

## Title

# The metagenomic data life-cycle: standards and best practices

## Authors

Petra ten Hoopen<sup>1</sup>, Robert D. Finn<sup>1</sup>, Lars Ailo Bongo<sup>2</sup>, Erwan Corre<sup>3</sup>, Bruno Fosso<sup>4</sup>, Folker Meyer<sup>5</sup>, Alex Mitchell<sup>1</sup>, Eric Pelletier<sup>6,7,8</sup>, Graziano Pesole<sup>4,9</sup>, Monica Santamaria<sup>4</sup>, Nils Peder Willassen<sup>2</sup>, and Guy Cochrane<sup>1\*</sup>

<sup>1</sup>European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, United Kingdom

<sup>2</sup>UiT The Arctic University of Norway, Tromsø N-9037, Norway

<sup>3</sup>CNRS-UPMC, FR 2424, Station Biologique, Roscoff 29680, France

<sup>4</sup>Institute of Biomembranes, Bioenergetics and Molecular Biotechnologies, CNR, Bari 70126, Italy

<sup>5</sup>Argonne National Laboratory, Argonne IL60439, USA

<sup>6</sup>Genoscope, CEA, Évry 91000, France

<sup>7</sup>CNRS / UMR-8030, Évry 91000, France

<sup>8</sup>Université Évry val d'Essonne, Évry 91000, France

<sup>9</sup>Department of Biosciences, Biotechnologies and Biopharmaceutics, University of Bari "A. Moro", Bari 70126, Italy

1  
2  
3  
4  
5 \* corresponding author  
6

7  
8 Guy Cochrane  
9

10 European Molecular Biology Laboratory  
11

12 European Bioinformatics Institute  
13

14  
15 Wellcome Genome Campus  
16

17  
18 Hinxton  
19

20 Cambridge CB10 1SD  
21

22  
23 United Kingdom  
24

25  
26 Tel: +44(0)1223 494444  
27

28 Email: [cochrane@ebi.ac.uk](mailto:cochrane@ebi.ac.uk)  
29  
30

31  
32  
33 Email addresses:  
34

35  
36 AM: [mitchell@ebi.ac.uk](mailto:mitchell@ebi.ac.uk)  
37

38 BF: [b.fosso@ibbe.cnr.it](mailto:b.fosso@ibbe.cnr.it)  
39

40 EC: [erwan.corre@sb-roscoff.fr](mailto:erwan.corre@sb-roscoff.fr)  
41

42 EP: [eric.pelletier@genoscope.cns.fr](mailto:eric.pelletier@genoscope.cns.fr)  
43

44 FM: [folker@anl.gov](mailto:folker@anl.gov)  
45

46 GC: [cochrane@ebi.ac.uk](mailto:cochrane@ebi.ac.uk)  
47

48 GP: [g.pesole@ibbe.cnr.it](mailto:g.pesole@ibbe.cnr.it)  
49

50 LAB: [lars.ailo.bongo@uit.no](mailto:lars.ailo.bongo@uit.no)  
51

52 MS: [m.santamaria@ibbe.cnr.it](mailto:m.santamaria@ibbe.cnr.it)  
53

54 NPW: [nils-peder.willassen@uit.no](mailto:nils-peder.willassen@uit.no)  
55

56 PTH: [petratenhoopen@yahoo.co.uk](mailto:petratenhoopen@yahoo.co.uk)  
57  
58  
59  
60  
61  
62  
63  
64  
65

RDF: [rdf@ebi.ac.uk](mailto:rdf@ebi.ac.uk)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

# Abstract

Metagenomics data analyses from independent studies can only be compared if the analysis workflows are described in a harmonised way. In this overview, we have mapped the landscape of data standards available for the description of essential steps in metagenomics: (1) material sampling, (2) material sequencing (3) data analysis and (4) data archiving & publishing.

Taking examples from marine research, we summarise essential variables used to describe material sampling processes and sequencing procedures in a metagenomics experiment.

These aspects of metagenomics dataset generation have been to some extent addressed by the scientific community but greater awareness and adoption is still needed.

We emphasise the lack of standards relating to reporting how metagenomics datasets are analysed and how the metagenomics data analysis outputs should be archived and published. We propose best practice as a foundation for a community standard to enable reproducibility and better sharing of metagenomics datasets, leading ultimately to greater metagenomics data reuse and repurposing.

## Keywords

Metagenomics, metadata, standard, best practice, sampling, sequencing, data analysis.

# Background

Recent technological advances allow researchers to examine communities of organisms using such methods as metagenomics (enumerating and exploring the genes within a community), metatranscriptomics (profiling and quantifying patterns of gene expression within a community), metabarcoding (profiling marker loci for species diversity and phylogenetic purposes) and metaproteomics (profiling the protein component of a community), enabling comprehensive insights into community composition and function (**Figure 1**). The increased popularity of these meta-omics methods, driven not least by ever decreasing cost, leads to increasing scale and complexity of experimental data and in approaches to their analysis. In addition, there is growing demand for comparisons between communities that have been studied independently, often using very different approaches. However, meaningful interpretation across studies (either through aggregation and interpretation of existing published analyses or through meta-analysis of published experimental data using a uniform method) is challenging. A number of reasons exist for this: (1) each 'omic' analysis workflow is a complex process, consisting of disparate and diverse tasks, ranging from sample collection and processing to data generation and analysis, where each task has many parameters that can affect analysis outputs (for example, it has been shown that a major factor explaining correlations within metagenomics datasets can be DNA preparation and sequencing, [1]); (2) each variable is frequently recorded in a non-standardised way, or not recorded at all; (3) presentation formats of the produced omics data are not unified; (4) omics experimental data and related analysis outputs are either dispersed in several public repositories, or not archived at all.

1  
2 Here, we review the workflow for metagenomics data generation and analysis. Where  
3  
4 possible, we specify essential parameters in the workflow and advise on standardised  
5  
6 systematic reporting of these as variables. We build on the expertise of major public  
7  
8 genomic and metagenomic resources: the European Nucleotide Archive (ENA) [2] and  
9  
10 EMBL-EBI Metagenomics (EMG) [3] at the EMBL European Bioinformatics Institute in UK;  
11  
12 MG-RAST [4] at Argonne National Laboratory in USA; and the extensive knowledge bases in  
13  
14 metagenomics available at research centers of excellence, the UiT in Norway, Genoscope in  
15  
16 France, SB-Roscoff in France and CNR in Italy.  
17  
18  
19  
20  
21  
22  
23  
24

25 For the purposes of this paper, we will predominately use marine metagenomics as a ‘use  
26  
27 case’ to highlight the standards environment that we describe. However, we believe that  
28  
29 these examples will broadly translate to all areas of metagenomics research, regardless of  
30  
31 the environment under study. From the outset, we stress that we do not wish to promote a  
32  
33 specific workflow, but rather to demonstrate the importance of having systematic reporting  
34  
35 conventions that accurately describe any chosen workflow, from sampling through to the  
36  
37 presentation of analysis outputs. Our aim is to describe conventions and standards that are  
38  
39 inclusive and extensible, and able to cope with evolving scientific developments in the field.  
40  
41 Furthermore, where a given standard has not emerged, we will point to, or propose, a  
42  
43 generalised ‘best practice’ that can be used in its place. While this may produce a  
44  
45 foundation from which a new standard could be proposed, any additional formal scientific  
46  
47 standards need to come from the community and be ratified by scientific bodies, such as the  
48  
49 Genomics Standards Consortium (GSC) [5].  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 For this paper, we have chosen a structure in which we introduce the generic data model  
2 that has been adopted by those working with metagenomic data and then move through  
3 the various practical steps – from sampling, through assay and analysis, to the archiving of  
4 analysis outputs - that a metagenomicist takes through a metagenomics investigation (see  
5 also **Figure 1**).

## 15 Overview of the metagenomics data model

19 The introduction of new generation sequencing technologies has enabled even small  
20 research groups to generate large-scale sequencing data. The resultant DNA sequences and  
21 associated information are typically captured in several interconnected objects (**Figure 2**),  
22 which represent the following concepts:

- 29 ● **Study:** Information about the scope of a sequencing effort that groups together all  
30 data of the project.
- 35 ● **Sample:** Information about provenance and characteristics of the sequenced  
36 samples.
- 41 ● **Experiment:** Information about the sequencing experiments, including library and  
42 instrument details.
- 46 ● **Run:** An output of a sequencing experiment containing sequencing reads  
47 represented in data files.
- 51 ● **Analysis:** A set of outputs computed from primary sequencing results, including  
52 sequence assemblies, functional and taxonomic annotations.

