



Title

The metagenomic data life-cycle: standards and best practices

Authors

Petra ten Hoopen¹, Robert D. Finn¹, Lars Ailo Bongo², Erwan Corre³, Bruno Fosso⁴, Folker Meyer⁵, Alex Mitchell¹, Eric Pelletier^{6,7,8}, Graziano Pesole^{4,9}, Monica Santamaria⁴, Nils Peder Willassen², and Guy Cochrane^{1*}

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

²UiT The Arctic University of Norway, Tromsø N-9037, Norway

³CNRS-UPMC, FR 2424, Station Biologique, Roscoff 29680, France

⁴Institute of Biomembranes and Bioenergetics, CNR, Bari 70126, Italy

⁵Argonne National Laboratory, Argonne IL60439, USA

⁶Genoscope, CEA, Évry 91000, France

⁷CNRS / UMR-8030, Évry 91000, France

⁸Université Évry val d'Essonne, Évry 91000, France

⁹University of Bari "A. Moro", Bari 70126, Italy

1 * corresponding author

2 Guy Cochrane

3 European Molecular Biology Laboratory

4 European Bioinformatics Institute

5 Wellcome Genome Campus

6 Hinxton

7 Cambridge CB10 1SD

8 United Kingdom

9 Tel: +44(0)1223 494444

10 Email: cochrane@ebi.ac.uk

11 Email addresses:

12 AM: mitchell@ebi.ac.uk

13 BF: b.fosso@ibbe.cnr.it

14 EC: erwan.corre@sb-roscoff.fr

15 EP: eric.pelletier@genoscope.cns.fr

16 FM: folker@anl.gov

17 GC: cochrane@ebi.ac.uk

18 GP: g.pesole@ibbe.cnr.it

19 LAB: lars.ailo.bongo@uit.no

20 MS: m.santamaria@ibbe.cnr.it

21 NPW: nils-peder.willassen@uit.no

22 PTH: petratenhoopen@yahoo.co.uk

23 RDF: rdf@ebi.ac.uk

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, metatranscriptomics and metaproteomics (**Figure 1**), enabling comprehensive insights into community composition and function. 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.

Here, we review the workflow for metagenomics data generation and analysis. Where possible, we specify essential parameters in the workflow and advise on standardised

1 systematic reporting of these as variables. We build on the expertise of major public
2 genomic and metagenomic resources: the European Nucleotide Archive (ENA) [2] and
3 EMBL-EBI Metagenomics (EMG) [3] at the EMBL European Bioinformatics Institute in UK;
4
5 MG-RAST [4] at Argonne National Laboratory in USA; and the extensive knowledge bases in
6
7 metagenomics available at research centers of excellence, the UiT in Norway, Genoscope in
8
9 France, SB-Roscoff in France and CNR in Italy.
10
11
12
13
14
15
16
17

18 For the purposes of this paper, we will predominately use marine metagenomics as a ‘use
19
20 case’ to highlight the standards environment that we describe. However, we believe that
21
22 these examples will broadly translate to all areas of metagenomics research, regardless of
23
24 the environment under study. From the outset, we stress that we do not wish to promote a
25
26 specific workflow, but rather to demonstrate the importance of having systematic reporting
27
28 conventions that accurately describe any chosen workflow, from sampling through to the
29
30 presentation of analysis outputs. Our aim is to describe conventions and standards that are
31
32 inclusive and extensible, and able to cope with evolving scientific developments in the field.
33
34 Furthermore, where a given standard has not emerged, we will point to, or propose, a
35
36 generalised ‘best practice’ that can be used in its place. While this may produce a
37
38 foundation from which a new standard could be proposed, any additional formal scientific
39
40 standards need to come from the community and be ratified by scientific bodies, such as the
41
42 Genomics Standards Consortium (GSC) [5].
43
44
45
46
47
48
49
50
51
52
53

54 **[Figure 1]**
55
56
57
58
59
60
61
62
63
64
65

Overview of the metagenomics data model

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

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

[Figure 2]

Information associated with DNA sequence is frequently referred to as '**metadata**'. This includes all information described in the study, sample, experiment and run data objects.

