



HAL
open science

A project-scale map of metadata to improve future data management

Mélanie Pétéra, Cecile Cabasson, Blandine Comte, Christophe Duperier, Olivier Filangi, Sylvain Prigent, Estelle Pujos-Guillot, Franck Giacomoni

► To cite this version:

Mélanie Pétéra, Cecile Cabasson, Blandine Comte, Christophe Duperier, Olivier Filangi, et al. A project-scale map of metadata to improve future data management. Analytics 2022, Sep 2022, Nantes, France. . hal-03776494

HAL Id: hal-03776494

<https://hal.inrae.fr/hal-03776494v1>

Submitted on 13 Sep 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A project-scale map of metadata to improve future data management

Mélanie PÉTÉRA^{1,5}; Cécile CABASSON^{2,3,5}; Blandine COMTE^{1,5}; Christophe DUPÉRIER^{1,5}; Olivier FILANGI^{4,5}; Sylvain PRIGENT^{2,3,5}; The MetaboHUB Consortium⁵; Estelle PUJOS-GUILLOT^{1,5}; Franck GIACOMONI^{1,5}

¹ Université Clermont Auvergne, INRAE, UNH, Plateforme d'Exploration du Métabolisme, MetaboHUB Clermont, F-63000 Clermont-Ferrand, France

² Univ. Bordeaux, INRAE, UMR1332 BFP, F-33382 Villenave d'Ornon, France

³ Bordeaux Metabolome, MetaboHUB, PHENOME-EMPHASIS, F-33140 Villenave d'Ornon, France

⁴ IGEPP, INRAE, Institut Agro, Université de Rennes, F-35653 Le Rheu, France

⁵ ANR INBS-0010, Toulouse, France

Introduction

Today, the intra-lab application of best practices in the metabolomics field usually guarantees an adequate data exploitation within a single lab. However, the growing interest in multi-analyses designs (e.g. complementary analytical platforms, variety of matrices, multi-omics), as well as the need of data sharing and reuse, increase the difficulty of data management. Indeed, managing the multiplicity and the heterogeneity of information involved is required to achieve relevant knowledge extraction from metabolomics data.

Within the MetaboHUB national infrastructure (MTH), one objective is to optimize data handling, especially metadata, to facilitate large-scale analyses, multi-platforms studies, and data FAIRisation (Findability, Accessibility, Interoperability, Reusability). In particular, this fits in the MetaboHUB scientific roadmap that promotes the open science development in the field of metabolomics.

Material & Method

In the context of metabolomic and lipidomic studies, data production and analysis come along with a large diversity of metadata (data of the data). To identify clearly-defined bottlenecks and targets for future improvement in data management, the objective of this work was to build a metadata map at the scale of a scientific project. Aiming for completeness, this map was constructed in a collaborative and multidisciplinary way involving chemists, biologists, data stewards as well as computer scientists, combining their respective experience and knowledge.

Materials:

- MetaboHUB expertise combining a large variety of skills, profiles and fields
- Data repository resource examples where metadata are central (e.g. EBI MetaboLights^[1])
- Publications on the data cycle (e.g. Savoi et al - Grapevine community^[2])
- Online ontology resources (e.g. BioPortal^[3], AgroPortal^[4], OLS ontology search^[5])

Methodology:

- Common meetings, including all area experts: **topics borders, constitutive of the future map**
- Iterative map building through various methods: face to face, by groups, by fields and scales: **we do not want to miss any metadata.**
- Work on information representation. How to represent the diversity and the complexity of different points of view and interests (e.g. chemistry vs. biology vs. bioinformatics)? Examples:
 - ✓ Which choice to make to group metadata in categories depending on how to use the mapping afterward?
 - ✓ How to combine multiple grouping, without losing essential information as well as keeping the map easy-to-handle for all actors?
- Work on how to transform this preliminary map into active and productive initiatives, useful for FAIR data and the open science movement.

Results

We constructed a project-oriented map of metadata that sorts and labels metadata through various angles. This complete map can be adapted in refined maps depending on the objective of representation. In this communication, we present a specific implementation of our first built map (Figure 1) where we display seven general topics to address through three temporal steps.