58 **[Figure 2]**

1  
2 Information associated with DNA sequence is frequently referred to as **'metadata'**. This  
3  
4 includes all information described in the study, sample, experiment and run data objects,  
5  
6 spanning sampling context, description of sample processing, experimental design, library  
7  
8 creation, sequencer configuration and provenance information required for attribution and  
9  
10 credit to comply with best scientific practice for publication in the academic literature and  
11  
12 to inform processes around Access and Benefit Sharing. **Primary data** represent, in this  
13  
14 context, primary "raw" experimental sequence reads produced by sequencing machines.  
15  
16 (On occasion, some basic data processing, such as quality control (filtering out of poor-  
17  
18 quality reads, clipping of low-quality regions, etc.), is applied to "raw" primary data and  
19  
20 these processed data are retained as primary; while it is preferable to retain true "raw"  
21  
22 primary data, perhaps in addition to these processed data, it is important to apply broadly  
23  
24 accepted processing methods and to describe these methods as part of the metadata  
25  
26 record.) Following this, for some metagenomics studies, the **primary data** are analysed  
27  
28 directly (e.g. 16S or 18S rRNA gene amplicon studies), while in others, they are assembled  
29  
30 into contigs before undergoing further analysis. Regardless of the approach, the output of  
31  
32 any computational analysis process (including assembly) on the primary data are here  
33  
34 referred to as **derived data**. We discuss derived data in more detail below, but the more  
35  
36 harmonised the formats and validations for data and metadata objects, the more easily the  
37  
38 generated data be shared, discovered, re-used and re-purposed.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53

54 Each metagenomics initiative has a scope, aim and one or more (human) contributors in  
55  
56 each step of the workflow, who may be distributed over a wide geographical area. It is  
57  
58 essential to capture contextual information regarding the contributors, since this supports  
59  
60  
61  
62  
63  
64  
65



1 appropriate attribution and credit, and clarifies the responsible parties for each step of the  
2 workflow. Contributors to (1) material sampling (2) primary data generation and (3) derived  
3 data generation should always be clearly presented in data records. Minimum metadata  
4  
5 checklists frequently do not specifically capture data generating or contributing institutions.  
6  
7 However, this information is frequently available and can be parsed from the registration  
8  
9 systems for reporting individual steps of the data generation workflow or from associated  
10  
11 peer-reviewed publications.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

## 22 Sampling

23  
24  
25  
26  
27 The method of collecting a sample (a fundamental unit of material isolated from the  
28 surrounding environment) is dictated by the nature of the community under investigation,  
29 the environment in which it is found and the type of 'omics' investigation being performed.  
30  
31 The slightest deviation in method, regardless of the protocol chosen, can have a profound  
32 impact on the final 'omics' analysis results. It is therefore essential that the details of the  
33 sampling process are captured accurately and in a standardised way.  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44

45 Domain experts are in the best position to formulate opinions on the general scope and  
46 content of contextual data (environmental characteristics observed or measured during  
47 sample collection) and methodological variables (such as sampling volume and filtration  
48 method). These opinions are conventionally formalised as data reporting standards by  
49 community initiatives such as the GSC for genomics data [5] - on several of which we expand  
50 below - the Proteomics Standards Initiative (PSI, [6]) for proteomics data or the Group on  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 Earth Observations Biodiversity Observation Network (GEO BON) for the various dimensions  
2 of biodiversity, including genetic variation and biodiversity data [7,8].  
3  
4  
5  
6

7 The Minimum Information about Metagenomic Sequence (MIMS, [9]) is a GSC-developed  
8 data-reporting standard, designed for accurate reporting of contextual information for  
9 samples associated with metagenomic sequencing, and is also largely applicable to  
10 metatranscriptomics studies. Minimum Information about a MARKer gene Sequence  
11 (MIMARKS, [10]) is another GSC-developed contextual data reporting standard for reporting  
12 information about a metabarcoding study, which is referred to in the standard as the  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23 'MIMARKS-survey investigation type'.  
24  
25  
26  
27

28 MIMS and MIMARKS are a part of a broader GSC standard, the Minimum Information of any  
29 (x) Sequence (MIxS) [11], which describes 15 different environmental packages that can be  
30 used to specify the environmental context of a sequenced microbial community, such as air,  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

66 To illustrate, **Table 1** summarises the minimum set of elements required for description of a  
67 metagenomic sample taken from an aquatic environment. It uses MIMS mandatory  
68 descriptors, combined with the mandatory descriptors of the Water Environment package.  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100

1 that support data reporting and presentation, such as the submission tools of the databases  
2 of the International Nucleotide Sequence Database Collaboration [13] and ISAtools [14]. It  
3 remains up to the experimentalist to choose the most appropriate package from within the  
4 checklist bundle for their study, thereby defining the list of fields that will be used to  
5 capture relevant metadata. Before embarking on a metagenomics study, we recommend  
6 that the appropriate checklist be identified, so that the appropriate metadata can be  
7 captured during the experiment, rather than retrospectively having to determine these  
8 metadata.  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

23 **[Table 1]**  
24  
25  
26  
27

28 The M2B3 data reporting and service standard [15] specifically addresses contextual data  
29 relating to marine microbial samples. It represents a common denominator of contextual  
30 data from data standards used in the public genomic data archives (MIxS, Version 4.0, [12]),  
31 pan-European network of oceanographic data archives (CDI schema, Version 3.0, [16]) and  
32 pan-European network of biodiversity data resources (OBIS schema, Version 1.1, [17]). This  
33 M2B3 unified data standard significantly simplifies contextual data reporting, since it  
34 provides an interoperable solution for sharing contextual data across data archives from  
35 different scientific domains. A minimum M2B3 checklist for reporting contextual data  
36 associated with marine microbial samples is summarised in **Table 2**.  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52

53 **[Table 2]**  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 For most adopted standards of this type, only a few fields of contextual data are mandatory,  
2 reflecting the balance between usability for the experimentalist reporting his/her science  
3 and consumers re-using this science; limiting the number of mandatory fields lowers the  
4 burden for experimentalists to comply with the standard, while a small number of  
5 parameters are universally, or near-universally, required for downstream analysis. The  
6 importance of the optional MlxS and M2B3 fields for metagenomic data analysis is detailed  
7 in **Table 3**.

8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20 We wish here to note a convention on the handling of replicate samples. Since biological  
21 replicates are separate physical entities, we recommend that multiple sample records are  
22 registered, one for each biological replicate, with reciprocal references represented as  
23 sample attribute with name 'biological replicate' and attribute value provided as the  
24 accession number(s) of the related biological sample(s). In contrast, 'technical replicates',  
25 for which only a single sample exists, are treated downstream in the workflow.

26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38 Consistent and rich contextual data can become a powerful tool for metagenomics data  
39 analysis. Two marine studies, the *TARA* Oceans sequencing study (PRJEB402, [18]) and  
40 Ocean Sampling Day (OSD, PRJEB5129, [19]) both use the same M2B3 contextual data  
41 reporting standard enabling comparison of data within and across studies. For instance,  
42 data from the *TARA* Oceans shotgun sequencing of the prokaryotic fraction filtered from  
43 seawater (PRJEB1787, [20]) can be compared to the shotgun data from OSD (PRJEB8682,  
44 [21]), enabling detailed or complex queries. Specifically, a taxonomic or functional profile  
45 from the *TARA* Oceans sample from depth 5m and salinity 38psu (SAMEA2591084, [22]) can  
46 be compared to profiles of the OSD sample from the depth 5m (SAMEA3275502, [23]) or the

1 OSD sample with the same salinity 38psu (SAMEA3275531, [24]). In contrast, very few  
2 conclusions can be drawn from a comparison to a sample with insufficient contextual  
3 information (SAMN00194025, [25]).  
4  
5  
6  
7  
8  
9

10 **[Table 3]**  
11

12  
13  
14  
15 Details of the project investigators are usually recorded in the Study metadata object and  
16 sampling contextual data are mostly captured in the Sample metadata object, **Figure 2**. A  
17 common way to standardise reporting of contextual data is via a checklist of key-value pairs,  
18 thereby ensuring parameters of a similar kind are described consistently. Furthermore,  
19 syntactic and semantic rules can be pre-defined in the checklist, enabling validation of  
20 compliance with these rules. For instance, automated checks can be applied to test whether  
21 a mandatory descriptor (key) in the checklist has a value and whether the value is in a  
22 specified format. Each element to be checked can be pre-defined as text, a class or term  
23 from an ontology, a controlled vocabulary or taxonomic index, or formulated as a regular  
24 expression. (Regular expressions can be used, for example, to check that the key 'collection  
25 date and time' complies with ISO 8601 standards and that numeric values lie within a  
26 defined range.)  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47