Primary data represent, in this context, primary experimental sequence reads produced by sequencing machines, sometimes with some basic level of quality control (filtering out of poor-quality reads, clipping of low-quality regions, etc.). Following this, for some

1 metagenomics studies, the **primary data** are analysed directly (e.g. 16S or 18S rRNA gene
2 amplicon studies), while in others, they are assembled into contigs before undergoing
3 further analysis. Regardless of the approach, the output of any computational analysis
4 process (including assembly) on the primary data are here referred to as **derived data**. We
5 discuss derived data in more detail below, but the more harmonised the formats and
6 validations for data and metadata objects, the more easily the generated data be shared,
7 discovered, re-used and re-purposed.
8
9

10 Each metagenomics study has a scope, aim and one or more (human) contributors in each
11 step of the workflow, who may be distributed over a wide geographical area. It is essential
12 to capture contextual information regarding the contributors, since this supports
13 appropriate attribution and credit , and clarifies the responsible parties for each step of the
14 workflow. Contributors to (1) material sampling (2) primary data generation and (3) derived
15 data generation should always be clearly presented in data records. Minimum metadata
16 checklists frequently do not specifically capture data generating or contributing institutions.
17 However, this information is frequently available and can be parsed from the registration
18 systems for reporting individual steps of the data generation workflow or from associated
19 peer-reviewed publications.
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 Sampling

51 The method of collecting a sample (a fundamental unit of material isolated from the
52 surrounding environment) is dictated by the nature of the community under investigation,
53
54
55
56
57
58
59
60
61
62
63
64
65

1 the environment in which it is found and the type of 'omics' investigation being performed.
2
3 The slightest deviation in method, regardless of the protocol chosen, can have a profound
4
5 impact on the final 'omics' analysis results. It is therefore essential that the details of the
6
7 sampling process are captured accurately and in a standardised way.
8
9

10
11
12 Domain experts are in the best position to formulate opinions on the general scope and
13
14 content of contextual data (environmental characteristics observed or measured during
15
16 sample collection) and methodological variables (such as sampling volume and filtration
17
18 method). These opinions are conventionally formalised as data reporting standards by
19
20 community initiatives such as the GSC for genomics data [5], the Proteomics Standards
21
22 Initiative (PSI, [6]) for proteomics data or Geo Bon [7] for biodiversity data [8].
23
24
25
26
27
28
29

30
31 The Minimum Information about Metagenomic Sequence (MIMS, [9]) is a GSC-developed
32
33 data-reporting standard, designed for accurate reporting of contextual information for
34
35 samples associated with metagenomic sequencing, and is also largely applicable to
36
37 metatranscriptomics studies. Minimum Information about a MARKer gene Sequence
38
39 (MIMARKS, [10]) is another GSC-developed contextual data reporting standard for reporting
40
41 information about a metabarcoding study, which is referred to in the standard as the
42
43 'MIMARKS-survey investigation type'.
44
45
46
47
48
49

50
51 MIMS and MIMARKS are a part of a broader GSC standard, the Minimum Information of any
52
53 (x) Sequence (MIxS) [11], which describes 15 different environmental packages that can be
54
55 used to specify the environmental context of a sequenced microbial community, such as air,
56
57
58
59
60
61
62
63
64
65

1 water or host organism-associated. The MIxS descriptors can be combined with any
2 environmental package and together provide rich information on sampling context.
3

4
5
6
7 To illustrate, **Table 1** summarises the minimum set of elements required for description of a
8 metagenomic sample taken from an aquatic environment. It uses MIMS mandatory
9 descriptors, combined with the mandatory descriptors of the Water Environment package.
10 Similarly, a sample taken from a gut of a fish host can be described using the MIMS core
11 descriptors, in combination with descriptors in the Host-associated Environment package.
12
13 The 15 different environmental packages defined by the GSC are available from the GSC
14 website as a single bundled download [12] and are presented in a host of informatics tools
15 that support data reporting and presentation, such as the submission tools of the databases
16 of the International Nucleotide Sequence Database Collaboration [13] and ISAtools [14]. It
17 remains up to the experimentalist to choose the most appropriate package from within the
18 checklist bundle for their study, thereby defining the list of fields that will be used to
19 capture relevant metadata. Before embarking on a metagenomics study, we recommend
20 that the appropriate checklist be identified, so that the appropriate metadata can be
21 captured during the experiment, rather than retrospectively having to determine these
22 metadata.
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 **[Table 1]**

