- Researcher
- Research group
- Rights holder





- Sponsor
- Supervisor
- Work Package Leader
- Other

#### Contributor / Contributor Identifier Scheme

- DOI
- idHAL
- info
- ISNI
- LCNA
- ORCID

Note : ORCID recommended

Language : ISO 639-1 standard

#### **Geospatial Metadata**

Geographic Coverage

#### Country / Nation : ISO 3166-1 standard

#### **Conformity / Degree**

- Conformant
- Not Conformant
- Not evaluated

#### Social Science and Humanities Metadata

Unit of Analysis : DDI Controlled Vocabulary for Analysis Unit

- Individual
- Organization
- Family
- Family: Household family
- Household
- Housing Unit
- Event/Process
- Geographic Unit
- Time Unit
- Text Unit
- Group
- Object
- Other

#### Time Method : DDI Controlled Vocabulary for Time Method

- Longitudinal
- Longitudinal: Cohort/Event-based
- Longitudinal: Trend/Repeated cross-section





- Longitudinal: Panel
- Longitudinal: Panel: Continuous
- Longitudinal: Panel: Interval
- Time Series
- TimeSeries: Continuous
- TimeSeries: Discrete
- Cross-section
- Cross-section ad-hoc follow-up
- Other

#### Sampling Procedure: DDI Controlled Vocabulary for Sampling Procedure

- Total universe/Complete enumeration
- Probability
- Probability: Simple random
- Probability: Systematic random
- Probability: Stratified
- Probability: Stratified: Proportional
- Probability: Stratified: Disproportional
- Probability: Cluster
- Probability: Cluster: Simple random
- Probability: Cluster: Stratified random
- Probability: Multistage
- Non-probability
- Non-probability: Availability
- Non-probability: Purposive
- Non-probability: Quota
- Non-probability: Respondent-assisted
- Mixed probability and non-probability

Collection Mode : DDI Controlled Vocabulary for Mode Of Collection

- Interview
- Face-to-face interview
- Face-to-face interview: CAPI
- Face-to-face interview: PAPI
- Telephone interview
- Telephone interview: CATI
- E-mail interview
- Web-based interview
- Self-administered questionnaire
- Fixed form self-administered questionnaire
- Fixed form self-administered questionnaire: E-mail
- Fixed form self-administered questionnaire: Paper
- Fixed form self-administered questionnaire: SMS/MMS
- Fixed form self-administered questionnaire: Web-based
- Interactive self-administered questionnaire
- Interactive self-administered questionnaire: CASI





- Interactive self-administered questionnaire: CASI: VCASI
- Interactive self-administered questionnaire: CASI: ACASI
- Interactive self-administered questionnaire: CASI: TACASI
- Interactive self-administered questionnaire: CAWI
- Focus group
- Face-to-face focus group
- Telephone focus group
- Online focus group
- Self-administered writings and/or diaries
- Self-administered writings and/or diaries: E-mail
- Self-administered writings and/or diaries: Paper
- Self-administered writings and/or diaries: Web-based
- Observation
- Field observation
- Participant field observation
- Non-participant field observation
- Laboratory observation
- Participant laboratory observation
- Non-participant laboratory observation
- Computer-based observation
- Experiment
- Laboratory experiment
- Field/Intervention experiment
- Web-based experiment
- Recording
- Content coding
- Transcription
- Compilation/Synthesis
- Summary
- Aggregation
- Simulation
- Measurements and tests
- Educational measurements and tests
- Physical measurements and tests
- Psychological measurements and tests
- Other

#### Life Sciences Metadata

Ontologies Experimental Factor Ontology (EFO); Ontology of Clinical Research (OCRE); Ontology for Biomedical Investigations; Infectious Disease Ontology (IDO); FlyBase Controlled Vocabulary (FBCV); BRENDA tissue / enzyme source (BTO); Chemical Methods Ontology (CHMO); Medical Subject Headings (MESH); eagle-i resource ontology

Design Type

- Case Control
- Cross Sectional
- Cohort Study
- Nested Case Control Design





- Not Specified
- Parallel Group Design
- Perturbation Design
- Randomized Controlled Trial
- Technological Design