48 The most common formats for sharing Study and Sample metadata are XML, TSV, ISA-tab or  
49 JSON formats. Examples of the Study and Sample XML are available from the European  
50 Nucleotide Archive [26], [27], where the files are also validated against the XML schema  
51 [28]. Regardless of the format used to supply the metadata, because they all use the same  
52 underlying standards, a simple translation between the formats enables different data to be  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 compared. This allows scientists to use different tools or approaches that they are most  
2 familiar with, whilst ensuring consistent delivery of the metadata.  
3  
4  
5  
6  
7  
8

## 9 Sequencing

10  
11  
12  
13  
14 Once a sample is collected and its provenance recorded, it is subjected to preparation steps  
15 for nucleotide sequence analysis. This may happen immediately after sampling, or in stages  
16 over many months. Processing steps cover all handling of the sample leading to the DNA  
17 isolation. Although MIxS covers some of the metadata fields for reporting the DNA  
18 extraction steps, it is extremely difficult to define a generic set of fields describing the DNA  
19 extraction method with a high granularity due to its complexity and diversity. For example,  
20 it might be relatively straightforward to identify variables for reporting isolation of DNA  
21 from a seawater sample but that will not suit the more complex DNA isolation procedure for  
22 a sediment sample. We suggest the best practice here is to use the existing MIxS fields, such  
23 as the *sample material processing*, *nucleic acid extraction* and *nucleic acid amplification* for  
24 concise description of the nucleic acid preparation. A detailed description, or a reference to  
25 the material preparation steps recorded in a data resource that specialises in protocol  
26 capture and dissemination, such as protocols.io [29], is important due to the significant  
27 influence this can have on the observed profile of the microbial community under  
28 investigation.  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 Equally critical for the downstream metagenomic data analysis and interpretation is the  
2 reporting of sequencing library preparation protocols and parameters as well as sequencing  
3 machine configurations.  
4  
5  
6  
7  
8  
9

10 **Table 4** shows mandatory descriptors for new generation nucleotide sequencing  
11 experiments as currently captured by INSDC databases. **Table 5** lists non-mandatory  
12 descriptors including MIMS sequence-related descriptors and provides our opinion on the  
13 importance of these descriptors for metagenomic data analysis. Note that while a number  
14 of controlled vocabularies have been developed for accurate recording of sequencing  
15 experiment parameters, the evolution of these constrained vocabularies is very dynamic  
16 and driven by technological advances.  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29

30 **[Table 4]**  
31  
32  
33  
34

35 **[Table 5]**  
36  
37  
38  
39

40 Variable parameters of the library preparation and instrumentation are captured in the  
41 metadata objects Experiment and Run (see **Figure 2**). Examples of the Experiment and Run  
42 XML are available, for example, from the ENA [30], [31]. Each Experiment should refer to  
43 Study and Sample objects, to provide context for the sequencing, and is referred to from the  
44 Run objects, which point to the primary sequencing reads.  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54

55 The **primary data** (the reads) are stored in files of various formats, which can be standard  
56 (BAM, CRAM or Fastq) or a platform-specific, as with SFF, PacBio, Oxford Nanopore or  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 Complete Genomics. Information on the read data format must be indicated in the  
2 description of sequencing.  
3  
4  
5  
6

7 The minimum information encapsulated in read data files are base calls with quality scores.  
8  
9 Quality requirements on read data files are file format-specific and are summarised, for  
10 example, in the ENA data submission documentation [32]. A freely available diagnostic tool  
11 for validation of CRAM and BAM files is the Picard ValidateSamFile [33]. Validation of Fastq  
12 files is less straightforward since there is no single FASTQ specification. Recommended  
13 usage of FASTQ can be found for instance in the ENA guidelines [34]. An open resource for  
14 managing next generation sequencing datasets is the NGSUtils [35], which also contains  
15 tools for operations with FASTQ files. As sequencing technologies change over time the  
16 formats and associated validation tools may well change, so a comprehensive list of formats  
17 and tools is likely to become outdated. The key point is to adopt a widely-used format and  
18 to check for file format and integrity (e.g. checksums).  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41

## 42 Analysis

43 Standards in metagenomics for the description of sampling and sequencing have grown out  
44 of those from more traditional genomics. While there are still some shortcomings in these  
45 standards, as highlighted in the previous sections, metadata concerning sampling and  
46 sequencing are commonly captured for metagenomics studies. Compliance is high partly  
47 due to the scientific journals requiring scientist to submit sequence data to an INSDC  
48 database prior to publication. However, there are currently no standards for reporting how  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



1 metagenomics datasets have been analysed. While systematic analysis workflows, such as  
2 those offered by EMG, IMG/M [36] META-pipe [37] and MG-RAST, provide a standard that is  
3 documented (albeit in different ways), many published datasets are analysed by in-house  
4 bespoke pipelines. Although many authors provide an outline in the 'materials and  
5 methods' or 'supplementary materials' section of their publications, it is rarely possible to  
6 reproduce the analysis from this alone, due to missing software parameters, lack of detail  
7 on software versions and ambiguous reference databases and their associated versions.  
8  
9

10 Typically, once the sequence read files have been produced, they are analysed using one or  
11 more workflows [38], with each workflow comprising different data processing or analysis  
12 components. Most workflows involve aspects such as quality control (for example, removing  
13 sequences that fail to meet predefined quality scores), assembly, sequence binning (e.g.  
14 identifying 16S rRNA genes or protein coding sequences), and taxonomic classification of  
15 sequences and/or functional prediction. However, each workflow will be tailored to how the  
16 sample has been processed and the question being addressed. For example, if a sample has  
17 been size-fractionated for viruses, using a 0.22  $\mu\text{m}$  filter, there would be little point analysing  
18 the data for eukaryotic 18S rRNA, as any eukaryotic organisms would have been physically  
19 removed from the sample before the DNA extraction process.  
20  
21

22 Analysis workflows typically have one or more of the following components:  
23  
24

- 25 (i) Central algorithmic software, which may be from a third-party source.
- 26 (ii) 'Glue' software that may ensure input/output formats, or split/join input files for  
27 parallelisation.  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 (iii) Reference datasets that are used by (i). For example, the Greengenes database of 16S  
2 rRNA genes [39], SILVA database of 16S/ 18S SSU ribosomal RNA genes [40] and the NCBI  
3 non-redundant database of non-identical protein sequences [41], for taxonomic or  
4 functional analysis.  
5  
6  
7  
8  
9

10  
11  
12 However, even knowing these elements may not be sufficient for analyses to be  
13 independently re-created. For example, the algorithm may accept a set of input parameters  
14 that can be used to fine-tune an analysis, such as selecting an E-value threshold for  
15 determining significance of a sequence match to a reference database. Other parameters  
16 may influence speed-performance, which allows the original analysis to complete in a timely  
17 fashion, but they may or may not have an effect on the results. For example, running  
18 *hmmsearch* from the HMMER package, changing the number of CPUs used will not change  
19 the results, but changing options on the heuristics such as the --F1 threshold (which controls  
20 the number of sequences passing the first heuristic state) may alter the output; both will  
21 potentially increase performance in terms of speed.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40

41 Capturing and reporting *all* provenance information is essential to understand exactly what  
42 analysis has been performed on the data, and to ensure reproducibility [42]. The use of  
43 publicly available analysis pipelines (such as EMG, MG-RAST or META-pipe) helps with this  
44 process, since analysis is performed using pre-defined components, settings and databases  
45 (or, in some cases, using user-selected components, selected from a predefined list of  
46 options). Nevertheless, capturing analysis metadata remains essential as, for example, MG-  
47 RAST allows the users to dynamically set E-value thresholds after the pipeline analysis has  
48 been performed. Furthermore, the tools, libraries, and reference databases used by the  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 pipelines are regularly updated, and thus capturing analysis provenance information is  
2 vitally important and should be systematically ‘tagged’ to the results.  
3  
4  
5  
6

7 To date, there is no universally endorsed ‘Analysis standard’ for describing and recreating a  
8 metagenomics analysis pipeline, and without this standard (and subsequent  
9 adoption/enforcement) it will continue to be difficult, if not impossible, to reproduce  
10 analysis workflows. However, all is not lost. ‘Workflows’ and their definitions is an active  
11 field of computer science research, and potential solutions are already available, including  
12 Common Workflow Language (CWL), Yet Another Workflow Language (YAWL), Business  
13 Process Execution Language (BPEL) and Microsoft Azure’s Workflow Definition Language to  
14 name but a few. Several of the co-authors for this publication already participate in the GSC  
15 M5 consultation group, which aims to define a standard enabling the recreation and  
16 exchange of metagenomics data sets. In the absence of a standard, we believe it is  
17 important to define some of the basic best practices, from which an accepted standard  
18 would formally encapsulate. For simplicity, we will focus on a single ‘best practice’ use case:  
19 the description of the analysis of a run. Other types of analysis, such as pooling of runs or  
20 comparing results between runs are beyond the scope of this article.  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45