50
51
52
53 The M2B3 data reporting and service standard [15] specifically addresses contextual data
54 relating to marine microbial samples. It represents a common denominator of contextual
55 data from data standards used in the public genomic data archives (MIxS, Version 4.0, [12]),
56
57
58
59
60
61
62
63
64
65

1 pan-European network of oceanographic data archives (CDI schema, Version 3.0, [16]) and
2 pan-European network of biodiversity data resources (OBIS schema, Version 1.1, [17]). This
3 M2B3 unified data standard significantly simplifies contextual data reporting, since it
4 provides an interoperable solution for sharing contextual data across data archives from
5 different scientific domains. A minimum M2B3 checklist for reporting contextual data
6 associated with marine microbial samples is summarised in **Table 2**.

7
8
9
10
11
12
13
14
15
16
17
18 **[Table 2]**

19
20
21
22 For most adopted standards of this type, only a few contextual data are mandatory,
23 reflecting the balance between usability for the experimentalist reporting his/her science
24 and consumers re-using this science; limiting the number of mandatory fields lowers the
25 burden for experimentalists to comply with the standard, while a small number of
26 parameters are universally, or near-universally, required for downstream analysis. The
27 importance of the optional MIxS and M2B3 fields for metagenomic data analysis is detailed
28 in **Table 3**.

29
30
31
32 For omics studies it can be useful to introduce a sample attribute, ‘biological replicate’,
33 where the attribute value is the accession number of the related biological sample. In
34 contrast, ‘technical replicates’, for which only a single sample exists, are treated
35 downstream in the workflow.

36
37
38
39 Consistent and rich contextual data can become a powerful tool for metagenomics data
40 analysis. Two marine studies, the *TARA* Oceans sequencing study (PRJEB402, [18]) and

1 Ocean Sampling Day (OSD, PRJEB5129, [19]) both use the same M2B3 contextual data
2 reporting standard enabling comparison of data within and across studies. For instance,
3 data from the *TARA* Oceans shotgun sequencing of the prokaryotic fraction filtered from
4 seawater (PRJEB1787, [20]) can be compared to the shotgun data from OSD (PRJEB8682,
5 [21]), enabling detailed or complex queries. Specifically, a taxonomic or functional profile
6 from the *TARA* Oceans sample from depth 5m and salinity 38psu (SAMEA2591084, [22]) can
7 be compared to profiles of the OSD sample from the depth 5m (SAMEA3275502, [23]) or the
8 OSD sample with the same salinity 38psu (SAMEA3275531, [24]). In contrast, very few
9 conclusions can be drawn from a comparison to a sample with insufficient contextual
10 information (SAMN00194025, [25]).
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27

28 **[Table 3]**
29
30
31
32

33 Details of the project investigators are usually recorded in the Study metadata object and
34 sampling contextual data are mostly captured in the Sample metadata object, **Figure 2**. A
35 common way to standardise reporting of contextual data is via a checklist of key-value pairs,
36 thereby ensuring parameters of a similar kind are described consistently. Furthermore,
37 syntactic and semantic rules can be pre-defined in the checklist, enabling validation of
38 compliance with these rules. For instance, automated checks can be applied to test whether
39 a mandatory descriptor (key) in the checklist has a value and whether the value is in a
40 specified format. Each element to be checked can be pre-defined as text, a class or term
41 from an ontology, a controlled vocabulary or taxonomic index, or formulated as a regular
42 expression. (Regular expressions can be used, for example, to check that the key 'collection
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 date and time' complies with ISO 8601 standards and that numeric values lie within a
2 defined range.)
3
4
5
6