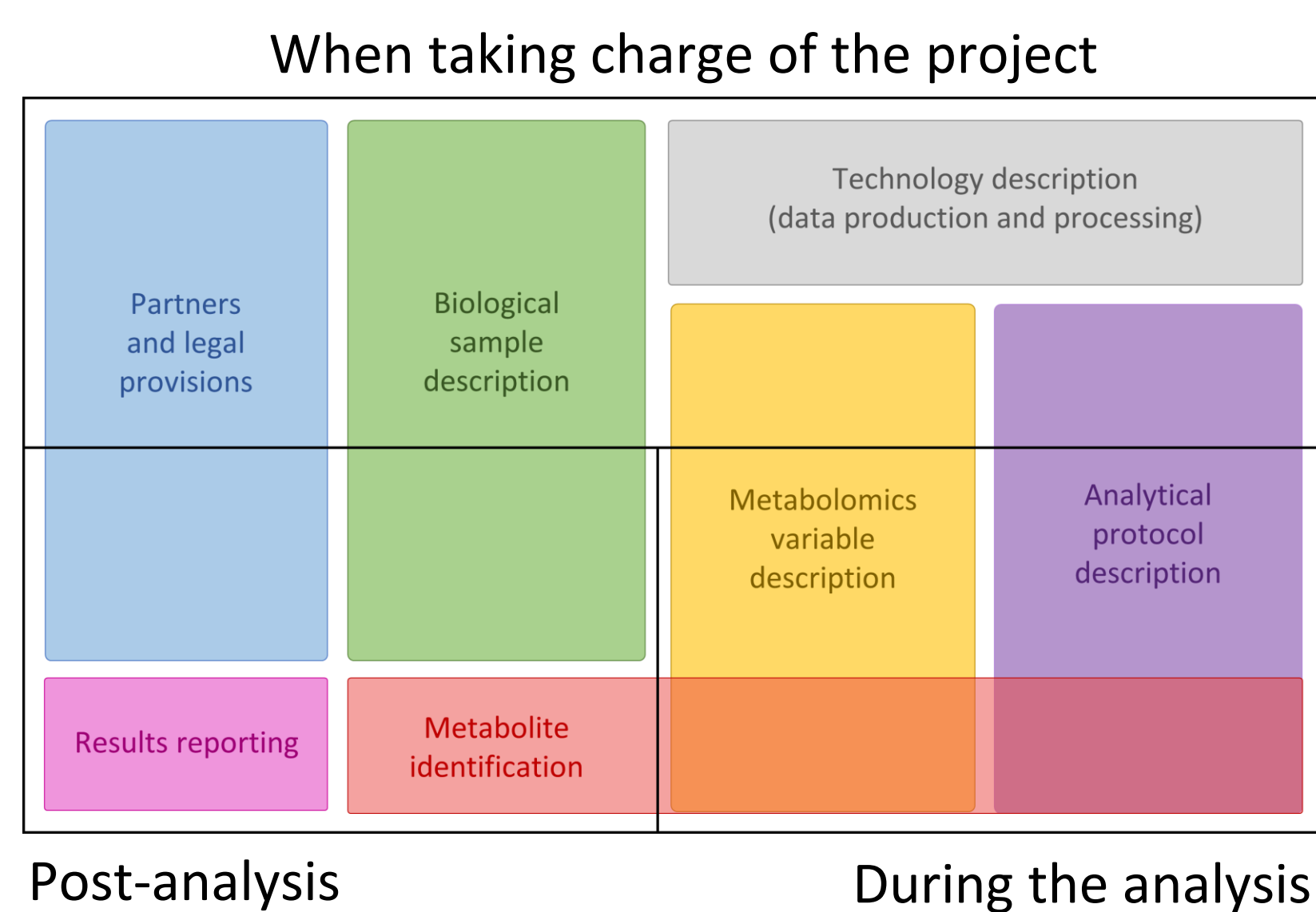


Figure 1: A map that links the timeline of metadata generation and the various types of metadata at the scale of a scientific project

Topics were also declined in more accurate aspects, illustrated in Figure 2. These aspects in each topic area, are source of metadata, usually difficult to store, retrieve and publish. Based on the resulting metadata map, targets (areas and topics) to be further investigated were identified, enabling the construction of transversal working groups at the MetaboHUB consortium scale. In particular, this work enables to focus efforts on clearly defined issues to improve standardisation of practices regarding data management and metadata documentation.

- Examples of potential working topics:
- 1 Naming standards inside the infrastructure
 - 2 Harmonised use of ontologies
- Examples of existing MTH task forces:
- 1 Task force about data integration
 - 2 Working force about reporting results