46 A schematic overview of a best practice for analysis metadata collection is shown in **Figure**  
47 **3A**. An overarching set of metadata relating to analysis should encapsulate generic  
48 information such as analysis centre, name of bioinformaticians, analysis objectives, name of  
49 overall analysis (if appropriate) and the date on which the analysis was performed. It should  
50 also contain appropriate pointers to the run data, run sequence metadata and associated  
51 sample data. Underneath the overarching analysis metadata is a collection of analysis  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 components, which describe each stage of the analysis **Figure 3B**. Each component can be  
2 divided into three sections: input(s), analysis algorithm and output(s).  
3  
4

5  
6  
7 The input section should describe the details of the various inputs to the analysis, which  
8 could be the raw sequence reads or the output of another analysis component, reference  
9 databases and their provenance data, such as version, where necessary. The analysis section  
10 should contain the algorithm tool, version, all parameters used and a basic description of  
11 the analysis. The output section should describe each output from the analysis, together  
12 with a description of contents and format.  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

23  
24  
25 Each analysis component could then be coupled to form an analysis workflow as shown in  
26 **Figure 3C**. The workflow may be in a portable intermediate format that can be submitted to  
27 a workflow manager for execution in a specific environment.  
28  
29  
30  
31  
32

33  
34  
35 **[Figure 3]**  
36

37  
38  
39  
40 This best practice framework is merely that - a best practice, and we have not touched on  
41 the technical issues of how to capture this information or on controlled vocabularies (since  
42 these need to come from the community). Furthermore, enforcing compliance and  
43 validation against the standard will also require a community effort. Complete validation  
44 would require the standard to be machine readable and deployable, with potentially the  
45 need to have small 'test' datasets and their associated results, to perform regression testing  
46 of the analysis metadata. However, who is responsible for validation and what happens if  
47 something fails after publication are open questions. This could arguably be a step too far;  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 currently sampling and sequence metadata are validated against the standard, but taken in  
2 good faith to be correct beyond this.  
3  
4  
5  
6  
7  
8  
9

## 10 Analysis Results Archiving - a final piece? 11

12 Having an analysis provenance standard would allow metagenomics analysis results to be  
13 recreated more readily. While this is undoubtedly an important and necessary step, it has a  
14 major limitation within the community. As indicated [43], the fraction of money spent on  
15 informatics from an overall project budget is increasing dramatically. Metagenomics  
16 datasets tend to be large, in the order of GB-TB and processing may take 1000s of CPU  
17 hours, restricting reanalysis to only those with significant compute resources. For example,  
18 the subset of the *TARA* Oceans Ocean Microbiome Project (PRJEB7988, [44]) that has been  
19 size-fractionated for prokaryotes comprises 135 samples with 248 runs containing 28.8 billion  
20 reads. The analysis output represents about 10TB of data with 23.2 billion predicted protein  
21 coding sequences. Thus, reanalysis would be costly and potentially wasteful if a particular  
22 workflow had already been run on the data. Therefore, a final step in a metagenomics  
23 analysis is the appropriate archiving of results. There is an obvious cost-benefit balance to  
24 be drawn here, as storing every intermediate of a workflow would lead to an explosion of  
25 data. Clearly, key intermediates and outputs of an analysis workflow need to be determined.  
26 These key archived components will be tailored to the analysis, but should at least include  
27 operational taxonomic unit (OTU) counts and assignments, functional assignment counts  
28 and read/sequence positional information for the aforementioned assignments. Such data  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 files are already made available from MG-RAST and EMG and those from other sources are  
2  
3 accepted for archiving within ENA.  
4

5  
6  
7 If metagenomic assemblies have been performed, then these should have an appropriate  
8  
9 structure of contigs, scaffolds or chromosomes with an appropriate format as detailed for  
10  
11 example in the ENA data submission documentation [45]. Due to the overheads of  
12  
13 producing an assembly, these should be archived, ideally with an INSDC database.  
14  
15

16  
17  
18 The data model for metagenomics, as described in **Figure 2**, represents metagenomic  
19  
20 analysis results in the data Analysis object with appropriate pointers to the corresponding  
21  
22 run sequence metadata and associated sample collection contextual data. While there is an  
23  
24 established practice to archive primary sequence data in the Run object and assemblies of  
25  
26 the primary sequences in the Analysis object, it is not a common practice to archive results  
27  
28 of functional and taxonomic metagenomic analysis of in-house bespoke pipelines. It would  
29  
30 be beneficial to the metagenomics community to include this into the best practice and such  
31  
32 data are accepted by ENA for archiving. The metagenomics standard environment reviewed  
33  
34 here as well as outcomes of the GSC M5 consultation group can contribute to defining  
35  
36 required descriptors of the Analysis object for archiving of metagenomics analysis results,  
37  
38 which can serve as a framework for exchange of metagenomics data sets on a routine basis,  
39  
40 similarly as is currently done for the primary sequence data.  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

# Future

One challenge over the next several years will be the validation of compliance across the entirety of the standards and best practise that we have covered. While validation tools and recommended practises exist for parts (e.g. contextual data descriptors using MIxS-compliant validation tools from ISA and experimental descriptors upon submission to an INSDC database), not all parts have such maturity (e.g. analysis descriptors) and there exists no overarching validation protocol for an entire metagenomics study. The GSC is aiming to contribute in this area with the introduction of MIxS “profiles”, to provide an overlay on top of MIxS environmental packages and the core MIxS fields. These profiles will enable the creation of tool suites for compliance checking. In addition (and perhaps more importantly) they will enable groups of researchers, institutes, funders and other communities to define levels of compliance for contextual data sets. Examples of this are the NSF NEON [46] and the NSF CZO [47] networks that are working with the GSC to establish silver, gold and platinum sets of parameters that need to be provided and validated for data sets to be compliant. A key moment in the acceptance of said new profiles will be the availability of tools support for data creators, end-users and portals. Imagine, for instance, a search for data sets in EMG or MG-RAST that allows restriction of the search to just platinum-level data sets. For the data consumer this will result in better ways of telling their science story using third party data and for the data creator this will provide guidance on what to create for a specific community. In addition funding agencies can require certain minimal compliance.

A further area to be addressed is that of standards around descriptions of metagenomics analyses. The creation of a lightweight data standard for an analysis object that allows easy

1 transfer of analyses is a key goal of the GSC M5 initiative but the complexity of the task and  
2 lack of dedicated resourcing has rendered progress slow; while frameworks and systems for  
3 recording analysis provenance need to be established, we have aimed to indicate in this  
4 publication a set of best practice that can form the foundation for a community standard  
5 enabling the recreation and exchange of metagenomics data sets. Improving  
6 standardisation will also help raise clarity in the literature around metagenomics through a  
7 tightening of language. For example, UProC [48] uses Pfam [49] matches as a reference  
8 library with the results being referred to as a 'Pfam hit'. However, this may not necessarily  
9 be a Pfam hit, as a Pfam hit is defined as a sequence match scoring greater than Pfam  
10 defined threshold to the Pfam profile Hidden Markov Model (i.e. the Pfam database  
11 method).

## 32 Conclusions

33 In this overview of the metagenomics standard environment we have outlined best practice  
34 for the reporting of metagenomics workflows. We have reviewed the essential steps: (1)  
35 material sampling, (2) material sequencing (3) data analysis and (4) data archiving, and  
36 highlighted essential variable parameters and common data formats in each step.