7 The most common formats for sharing Study and Sample metadata are XML, TSV, ISA-tab or
8 JSON formats. Examples of the Study and Sample XML are available from the European
9 Nucleotide Archive [26], [27], where the files are also validated against the XML schema
10 [28]. Regardless of the format used to supply the metadata, because they all use the same
11 underlying standards, a simple translation between the formats enables different data to be
12 compared. This allows scientists to use different tools or approaches that they are most
13 familiar with, whilst ensuring consistent delivery of the metadata.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

30 Sequencing

31
32
33
34 Once a sample is collected and its provenance recorded, it is subjected to preparation steps
35 for nucleotide sequence analysis. This may happen immediately after sampling, or in stages
36 over many months. Processing steps cover all handling of the sample leading to the DNA
37 isolation. Although MIxS covers some of the metadata fields for reporting the DNA
38 extraction steps, it is extremely difficult to define a generic set of fields describing the DNA
39 extraction method with a high granularity due to its complexity and diversity. For example,
40 it might be relatively straightforward to identify variables for reporting isolation of DNA
41 from a seawater sample but that will not suit the more complex DNA isolation procedure for
42 a sediment sample. We suggest the best practice here is to use the existing MIxS fields, such
43 as the *sample material processing*, *nucleic acid extraction* and *nucleic acid amplification* for
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 concise description of the nucleic acid preparation. A detailed description, or a reference to
2 the material preparation steps, is important due to the significant influence this can have on
3 the observed profile of the microbial community under investigation.
4
5
6
7
8
9

10 Equally critical for the downstream metagenomic data analysis and interpretation is the
11 reporting of sequencing library preparation protocols and parameters as well as sequencing
12 machine configurations.
13
14
15
16
17
18
19

20 **Table 4** shows mandatory descriptors for new generation nucleotide sequencing
21 experiments as currently captured by INSDC databases. **Table 5** lists non-mandatory
22 descriptors including MIMS sequence-related descriptors and provides our opinion on the
23 importance of these descriptors for metagenomic data analysis. Note that while a number
24 of controlled vocabularies have been developed for accurate recording of sequencing
25 experiment parameters, the evolution of these constrained vocabularies is very dynamic
26 and driven by technological advances.
27
28
29
30
31
32
33
34
35
36
37
38
39

40
41 **[Table 4]**
42
43
44

45
46 **[Table 5]**
47
48
49

50 Variable parameters of the library preparation and instrumentation are captured in the
51 metadata objects Experiment and Run (see **Figure 2**). Examples of the Experiment and Run
52 XML are available, for example, from the ENA [29],[30]. Each Experiment should refer to
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Study and Sample objects, to provide context for the sequencing, and is referred to from the
2
3 Run objects, which point to the primary sequencing reads.
4
5
6

7 The **primary data** (the reads) are stored in files of various formats, which can be standard
8
9 (BAM, CRAM or Fastq) or a platform-specific, as with SFF, PacBio, Oxford Nanopore or
10
11 Complete Genomics. Information on the read data format must be indicated in the
12
13 description of sequencing.
14
15
16

17
18
19
20 The minimum information encapsulated in read data files are base calls with quality scores.
21
22
23 Quality requirements on read data files are file format-specific and are summarised, for
24
25 example, in the ENA data submission documentation [31]. A freely available diagnostic tool
26
27 for validation of CRAM and BAM files is the Picard ValidateSamFile [32]. Validation of Fastq
28
29 files is less straightforward since there is no single FASTQ specification. Recommended
30
31 usage of FASTQ can be found for instance in the ENA guidelines [33]. An open resource for
32
33 managing next generation sequencing datasets is the NGSUtils [34], which also contains
34
35 tools for operations with FASTQ files. As sequencing technologies change over time the
36
37 formats and associated validation tools may well change, so a comprehensive list of formats
38
39 and tools is likely to become outdated. The key point is to adopt a widely used format and
40
41 to check for file format and integrity (e.g. checksums).
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Analysis