#### **Factor Type**

- Age
- Biomarkers
- Cell Surface Markers
- Cell Type/Cell Line
- Developmental Stage
- Disease State
- Drug Susceptibility
- Extract Molecule
- Genetic Characteristics
- Immunoprecipitation Antibody
- Organism
- Other
- Passages
- Platform
- Sex
- Strain
- Time Point
- Tissue Type
- Treatment Compound
- Treatment Type

#### Organism : NCBI Taxonomy

- Arabidopsis thaliana
- Bos taurus
- Caenorhabditis elegans
- Chlamydomonas reinhardtii
- Danio rerio (zebrafish)
- Dictyostelium discoideum
- Drosophila melanogaster
- Escherichia coli
- Hepatitis C virus
- Homo sapiens
- Mus musculus
- Mycobacterium africanum
- Mycobacterium canetti
- Mycobacterium tuberculosis
- Mycoplasma pneumoniae
- Oryza sativa
- Plasmodium falciparum
- Pneumocystis carinii





- Rattus norvegicus
- Saccharomyces cerevisiae (brewer's yeast)
- Schizosaccharomyces pombe
- Takifugu rubripes
- Xenopus laevis
- Zea mays
- Other

#### Measurement Type

- cell counting
- cell sorting
- clinical chemistry analysis
- copy number variation profiling
- DNA methylation profiling
- DNA methylation profiling (Bisulfite-Seq)
- DNA methylation profiling (MeDIP-Seq)
- drug susceptibility
- environmental gene survey
- genome sequencing
- hematology
- histology
- Histone Modification (ChIP-Seq)
- loss of heterozygosity profiling
- metabolite profiling
- metagenome sequencing
- protein expression profiling
- protein identification
- protein-DNA binding site identification
- protein-protein interaction detection
- protein-RNA binding (RIP-Seq)
- SNP analysis
- targeted sequencing
- transcription factor binding (ChIP-Seq)
- transcription factor binding site identification
- transcription profiling
- transcription profiling
- transcription profiling (Microarray)
- transcription profiling (RNA-Seq)
- TRAP translational profiling
- Other

#### Technology Type

- culture based drug susceptibility testing, single concentration
- culture based drug susceptibility testing, two concentrations
- culture based drug susceptibility testing, three or more concentrations (minimium inhibitory concentration measurement)
- DNA microarray





- flow cytometry
- gel electrophoresis
- mass spectrometry
- NMR spectroscopy
- nucleotide sequencing
- protein microarray
- real time PCR
- no technology required
- Other

#### **Technology Platform**

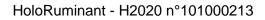
- 210-MS GC Ion Trap (Varian)
- 220-MS GC Ion Trap (Varian)
- 225-MS GC Ion Trap (Varian)
- 240-MS GC Ion Trap (Varian)
- 300-MS quadrupole GC/MS (Varian)
- 320-MS LC/MS (Varian)
- 325-MS LC/MS (Varian)
- 320-MS GC/MS (Varian)
- 500-MS LC/MS (Varian)
- 800D (Jeol)
- 910-MS TQ-FT (Varian)
- 920-MS TQ-FT (Varian)
- 3100 Mass Detector (Waters)
- 6110 Quadrupole LC/MS (Agilent)
- 6120 Quadrupole LC/MS (Agilent)
- 6130 Quadrupole LC/MS (Agilent)
- 6140 Quadrupole LC/MS (Agilent)
- 6310 Ion Trap LC/MS (Agilent)
- 6320 Ion Trap LC/MS (Agilent)
- 6330 Ion Trap LC/MS (Agilent)
- 6340 Ion Trap LC/MS (Agilent)
- 6410 Triple Quadrupole LC/MS (Agilent)
- 6430 Triple Quadrupole LC/MS (Agilent)
- 6460 Triple Quadrupole LC/MS (Agilent)
- 6490 Triple Quadrupole LC/MS (Agilent)
- 6530 Q-TOF LC/MS (Agilent)
- 6540 Q-TOF LC/MS (Agilent)
- 6210 TOF LC/MS (Agilent)
- 6220 TOF LC/MS (Agilent)
- 6230 TOF LC/MS (Agilent)
- 7000B Triple Quadrupole GC/MS (Agilent)
- AccuTO DART (Jeol)
- AccuTOF GC (Jeol)
- AccuTOF LC (Jeol)
- ACQUITY SQD (Waters)
- ACQUITY TQD (Waters)
- Agilent