37 Reporting on the provenance of a sample and associated nucleotide sequence data is largely  
38 established by public sequence data repositories and is also being addressed by contextual  
39 data standardisation initiatives. In contrast, a reporting standard on metagenomics data  
40 analysis is absent, yet the high complexity of metagenomics creates a pressing demand for



1 establishing such a practice. Capturing key metadata relating to analysis would greatly  
2 improve reproducibility. Archiving key results of the metagenomics data analysis would  
3 allow a more accurate evaluation of the benefits of reproducing the analysis.  
4

5  
6  
7 Only by adopting these standards and best practices can metagenomics data be assessed  
8 against the FAIR (Findable, Accessible, Interoperable and Reusable) principles that should be  
9 applied to any scientific dataset [50].  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

## 22 List of abbreviations

23  
24  
25  
26  
27 BAM: Binary Alignment/Map

28  
29 BPEL: Business Process Execution Language

30  
31  
32 CDI: Common Data Index

33  
34 CHEBI: Chemical Entities of Biological Interest

35  
36  
37 CNR: Consiglio Nazionale delle Ricerche

38  
39 CPU: Central Processing Unit

40  
41  
42 CRAM: Compression Reduced Alignment/Map

43  
44  
45 CWL: Common Workflow Language

46  
47 EBI: European Bioinformatics Institute

48  
49  
50 EMBL: European Molecular Biology Laboratory

51  
52  
53 EMG: EBI Metagenomics

54  
55  
56 ENA : European Nucleotide Archive

57  
58 ENVO: Environment Ontology  
59  
60  
61  
62  
63  
64  
65

1 GEO BON: Group on Earth Observations Biodiversity Observation Network

2 GO: Gene Ontology

3 GSC: Genomic Standards Consortium

4 GSC M5: Genomic Standards Consortium M5 working group

5 IMG/M: Integrated Microbial Genomes with Microbiomes

6 INSDC: International Nucleotide Sequence Database Collaboration

7 ISO: International Standards Organisation

8 JSON: JavaScript Object Notation

9 MG-RAST: Metagenomic Rapid Annotations using Subsystems Technology

10 MIMS: Minimum Information about a Metagenome Sequence

11 MIMARKS: Minimum Information about a MARKer gene Sequence

12 MixS: Minimum Information about any (x) Sequence

13 M2B3: Marine Microbial Biodiversity Bioinformatics and Biotechnology

14 NCBI: National Center for Biotechnology Information

15 NSF NEON: National Science Foundation National Ecological Observatory Network

16 NSF CZO: National Science Foundation Critical Zone Observatory

17 OBIS: Ocean Biogeographic Information System

18 OSD: Ocean Sampling Day

19 OTU: Operational Taxonomic Unit

20 PSI: Proteomics Standards Initiative

21 SDN: SeaDataNet

22 SFF: Standard Flowgram Format

23 SSU: Small Subunit ribosomal RNA

24 TB: Terabyte

1 TSV: Tab Separated Values

2 UiT: Universitetet i Tromsø

3 UProC: Ultrafast Protein Classification

4 UTC: Coordinated Universal Time

5 XML: Extensible Markup Language

6 YAWL: Yet Another Workflow Language

## 7

## 8

## 9

## 10

## 11

## 12

## 13

## 14

## 15

## 16

## 17

## 18

## 19

## 20

## 21

## 22

## 23

## 24

## 25

## 26

## 27

## 28

## 29

## 30

## 31

## 32

## 33

## 34

## 35

## 36

## 37

## 38

## 39

## 40

## 41

## 42

## 43

## 44

## 45

## 46

## 47

## 48

## 49

## 50

## 51

## 52

## 53

## 54

## 55

## 56

## 57

## 58

## 59

## 60

## 61

## 62

## 63

## 64

## 65

# Declarations

### Ethics approval and consent to participate

Not applicable

### Consent for publication

Not applicable

### Competing interests

The authors declare that they have no competing interests.

### Funding

This work was supported by the ELIXIR-EXCELERATE funded by the European Commission within the Research Infrastructures programme of Horizon 2020, grant agreement number 676559.

## Author's contribution

PtH and GC conceived the study; PtH and RDF drafted the manuscript; all authors contributed further and revised. All authors have approved the final manuscript.

## References

- [1] Salter SJ, Cox MJ, Turek EM, Calus ST, Cookson WO, Moffatt MF, et al. Reagent and laboratory contamination can critically impact sequence-based microbiome analysis. *BMC Biology*. 2014;12:87.
- [2] Toribio AL, Alako B, Amid C, Cerdeño-Tárraga A, Clarke L, Cleland I, et al. European Nucleotide Archive in 2016. *Nucleic Acids Res*. 2016; doi:10.1093/nar/gkw1106.
- [3] Mitchell A, Bucchini F, Cochrane G, Denise H, Hoopen PT, Fraser M, et al. EBI metagenomics in 2016 - an expanding and evolving resource for the analysis and archiving of metagenomic data. *Nucleic Acid Res*. 2016;44:D595-D603.
- [4] Meyer F, Paarmann D, D'Souza M, Olson R , Glass EM, Kubal M, et al. The Metagenomics RAST server – a public resource for the automatic phylogenetic and functional analysis of metagenomes. *BMC Bioinformatics*. 2008;9:386.

1 [5] Field D, Amaral-Zettler L, Cochrane G, Cole JR, Dawyndt P, Garrity GM et al. The  
2 Genomic Standards Consortium. PLoS Biol. 2011;9(6):e1001088.  
3

4  
5  
6  
7 [6] Orchard S., Hermjakob H. and Apweiler R. (2003) The Proteomics Standards  
8 Initiative. Proteomics. 2003;3(7):1374-6.  
9

10  
11  
12  
13  
14  
15 [7] The Group on Earth Observations Biodiversity Observation Network.  
16 <http://geobon.org/>, accessed 5. Jun 2017.  
17

18  
19  
20  
21  
22 [8] Bruford MW, Davies N, Dooloo ME, Faith DP, Walters M: Monitoring Changes in  
23 Genetic Diversity. In The GEO Handbook on Biodiversity Observation Networks. Edited by  
24 Walters M, Scholes RJ: Springer; 2017: 107-128  
25  
26  
27  
28  
29

30  
31  
32  
33 [9] The Minimum Information about a Metagenome Sequence.  
34 <http://wiki.gensc.org/index.php?title=MIGS/MIMS>. Accessed 19 Jan 2017.  
35  
36  
37  
38  
39

40  
41 [10] The Minimum Information about a Marker Gene Sequence.  
42 <http://wiki.gensc.org/index.php?title=MIMARKS>. Accessed 19 Jan 2017.  
43  
44  
45  
46  
47

48  
49 [11] Yilmaz, P, Gilbert, JA, Knight, R, Amaral-Zettler L, Karsch-Mizrachi I, Cochrane G et al.  
50 The Genomic Standards Consortium: bringing standards to life for microbial ecology. The  
51 ISME J. 2011;5,1565-7.  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 [12] The Minimum Information about any (x) Sequence, Version 4.0. [http:// wiki. genesc.](http://wiki.genesc.org/index.php?title=MlxS)  
2 [org/ index. php? title= MlxS](http://wiki.genesc.org/index.php?title=MlxS). Accessed 19 Jan 2017.  
3  
4

5  
6  
7 [13] Cochrane G, Karsch-Mizrachi I, Takagi T and International Nucleotide Sequence  
8 Database Collaboration. The International Nucleotide Sequence Database Collaboration.  
9  
10 Nucleic Acid Res. 2016;44:D40-D50.  
11  
12  
13

14  
15  
16  
17 [14] The ISA framework and tools. <http://isa-tools.org/>. Accessed 7 February 2017.  
18  
19  
20

21  
22 [15] Ten Hoopen P, Pesant S, Kottman R, Kopf A, Bickel M, Claus S et al. Marine microbial  
23 biodiversity, bioinformatics and biotechnology (M2B3) data reporting and service standards.  
24  
25 The Standards in Genomic Sciences 2015;10:20.  
26  
27  
28

29  
30  
31 [16] The Common Data Index, Version 3.0. [http:// www. seadatanet. org/ Data-Access/](http://www.seadatanet.org/Data-Access/Common-Data-Index-CDI)  
32 [Common-Data-Index-CDI](http://www.seadatanet.org/Data-Access/Common-Data-Index-CDI). Accessed 19 Jan 2017.  
33  
34  
35  
36

37  
38 [17] The Ocean Biogeographic Information System data standard, Version 1.1.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

66 [18] The *TARA* Oceans umbrella project record of barcoding and shotgun sequencing.  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100  
<http://www.ebi.ac.uk/ena/data/view/PRJEB402>. Accessed 7 Feb 2017.