Standards in metagenomics for the description of sampling and sequencing have grown out of those from more traditional genomics. While there are still some shortcomings in these standards, as highlighted in the previous sections, metadata concerning sampling and sequencing are commonly captured for metagenomics studies. Compliance is high partly due to the scientific journals requiring scientist to submit sequence data to an INSDC database prior to publication. However, there are currently no standards for reporting how metagenomics datasets have been analysed. While systematic analysis workflows, such as those offered by IMG, IMG/M [35] META-pipe [36] and MG-RAST, provide a standard that is documented (albeit in different ways), many published datasets are analysed by in-house bespoke pipelines. Although many authors provide an outline in the 'materials and methods' or 'supplementary materials' section of their publications, it is rarely possible to reproduce the analysis from this alone, due to missing software parameters, lack of detail on software versions and ambiguous reference databases and their associated versions.

Typically, once the sequence read files have been produced, they are analysed using one or more workflows [37], with each workflow comprising different data processing or analysis components. Most workflows involve aspects such as quality control (for example, removing sequences that fail to meet predefined quality scores), assembly, sequence binning (e.g. identifying 16S rRNA genes or protein coding sequences), and taxonomic classification of sequences and/or functional prediction. However, each workflow will be tailored to how the sample has been processed and the question being addressed. For example, if a sample has been size-fractionated for viruses, using a 0.22 μm filter, there would be little point analysing

1 the data for eukaryotic 18S rRNA, as any eukaryotic organisms would have been physically
2 removed from the sample before the DNA extraction process.
3
4

5
6
7 Analysis workflows typically have one or more of the following components:
8

9
10 (i) Central algorithmic software, which may be from a third party source.
11

12 (ii) 'Glue' software that may ensure input/output formats, or split/join input files for
13 parallelisation.
14
15

16 (iii) Reference datasets that are used by (i). For example, the Greengenes database of 16S
17 rRNA genes [38], SILVA database of 16S/ 18S SSU ribosomal RNA genes [39] and the NCBI
18 non-redundant database of non-identical protein sequences [40], for taxonomic or
19 functional analysis.
20
21
22
23
24
25
26

27
28 However, even knowing these elements may not be sufficient for analyses to be
29 independently re-created. For example, the algorithm may accept a set of input parameters
30 that can be used to fine-tune an analysis, such as selecting an E-value threshold for
31 determining significance of a sequence match to a reference database. Other parameters
32 may influence speed-performance, which allows the original analysis to complete in a timely
33 fashion, but they may or may not have an effect on the results. For example, running
34 *hmmsearch* from the HMMER package, changing the number of CPUs used will not change
35 the results, but changing options on the heuristics such as the --F1 threshold (which controls
36 the number of sequences passing the first heuristic state) may alter the output; both will
37 potentially increase performance in terms of speed.
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 Capturing and reporting *all* provenance information is essential to understand exactly what
2 analysis has been performed on the data, and to ensure reproducibility [41]. The use of
3 publicly available analysis pipelines (such as EMG, MG-RAST or META-pipe) helps with this
4 process, since analysis is performed using pre-defined components, settings and databases
5 (or, in some cases, using user-selected components, selected from a predefined list of
6 options). Nevertheless, capturing analysis metadata remains essential as, for example, MG-
7 RAST allows the users to dynamically set E-value thresholds after the pipeline analysis has
8 been performed. Furthermore, the tools, libraries, and reference databases used by the
9 pipelines are regularly updated, and thus capturing analysis provenance information is
10 vitally important and should be systematically ‘tagged’ to the results.
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27