- Agilent 5975E GC/MSD (Agilent)
- Agilent 5975T LTM GC/MSD (Agilent)
- 5975C Series GC/MSD (Agilent)
- Affymetrix
- amaZon ETD ESI Ion Trap (Bruker)
- amaZon X ESI Ion Trap (Bruker)
- apex-ultra hybrid Qq-FTMS (Bruker)
- API 2000 (AB Sciex)
- API 3200 (AB Sciex)
- API 3200 QTRAP (AB Sciex)
- API 4000 (AB Sciex)
- API 4000 QTRAP (AB Sciex)
- API 5000 (AB Sciex)
- API 5500 (AB Sciex)
- API 5500 QTRAP (AB Sciex)
- Applied Biosystems Group (ABI)
- AQI Biosciences
- Atmospheric Pressure GC (Waters)
- autoflex III MALDI-TOF MS (Bruker)
- autoflex speed(Bruker)
- AutoSpec Premier (Waters)
- AXIMA Mega TOF (Shimadzu)
- AXIMA Performance MALDI TOF/TOF (Shimadzu)
- A-10 Analyzer (Apogee)
- A-40-MiniFCM (Apogee)
- Bactiflow (Chemunex SA)
- Base4innovation
- BD BACTEC MGIT 320
- BD BACTEC MGIT 960
- BD Radiometric BACTEC 460TB
- BioNanomatrix
- Cell Lab Quanta SC (Becman Coulter)
- Clarus 560 D GC/MS (PerkinElmer)
- Clarus 560 S GC/MS (PerkinElmer)
- Clarus 600 GC/MS (PerkinElmer)
- Complete Genomics
- Cyan (Dako Cytomation)
- CyFlow ML (Partec)
- Cyow SL (Partec)
- CyFlow SL3 (Partec)
- CytoBuoy (Cyto Buoy Inc)
- CytoSence (Cyto Buoy Inc)
- CytoSub (Cyto Buoy Inc)
- Danaher
- DFS (Thermo Scientific)
- Exactive(Thermo Scientific)
- FACS Canto (Becton Dickinson)
- FACS Canto2 (Becton Dickinson)
- FACS Scan (Becton Dickinson)

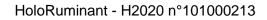






- FC 500 (Becman Coulter)
- GCmate II GC/MS (Jeol)
- GCMS-QP2010 Plus (Shimadzu)
- GCMS-QP2010S Plus (Shimadzu)
- GCT Premier (Waters)
- GENEQ
- Genome Corp.
- GenoVoxx
- GnuBio
- Guava EasyCyte Mini (Millipore)
- Guava EasyCyte Plus (Millipore)
- Guava Personal Cell Analysis (Millipore)
- Guava Personal Cell Analysis-96 (Millipore)
- Helicos BioSciences
- Illumina
- Indirect proportion method on LJ medium
- Indirect proportion method on Middlebrook Agar 7H9
- Indirect proportion method on Middlebrook Agar 7H10
- Indirect proportion method on Middlebrook Agar 7H11
- inFlux Analyzer (Cytopeia)
- Intelligent Bio-Systems
- ITQ 700 (Thermo Scientific)
- ITQ 900 (Thermo Scientific)
- ITQ 1100 (Thermo Scientific)
- JMS-53000 SpiralTOF (Jeol)
- LaserGen
- LCMS-2020 (Shimadzu)
- LCMS-2010EV (Shimadzu)
- LCMS-IT-TOF (Shimadzu)
- Li-Cor
- Life Tech
- LightSpeed Genomics
- LCT Premier XE (Waters)
- LCQ Deca XP MAX (Thermo Scientific)
- LCQ Fleet (Thermo Scientific)
- LXQ (Thermo Scientific)
- LTQ Classic (Thermo Scientific)
- LTQ XL (Thermo Scientific)
- LTQ Velos (Thermo Scientific)
- LTQ Orbitrap Classic (Thermo Scientific)
- LTQ Orbitrap XL (Thermo Scientific)
- LTQ Orbitrap Discovery (Thermo Scientific)
- LTQ Orbitrap Velos (Thermo Scientific)
- Luminex 100 (Luminex)
- Luminex 200 (Luminex)
- MACS Quant (Miltenyi)
- MALDI SYNAPT G2 HDMS (Waters)
- MALDI SYNAPT G2 MS (Waters)
- MALDI SYNAPT HDMS (Waters)