101 [19] The Ocean Sampling Day umbrella project record of amplicon and metagenome  
102 sequencing. <http://www.ebi.ac.uk/ena/data/view/PRJEB5129>. Accessed 7 Feb 2017.  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165

1  
2 [20] The record of the *TARA* Oceans shotgun sequencing project of the prokaryotic  
3 fraction filtered from seawater. <http://www.ebi.ac.uk/ena/data/view/PRJEB1787>. Accessed  
4  
5 7 Feb 2017.  
6  
7

8  
9  
10  
11  
12 [21] The record of the Ocean Sampling Day shotgun sequencing project from the year  
13 2014. <http://www.ebi.ac.uk/ena/data/view/PRJEB8682>. Accessed 7 Feb 2017.  
14  
15  
16  
17

18  
19  
20 [22] The record of a *TARA* Oceans sample from depth 5m and salinity 38psu.  
21  
22 <https://www.ebi.ac.uk/metagenomics/projects/ERP001736/samples/ERS477979>. Accessed  
23  
24 7 Feb 2017.  
25  
26  
27

28  
29  
30 [23] The record of an Ocean Sampling Day sample from depth 5m.  
31  
32 <https://www.ebi.ac.uk/metagenomics/projects/ERP009703/samples/ERS667511>. Accessed  
33  
34 7 Feb 2017.  
35  
36  
37

38  
39  
40 [24] The record of an Ocean Sampling Day sample with salinity 38psu.  
41  
42 <https://www.ebi.ac.uk/metagenomics/projects/ERP009703/samples/ERS667548>. Accessed  
43  
44 7 Feb 2017.  
45  
46  
47

48  
49  
50 [25] The record of an oil spill water sample from Gulfport.  
51  
52 <http://www.ebi.ac.uk/ena/data/view/SAMN00194025>. Accessed 7 Feb 2017.  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 [26] An example of a study XML. [http://www.ebi.ac.uk/ena/submit/preparing-](http://www.ebi.ac.uk/ena/submit/preparing-xmls#study)  
2 [xmls#study](http://www.ebi.ac.uk/ena/submit/preparing-xmls#study). Accessed 19 Jan 2017.  
3

4  
5  
6  
7 [27] An example of a sample XML. [http://www.ebi.ac.uk/ena/submit/preparing-](http://www.ebi.ac.uk/ena/submit/preparing-xmls#sample)  
8 [xmls#sample](http://www.ebi.ac.uk/ena/submit/preparing-xmls#sample). Accessed 19 Jan 2017.  
9

10  
11  
12  
13  
14  
15 [28] The validating XMLs document. <http://www.ebi.ac.uk/ena/submit/validating-xmls>.  
16  
17 Accessed 19 Jan 2017.  
18

19  
20  
21  
22  
23 [29] Protocols.io. <https://www.protocols.io/> accessed 5 Jun 2017  
24

25  
26  
27  
28 [30] An example of an experiment XML. [http://www.ebi.ac.uk/ena/submit/preparing-](http://www.ebi.ac.uk/ena/submit/preparing-xmls#experiment)  
29 [xmls#experiment](http://www.ebi.ac.uk/ena/submit/preparing-xmls#experiment). Accessed 19 Jan 2017.  
30

31  
32  
33  
34  
35  
36 [31] An example of a run XML. <http://www.ebi.ac.uk/ena/submit/preparing-xmls#run>.  
37  
38 Accessed 19 Jan 2017.  
39

40  
41  
42  
43 [32] The document on the ENA-supported read file formats.  
44  
45 <http://www.ebi.ac.uk/ena/submit/read-file-formats>. Accessed 19 Jan 2017.  
46  
47

48  
49  
50  
51 [33] The document on the Picard set of command line tools.  
52  
53 <http://broadinstitute.github.io/picard/>. Accessed 19 Jan 2017.  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



1 [34] The document on recommended usage of FASTQ files.

2 <http://www.ebi.ac.uk/ena/submit/read-data-format>. Accessed 7 February 2017.

3  
4  
5  
6  
7 [35] The document on the NGSUtils tools for next-generation sequencing analysis.

8 <http://ngsutils.org/>. Accessed 19 Jan 2017.

9  
10  
11  
12  
13  
14 [36] Markowitz VM, Chen IM, Chu K, Szeto E, Palaniappan K, Pillay M, Ratner A et al.

15 IMG/M 4 version of the integrated metagenome comparative analysis system. Nucleic Acid  
16 Res. 2015;42:D568-73.

17  
18  
19  
20  
21 [37] Robertsen EM, Kahlke T, Raknes IA, Pedersen E, Semb EK, Ernstsens M et al. META-

22 pipe – pipeline annotation, analysis and visualisation of marine metagenomic sequence

23 data. <https://arxiv.org/abs/1604.04103v1>. Accessed 19 Jan 2017.

24  
25  
26 [38] Leipzig J. A review of bioinformatic pipeline frameworks. Briefings in Bioinformatics.

27 2016; doi:[10.1093/bib/bbw020](https://doi.org/10.1093/bib/bbw020).

28  
29  
30  
31  
32  
33  
34 [39] The Greengenes Database. <http://greengenes.secondgenome.com/>. Accessed 19 Jan

35 2017.

36  
37  
38  
39  
40  
41 [40] The Silva Database. <https://www.arb-silva.de/>. Accessed 19 Jan 2017.

42  
43  
44  
45  
46  
47 [41] The NCBI nr Database. <https://www.ncbi.nlm.nih.gov/refseq/>. Accessed 19 Jan 2017.

1 [42] Huntemann M, Ivanova NN, Mavromatis K, Tripp HJ, Paez-Espino D, Tennessen K et  
2 al. The standard operating procedure of the DOE-JGI Metagenome Annotation Pipeline  
3 (MAPv v4.). Standards in Genomic Sciences. 2016; doi:10.1186/s40793-016-0138-x.  
4  
5  
6

7  
8  
9  
10 [43] Sboner A, Mu XJ, Greenbaum D, Auerbach RK and Gerstein MB. The real costs of  
11 sequencing: higher than you think! Genome Biology. 2011;12:125.  
12  
13  
14

15  
16  
17  
18 [44] The record of the TARA Oceans Ocean Microbiome Project.  
19  
20 <http://www.ebi.ac.uk/ena/data/view/PRJEB7988>. Accessed 7 Feb 2017.  
21  
22  
23

24  
25  
26 [45] The document on a genome assembly submission to the ENA.  
27  
28 <http://www.ebi.ac.uk/ena/submit/genomes-sequence-submission>. Accessed 19 Jan 2017.  
29  
30  
31

32  
33  
34 [46] The National Science Foundation National Ecological Observatory Network.  
35  
36 <http://www.neonscience.org/>. Accessed 19 Jan 2017.  
37  
38  
39

40  
41 [47] The National Science Foundation Critical Zone Observatory.  
42  
43 <http://criticalzone.org/national/>. Accessed 19 Jan 2017.  
44  
45  
46

47  
48  
49 [48] Meinicke P. UProC: tools for ultra-fast protein domain classification. Bioinformatics.  
50 2014;31(9):1382-8.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 [49] Finn RD, Coggill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL et al. The Pfam  
2 protein families database: towards a more sustainable future. Nucleic Acid Res.  
3  
4  
5 2016;44:D279-85.  
6

7  
8  
9  
10 [50] Wilkinson MD, Dumontier M, Aalbersberg IJJ, Appleton G, Axton M, Baak A et al. The  
11 FAIR Guiding Principles for scientific data management and stewardship. Sci. Data 2016;  
12  
13 doi:10.1038/sdata.2016.18 (2016).  
14  
15

16  
17  
18  
19  
20 [51] The controlled vocabulary for the INSDC country qualifier.  
21  
22 <http://www.insdc.org/country.html>. Accessed 19 Jan 2017.  
23  
24

25  
26  
27  
28 [52] The Environment Ontology browser. [http://www.environmentontology.org/Browse-  
29 EnvO](http://www.environmentontology.org/Browse-EnvO). Accessed 19 Jan 2017.  
30  
31

32  
33  
34  
35  
36 [53] The SeaDataNet L06 controlled vocabulary of platform categories.  
37  
38 [http://seadatanet.maris2.nl/v\\_bodc\\_vocab\\_v2/search.asp?lib=L06](http://seadatanet.maris2.nl/v_bodc_vocab_v2/search.asp?lib=L06). Accessed 19 Jan 2017.  
39  
40

41  
42  
43  
44 [54] The SeaDataNet P02 controlled vocabulary of parameters.  
45  
46 [http://seadatanet.maris2.nl/v\\_bodc\\_vocab\\_v2/search.asp?lib=P02](http://seadatanet.maris2.nl/v_bodc_vocab_v2/search.asp?lib=P02). Accessed 19 Jan 2017.  
47  
48

49  
50  
51 [55] The controlled vocabulary of BODC data storage units.  
52  
53 [http://seadatanet.maris2.nl/v\\_bodc\\_vocab\\_v2/search.asp?lib=P06](http://seadatanet.maris2.nl/v_bodc_vocab_v2/search.asp?lib=P06). Accessed 19 Jan 2017.  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 [56] The CHEBI ontological classification of small chemical compounds.

2 <http://www.ebi.ac.uk/chebi/init.do>. Accessed 19 Jan 2017.

3  
4  
5  
6  
7 [57] The National Center for Biotechnology Information taxonomy index.

8 <https://www.ncbi.nlm.nih.gov/taxonomy>. Accessed 19 Jan 2017.

9  
10  
11  
12  
13  
14 [58] The controlled vocabulary for sequencing instrument models.

15 <http://www.ebi.ac.uk/ena/submit/preparing-xm1s#experiment>. Accessed 19 Jan 2017.

16  
17  
18  
19  
20  
21 [59] The controlled vocabulary for the library source.

22 <http://www.ebi.ac.uk/ena/submit/preparing-xm1s#experiment>. Accessed 19 Jan 2017.