28 To date, there is no universally endorsed ‘Analysis standard’ for describing and recreating a
29 metagenomics analysis pipeline, and without this standard (and subsequent
30 adoption/enforcement) it will continue to be difficult, if not impossible, to reproduce
31 analysis workflows. However, all is not lost. ‘Workflows’ and their definitions is an active
32 field of computer science research, and potential solutions are already available, including
33 Common Workflow Language (CWL), Yet Another Workflow Language (YAWL), Business
34 Process Execution Language (BPEL) and Microsoft Azure’s Workflow Definition Language to
35 name but a few. Several of the co-authors for this publication already participate in the GSC
36 M5 consultation group, which aims to define a standard enabling the recreation and
37 exchange of metagenomics data sets. In the absence of a standard, we believe it is
38 important to define some of the basic best practices, from which an accepted standard
39 would formally encapsulate. For simplicity, we will focus on a single ‘best practice’ use case:
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 the description of the analysis of a run. Other types of analysis, such as pooling of runs or
2
3 comparing results between runs are beyond the scope of this article.
4
5
6

7 A schematic overview of a best practice for analysis metadata collection is shown in **Figure**
8
9 **3A**. An overarching set of metadata relating to analysis should encapsulate generic
10 information such as analysis centre, name of bioinformaticians, analysis objectives, name of
11 overall analysis (if appropriate) and the date on which the analysis was performed. It should
12 also contain appropriate pointers to the run data, run sequence metadata and associated
13 sample data. Underneath the overarching analysis metadata is a collection of analysis
14 components, which describe each stage of the analysis **Figure 3B**. Each component can be
15 divided into three sections: input(s), analysis algorithm and output(s).
16
17
18
19
20
21
22
23
24
25
26
27
28
29

30 The input section should describe the details of the various inputs to the analysis, which
31 could be the raw sequence reads or the output of another analysis component, reference
32 databases and their provenance data, such as version, where necessary. The analysis section
33 should contain the algorithm tool, version, all parameters used and a basic description of
34 the analysis. The output section should describe each output from the analysis, together
35 with a description of contents and format.
36
37
38
39
40
41
42
43
44
45
46
47

48 Each analysis component could then be coupled to form an analysis workflow as shown in
49 **Figure 3C**. The workflow may be in a portable intermediate format that can be submitted to
50 a workflow manager for execution in a specific environment.
51
52
53
54
55
56
57

58 **[Figure 3]**
59
60
61
62
63
64
65

1
2 This best practice framework is merely that - a best practice, and we have not touched on
3
4 the technical issues of how to capture this information or on controlled vocabularies (since
5
6 these need to come from the community). Furthermore, enforcing compliance and
7
8 validation against the standard will also require a community effort. Complete validation
9
10 would require the standard to be machine readable and deployable, with potentially the
11
12 need to have small 'test' datasets and their associated results, to perform regression testing
13
14 of the analysis metadata. However, who is responsible for validation and what happens if
15
16 something fails after publication are open questions. This could arguably be a step too far;
17
18 currently sampling and sequence metadata are validated against the standard, but taken in
19
20 good faith to be correct beyond this.
21
22
23
24
25
26
27
28
29
30
31
32
33
34

35 Analysis Results Archiving - a final piece? 36 37 38

39 Having an analysis provenance standard would allow metagenomics analysis results to be
40
41 recreated more readily. While this is undoubtedly an important and necessary step, it has a
42
43 major limitation within the community. As indicated [42], the fraction of money spent on
44
45 informatics from an overall project budget is increasing dramatically. Metagenomics
46
47 datasets tend to be large, in the order of GB-TB and processing may take 1000s of CPU
48
49 hours, restricting reanalysis to only those with significant compute resources. For example,
50
51 the subset of the *TARA* Oceans Ocean Microbiome Project (PRJEB7988, [43]) that has been
52
53 size-fractioned for prokaryotes comprises 135 samples with 248 runs containing 28.8 billion
54
55 reads. The analysis output represents about 10TB of data with 23.2 billion predicted protein
56
57
58
59
60
61
62
63
64
65