- MALDI SYNAPT MS (Waters)
- MALDI micro MX (Waters)
- maXis (Bruker)
- maXis G4 (Bruker)
- microflex LT MALDI-TOF MS (Bruker)
- microflex LRF MALDI-TOF MS (Bruker)
- microflex III MALDI-TOF MS (Bruker)
- micrOTOF II ESI TOF (Bruker)
- micrOTOF-Q II ESI-Qq-TOF (Bruker)
- microplate Alamar Blue (resazurin) colorimetric method
- Mstation (Jeol)
- MSQ Plus (Thermo Scientific)
- NABsys
- Nanophotonics Biosciences
- Network Biosystems
- Nimblegen
- Oxford Nanopore Technologies
- Pacific Biosciences
- Population Genetics Technologies
- Q1000GC UltraQuad (Jeol)
- Quattro micro API (Waters)
- Quattro micro GC (Waters)
- Quattro Premier XE (Waters)
- QSTAR (AB Sciex)
- Reveo
- Roche
- Seirad
- solariX hybrid Qq-FTMS (Bruker)
- Somacount (Bently Instruments)
- SomaScope (Bently Instruments)
- SYNAPT G2 HDMS (Waters)
- SYNAPT G2 MS (Waters)
- SYNAPT HDMS (Waters)
- SYNAPT MS (Waters)
- TripleTOF 5600 (AB Sciex)
- TSQ Quantum Ultra (Thermo Scientific)
- TSQ Quantum Access (Thermo Scientific)
- TSQ Quantum Access MAX (Thermo Scientific)
- TSQ Quantum Discovery MAX (Thermo Scientific)
- TSQ Quantum GC (Thermo Scientific)
- TSQ Quantum XLS (Thermo Scientific)
- TSQ Vantage (Thermo Scientific)
- ultrafleXtreme MALDI-TOF MS (Bruker)
- VisiGen Biotechnologies
- Xevo G2 QTOF (Waters)
- Xevo QTof MS (Waters)
- Xevo TQ MS (Waters)
- Xevo TQ-S (Waters)
- Other





Journal Matadata Type of Article

- abstract
- addendum
- announcement
- article-commentary
- book review
- books received
- brief report
- calendar
- case report
- collection
- correction
- data paper
- discussion
- dissertation
- editorial
- in brief
- introduction
- letter
- meeting report
- news
- obituary
- oration
- partial retraction
- product review
- rapid communication
- reply
- reprint
- research article
- retraction
- review article
- translation
- other

Semantic resource Type of semantic resource

- Classification scheme
- Dictionary
- Gazetteer
- Glossary
- List
- Name authority list
- Ontology
- Semantic network
- Subject heading scheme





- Taxonomy
- Terminology
- Thesaurus
- Other

Representation syntax

- JSON-LD
- LD Patch
- Microdata
- N3
- N-Quads
- OBO
- OWL Functional Syntax
- OWL Manchester Syntax
- OWL XML Serialization
- POWDER
- POWDER-S
- PROV-N
- PROV-XML
- RDF/JSON
- RDF/XML
- RDFa
- RIF XML Syntax
- TriG
- Turtle
- Other

Nature of the semantic resource

- Application Ontology
- Core Ontology
- Domain Ontology
- Task Ontology
- Upper Level Ontology
- Vocabulary
- Other

Designed for

- Annotation Task
- Configuration Task
- Filtering Task
- Indexing Task
- Integration Task
- Matching Task
- Mediation Task
- Personalization Task





- Query Formulation Task
- Query Rewriting Task
- Search Task

Information on current version / Version status

- alpha
- beta
- production
- retired