23  
24  
25  
26  
27 [60] The controlled vocabulary for the library strategy.

28 <http://www.ebi.ac.uk/ena/submit/preparing-xm1s#experiment>. Accessed 19 Jan 2017.

29  
30  
31  
32  
33 [61] The controlled vocabulary for the library selection.

34 <http://www.ebi.ac.uk/ena/submit/preparing-xm1s#experiment>. Accessed 19 Jan 2017.

35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55 **Figure Legends**  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 **Figure 1:** A generalised metagenomics data analysis workflow in the context of other ‘omics’  
2  
3 approaches.

4  
5  
6  
7 **Figure 2:** A common data model for read data and associated metadata.  
8  
9

10  
11  
12 **Figure 3:** Schematic overview of best practice for analysis metadata collection with example  
13  
14 fields.  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55

56 **Tables**  
57  
58  
59  
60  
61  
62  
63  
64  
65

**Table 1:** Checklist of MIMS mandatory descriptors for a sample taken from an aquatic environment and associated with a metagenomic sequencing experiment.

MIMS-mandatory water sample provenance descriptors	descriptor format
submitted to INSDC	boolean
project name	text
investigation type	fixed value: 'metagenome'
geographic location (latitude and longitude)	decimal degrees in WGS84 system
depth	metres: positive below the sea surface
geographic location (country and/or sea region)	INSDC country list [51]
collection date	ISO8601 date and time
environment (biome)	ENVO class [52]
environment (feature)	ENVO class
environment (material)	ENVO class
environment package	MIxS controlled vocabulary [12]

**Table 2:** Checklist of M2B3 mandatory descriptors for a microbial sample taken from a saline water environment and associated with a metagenomic sequencing experiment.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

M2B3-mandatory saline water sample provenance descriptors	descriptor format
INVESTIGATION_campaign	text
INVESTIGATION_site	text
INVESTIGATION_platform	SDN:L06 controlled vocabulary [53]
EVENT_latitude	decimal degrees in WGS84 system
EVENT_longitude	decimal degrees in WGS84 system
EVENT_date/time	ISO8601 date and time in UTC
SAMPLE_title	text
SAMPLE_protocol label	text
SAMPLE_depth	metres; positive below the sea surface
ENVIRONMENT_environment (biome)	ENVO class
ENVIRONMENT_environment (feature)	ENVO class
ENVIRONMENT_environment (material)	ENVO class
ENVIRONMENT_temperature	SDN:P02 [54], SDN:P06 [55] controlled vocab.
ENVIRONMENT_salinity	SDN:P02, SDN:P06 controlled vocab.

**Table 3:** Selection of non-mandatory MlxS and M2B3 descriptors (column B) and formats (column D). These descriptors cover such areas as the structure or viability of the community under investigation and sample pooling procedures. Column A groups descriptors that are related conceptually (1 – sample collection method & device, 2 – sample processing, 3 – sample quantity, 4 – storage container, 5 – storage duration, 6 – storage temperature, 7 – chemical treatment, 8 – microbial fraction thresholds, 9 – sample content, 10 – pigment concentration, 11 – fluorescence, 12 – density, 13 – organism abundance, 14 – primary production, 15 – bacterial production, 16 – organism biomass, 17 – organism biovolume, 18 – organism size, 19 – investigation contributors, 20 – unique taxonomic index identifier for organism host). Column C shows the descriptor association with the respective contextual data reporting standard suitable for marine metagenomic data. Column E suggests the descriptor’s importance for metagenomic data analysis (H – high relevance, M – medium relevance, L – low relevance).

A group	B non-mandatory sample provenance descriptors	C standard	D descriptor format	E value for analysis (H/M/L)
1	sample collection device or method	MlxS(MIMS)	text	H
1	EVENT_device	M2B3	text	H
1	EVENT_method	M2B3	text	H
2	sample material processing	MlxS(MIMS)	text	H
3	amount or size of sample collected	MlxS(MIMS)	numeric & unit	H
3	SAMPLE_quantity (e.g. length, mass)	M2B3	text	H



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

4	sample storage location	MlxS(water)	text	L
4	SAMPLE_container (e.g. storage container)	M2B3	text	L
5	sample storage duration	MlxS(water)	interval	H
6	sample storage temperature	MlxS(water)	numeric & unit	H
6	SAMPLE_treatment_storage (e.g. temperature)	M2B3	text	H
7	chemical administration	MlxS(water)	CHEBI ontology [56]	M
7	SAMPLE_treatment_chemicals	M2B3	CHEBI ontology	M
8	SAMPLE_size_fraction_upper_threshold	M2B3	text	H
8	SAMPLE_size_fraction_lower_threshold	M2B3	text	H
9	SAMPLE_content (e.g. 0.22 µm filter, 20mL water)	M2B3	text	H
10	concentration of chlorophyll	MlxS(water)	numeric & unit	HM
10	ENVIRONMENT_ecosystem_pigment concentration	M2B3	SDN:P02, SDN:P06 controlled vocab.	HM
11	Fluorescence	MlxS(water)	numeric & unit	HM
11	ENVIRONMENT_ecosystem_fluorescence	M2B3	SDN:P02, SDN:P06 controlled vocab.	HM
12	density	MlxS(water)	numeric & unit	M
13	organism count	MlxS(water)	numeric & unit	ML
13	ENVIRONMENT_ecosystem_picoplankton (flow cytometry) abundance	M2B3	SDN:P02, SDN:P06 controlled vocab.	ML

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

13	ENVIRONMENT_ecosystem_nano/microplankton abundance	M2B3	SDN:P02, SDN:P06 controlled vocab.	ML
13	ENVIRONMENT_ecosystem_meso/macroplankton abundance	M2B3	SDN:P02, SDN:P06 controlled vocab.	ML
14	primary production	MixS(water)	numeric & unit	M
14	ENVIRONMENT_ecosystem_primary production	M2B3	SDN:P02, SDN:P06 controlled vocab.	M
15	bacterial production	MixS(water)	numeric & unit	M
15	ENVIRONMENT_ecosystem_bacterial production	M2B3	SDN:P02, SDN:P06 controlled vocab.	M
16	biomass	MixS(water)	numeric & unit	ML
16	ORGANISM_biomass	M2B3	numeric & unit & method	ML
17	ORGANISM_biovolume	M2B3	numeric & unit & method	L
18	ORGANISM_size	M2B3	numeric & unit & method	L
19	INVESTIGATION_authors	M2B3	text	M
20	host taxid	MixS (host associated)	NCBI Taxonomy identifier [57]	M

**Table 4:** Mandatory descriptors for sequencing.

mandatory descriptors of sequencing provenance	descriptor format
instrument platform	controlled vocabulary [illumina, oxford nanopore, pacbio smrt, ion torrent, ls454, complete genomics, capillary]
instrument model	controlled vocabulary [58]
library source	controlled vocabulary [59]
library strategy	controlled vocabulary [60]
library selection	controlled vocabulary [61]
library layout	controlled vocabulary [single, paired]
read file name	text
read file md5 checksum	32-digit hexadecimal number
second read file name (for paired Fastq files)	text
Second read file md5 checksum (for paired Fastq files)	32-digit hexadecimal number