1 coding sequences. Thus, reanalysis would be costly and potentially wasteful if a particular
2 workflow had already been run on the data. Therefore, a final step in a metagenomics
3 analysis is the appropriate archiving of results. There is an obvious cost-benefit balance to
4 be drawn here, as storing every intermediate of a workflow would lead to an explosion of
5 data. Clearly, key intermediates and outputs of an analysis workflow need to be determined.
6
7 These key archived components will be tailored to the analysis, but should at least include
8 operational taxonomic unit (OTU) counts and assignments, functional assignment counts
9 and read/sequence positional information for the aforementioned assignments. Such data
10 files are already made available from MG-RAST and EMG and those from other sources are
11 accepted for archiving within ENA.
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27

28 If metagenomic assemblies have been performed, then these should have an appropriate
29 structure of contigs, scaffolds or chromosomes with an appropriate format as detailed for
30 example in the ENA data submission documentation [44]. Due to the overheads of
31 producing an assembly, these should be archived, ideally with an INSDC database.
32
33
34
35
36
37
38
39
40

41 The data model for metagenomics, as described in **Figure 2**, represents metagenomic
42 analysis results in the data Analysis object with appropriate pointers to the corresponding
43 run sequence metadata and associated sample collection contextual data. While there is an
44 established practice to archive primary sequence data in the Run object and assemblies of
45 the primary sequences in the Analysis object, it is not a common practice to archive results
46 of functional and taxonomic metagenomic analysis of in-house bespoke pipelines. It would
47 be beneficial to the metagenomics community to include this into the best practice and such
48 data are accepted by ENA for archiving. The metagenomics standard environment reviewed
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 here as well as outcomes of the GSC M5 consultation group can contribute to defining
2 required descriptors of the Analysis object for archiving of metagenomics analysis results,
3
4 which can serve as a framework for exchange of metagenomics data sets on a routine bases,
5
6
7 similarly as is currently done for the primary sequence data.
8
9

10 11 12 13 14 15 Future

16
17
18
19 One challenge over the next several years will be the validation of compliance across the
20
21 entirety of the standards and best practise that we have covered. While validation tools and
22
23 recommended practises exist for parts (e.g. contextual data descriptors using MIxS-
24
25 compliant validation tools from ISA and experimental descriptors upon submission to an
26
27 INSDC database), not all parts have such maturity (e.g. analysis descriptors) and there exists
28
29 no overarching validation protocol for an entire metagenomics study. The GSC is aiming to
30
31 contribute in this area with the introduction of MIxS “profiles”, to provide an overlay on top
32
33 of MIxS environmental packages and the core MIxS fields. These profiles will enable the
34
35 creation of tool suites for compliance checking. In addition (and perhaps more importantly)
36
37 they will enable groups of researchers, institutes, funders and other communities to define
38
39 levels of compliance for contextual data sets. Examples of this are the NSF NEON [45] and
40
41 the NSF CZO [46] networks that are working with the GSC to establish silver, gold and
42
43 platinum sets of parameters that need to be provided and validated for data sets to be
44
45 compliant. A key moment in the acceptance of said new profiles will be the availability of
46
47 tools support for data creators, end-users and portals. Imagine, for instance, a search for
48
49 data sets in EMG or MG-RAST that allows restriction of the search to just platinum-level data
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 sets. For the data consumer this will result in better ways of telling their science story using
2 third party data and for the data creator this will provide guidance on what to create for a
3 specific community. In addition funding agencies can require certain minimal compliance.
4
5
6
7
8
9

10 A further area to be addressed is that of standards around descriptions of metagenomics
11 analyses. The creation of a lightweight data standard for an analysis object that allows easy
12 transfer of analyses is a key goal of the GSC M5 initiative but the complexity of the task and
13 lack of dedicated resourcing has rendered progress slow; while frameworks and systems for
14 recording analysis provenance need to be established, we have aimed to indicate in this
15 publication a set of best practice that can form the foundation for a community standard
16 enabling the recreation and exchange of metagenomics data sets. Improving
17 standardisation will also help raise clarity in the literature around metagenomics through a
18 tightening of language. For example, UProC [47] uses Pfam [48] matches as a reference
19 library with the results being referred to as a 'Pfam hit'. However, this may not necessarily
20 be a Pfam hit, as a Pfam hit is defined as a sequence match scoring greater than Pfam
21 defined threshold to the Pfam profile Hidden Markov Model (i.e. the Pfam database
22 method).
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 Conclusions