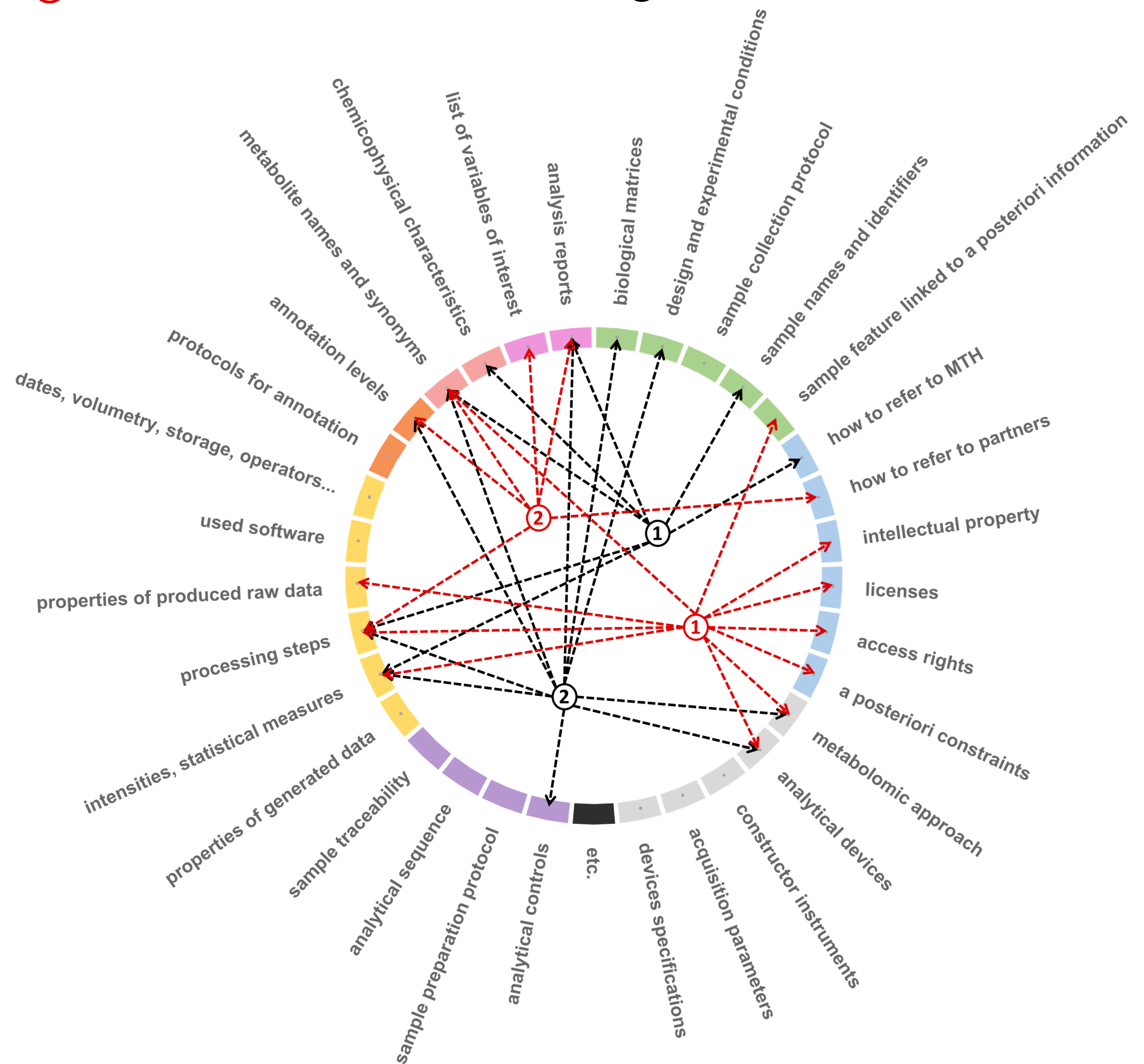


Figure 2: Examples of metadata types listed by fields identified in Figure 1 and examples of mapping with possible working topics and MetaboHUB task forces activities

Example of mapping of a concrete MTH development project:

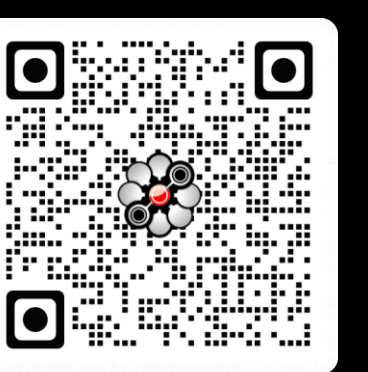
The project Metabolomics Semantics Data Lake aims to build and provide a large-scale distributed infrastructure for data-processing and massive integration of semantic information about metabolomics data studies. This system is strongly structured from semantic web technologies to maximize the reuse of existent ontologies and knowledge (NCBI, EBI) and finally manage heterogeneous metabolomics content from the MetaboHUB consortium.

Interacting fields identified with the scope of routine usage:

- processing step
- used software
- dates, volumetry, storage
- analyses reports
- sample-collection processing
- sample names and identifiers
- sample feature linked to a posteriori information

Conclusion

In conclusion, this collaborative map construction has been shown to be an efficient tool to draw a clear « where do we stand / where do we go » picture inside a national infrastructure like **MetaboHUB** regarding project-scale metadata. This facilitates the definition of a precise data management. Such an approach could be translated within other infrastructures, consortia and/or communities.



SCAN ME