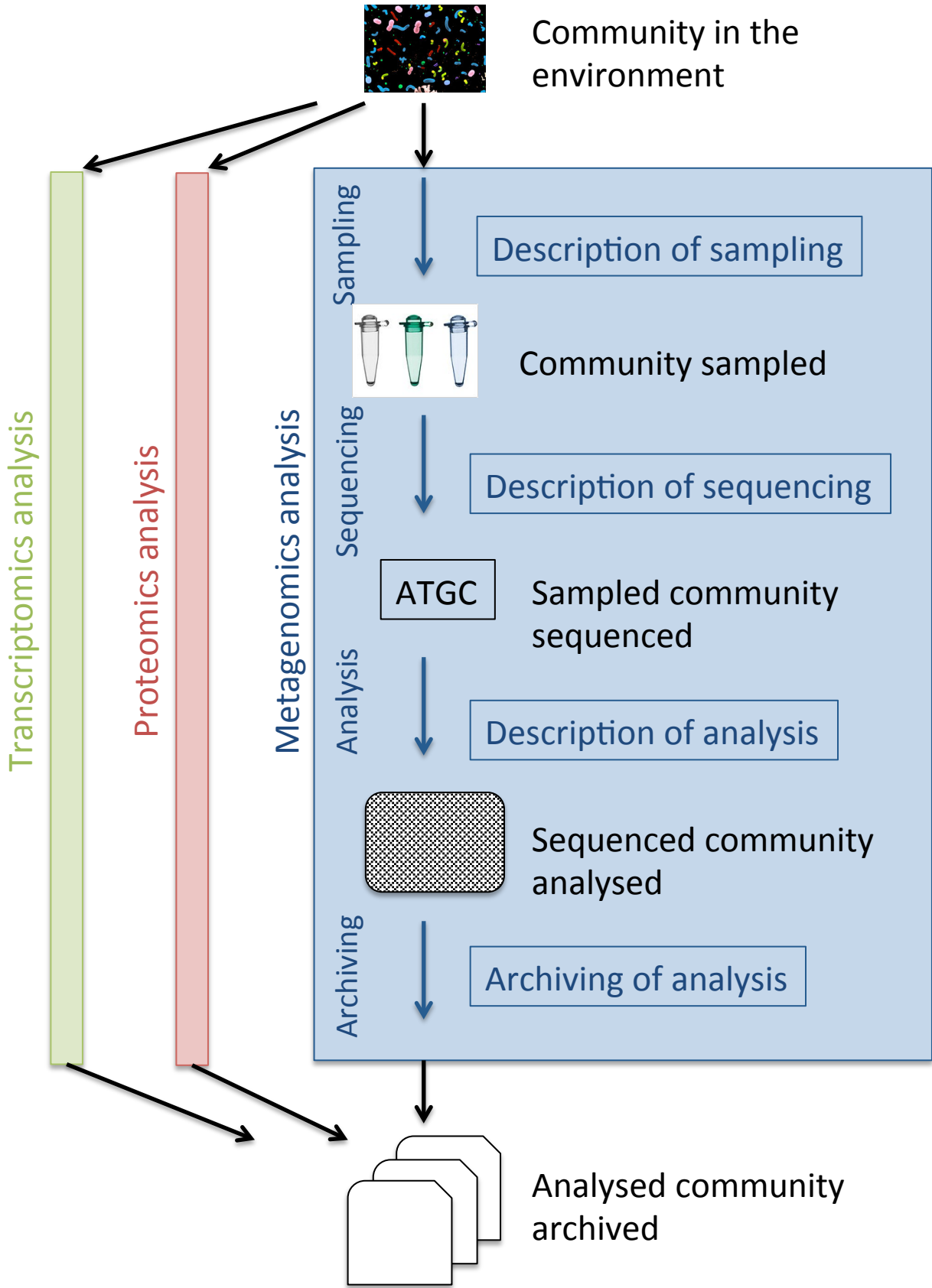
**Table 5:** Non-mandatory sequencing descriptors (column A) and formats (column B).

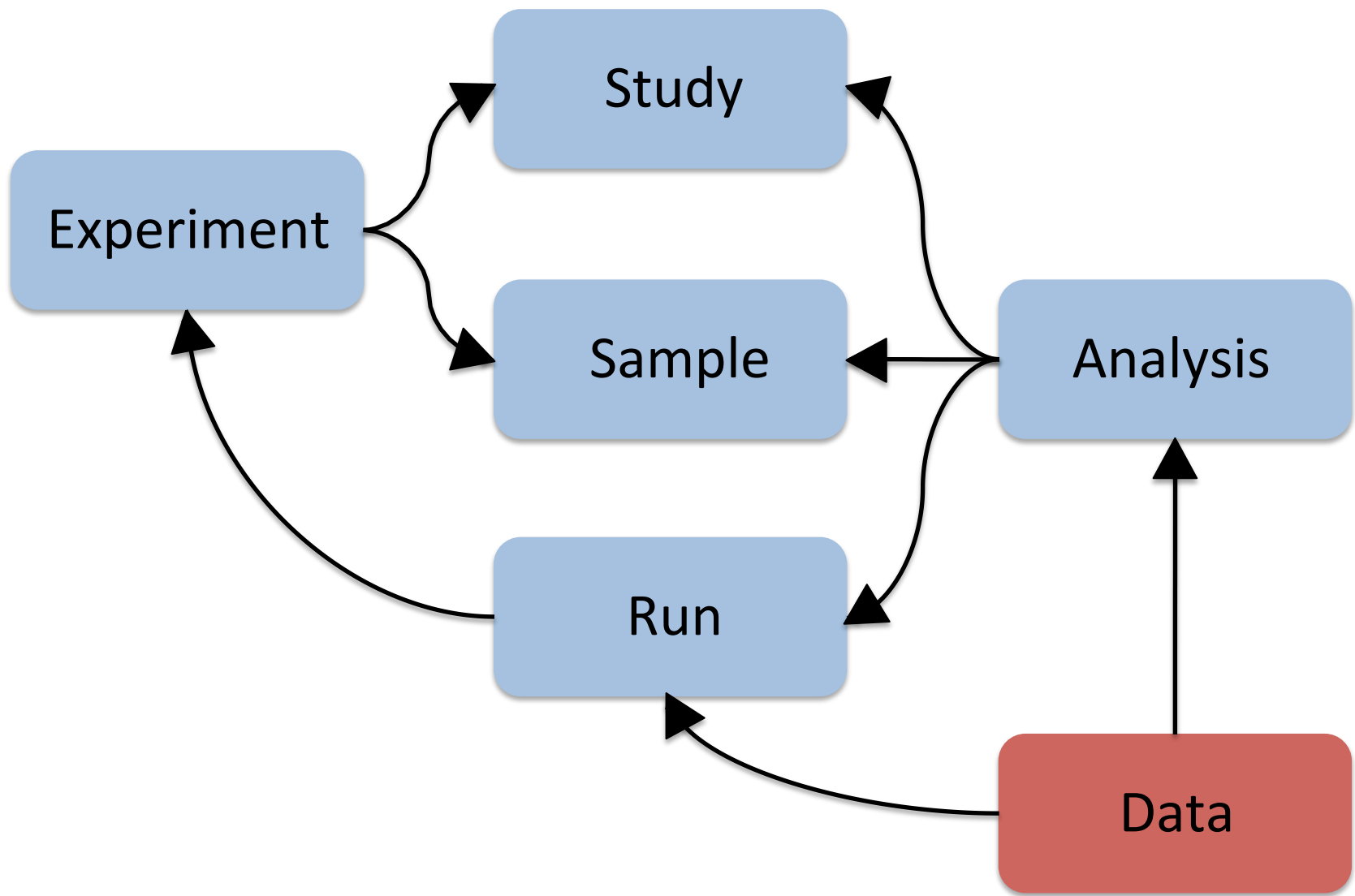
Column C suggests the descriptor’s potential importance for metagenomic data analysis (H – high relevance, M – medium relevance, L – low relevance).

<b>A</b> non-mandatory descriptors of sequencing provenance	<b>B</b> descriptor format	<b>C</b> value for analysis (H/M/L)
sequencing centre contact	text	M
sequencing experiment name	text	L
library name	text	L
library description	text	L
library construction protocol	text	M
library construction method (MIMS)	text	M
library size (MIMS)	numeric	M
library reads sequenced (MIMS)	numeric	M
library vector (MIMS)	text	M
library screening strategy (MIMS)	text	M
insert size (for paired read files)	numeric	M
spot layout (for SFF read files)	controlled vocabulary (single, paired FF, paired FR)	M
linker sequence (for SFF read files)	sequence of nucleotides	H
multiplex identifiers (MIMS)	sequence of nucleotides	H

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

adapters (MIMS)	sequence of nucleotides	H
quality scoring system (for Fastq files)	controlled vocabulary (phred, log-odds)	H
quality encoding (for Fastq files)	controlled vocabulary (ascii, decimal, hexadecimal)	H
ascii offset (for Fastq files)	controlled vocabulary (!, @)	H
nucleic acid extraction SOP (MIMS)	text	H
nucleic acid amplification SOP (MIMS)	text	H
sequencing coverage	numeric	H





**A**

## Metagenomic Analysis Metadata & Example

**Name:** EBI metagenomics

**Type (objective):** Functional and taxonomic analysis

**Centre:** EMBL-EBI

**Date:** 18/07/2016

**Total CPU time (hrs):** 102

**Max memory (GB):** 18

**Sequence Reference:** ERP00001

...



Component 3

Component 2

Component 1



B

Component Metadata & Example	
<b>Name:</b> InterProScan	
Input(s)	Output 1 Component: FragGeneScan
	<b>Reference DBs:</b> InterPro 58.0 signature lib.
Tools	<b>Name:</b> InterPro
	<b>Version:</b> 5.19-58.0
	<b>Source Code:</b> <a href="https://www.ebi.ac.uk/interpro/interproscan.html">https://www.ebi.ac.uk/interpro/interproscan.html</a>
	<b>Parameters</b>
Output(s)	<ol style="list-style-type: none"><li>1. Hit XML</li><li>2. GO slim</li></ol>

C

### Component Dependencies