49 In this overview of the metagenomics standard environment we have outlined best practice
50 for the reporting of metagenomics workflows. We have reviewed the essential steps: (1)
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 material sampling, (2) material sequencing (3) data analysis and (4) data archiving, and
2 highlighted essential variable parameters and common data formats in each step.
3
4
5

6
7 Reporting on the provenance of a sample and associated nucleotide sequence data is largely
8 established by public sequence data repositories and is also being addressed by contextual
9 data standardisation initiatives. In contrast, a reporting standard on metagenomics data
10 analysis is absent, yet the high complexity of metagenomics creates a pressing demand for
11 establishing such a practice. Capturing key metadata relating to analysis would greatly
12 improve reproducibility. Archiving key results of the metagenomics data analysis would
13 allow a more accurate evaluation of the benefits of reproducing the analysis.
14
15

16 Only by adopting these standards and best practices can metagenomics data be assessed
17 against the FAIR (Findable, Accessible, Interoperable and Reusable) principles that should be
18 applied to any scientific dataset [49].
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

40 List of abbreviations 41 42 43

44 BAM: Binary Alignment/Map

45
46 BPEL: Business Process Execution Language

47
48
49 CDI: Common Data Index

50
51 CHEBI: Chemical Entities of Biological Interest

52
53
54 CNR: Consiglio Nazionale delle Ricerche

55
56
57 CPU: Central Processing Unit
58
59
60
61
62
63
64
65

1 CRAM: Compression Reduced Alignment/Map
2
3 CWL: Common Workflow Language
4
5 EBI: European Bioinformatics Institute
6
7 EMBL: European Molecular Biology Laboratory
8
9
10 EMG: EBI Metagenomics
11
12 ENA : European Nucleotide Archive
13
14
15 ENVO: Environment Ontology
16
17
18 GEO BON: Group on Earth Observations Biodiversity Observation Network
19
20
21 GO: Gene Ontology
22
23 GSC: Genomic Standards Consortium
24
25
26 GSC M5: Genomic Standards Consortium M5 working group
27
28
29 IMG/M: Integrated Microbial Genomes with Microbiomes
30
31 INSDC: International Nucleotide Sequence Database Collaboration
32
33
34 ISO: International Standards Organisation
35
36
37 JSON: JavaScript Object Notation
38
39
40
41 MG-RAST: Metagenomic Rapid Annotations using Subsystems Technology
42
43
44 MIMS: Minimum Information about a Metagenome Sequence
45
46
47 MIMARKS: Minimum Information about a MARKer gene Sequence
48
49
50
51
52 M2B3: Marine Microbial Biodiversity Bioinformatics and Biotechnology
53
54
55
56
57
58
59
60
61
62
63
64
65

1 OSD: Ocean Sampling Day
2
3 OTU: Operational Taxonomic Unit
4
5 PSI: Proteomics Standards Initiative
6
7 SDN: SeaDataNet
8
9
10 SFF: Standard Flowgram Format
11
12
13 SSU: Small Subunit ribosomal RNA
14
15
16 TB: Terabyte
17
18 TSV: Tab Separated Values
19
20
21 UiT: Universitetet i Tromsø
22
23 UProC: Ultrafast Protein Classification
24
25
26 UTC: Coordinated Universal Time
27
28 XML: Extensible Markup Language
29
30
31 YAWL: Yet Another Workflow Language
32
33
34
35
36
37

38 Declarations

39
40
41
42
43

44 Ethics approval and consent to participate

45
46
47
48 Not applicable
49
50
51

52 Consent for publication

53
54
55
56 Not applicable
57
58
59
60
61
62
63
64
65

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.

1
2
3 [3] Mitchell A, Bucchini F, Cochrane G, Denise H, Hoopen PT, Fraser M, et al. EBI
4 metagenomics in 2016 - an expanding and evolving resource for the analysis and archiving
5 of metagenomic data. Nucleic Acid Res. 2016;44:D595-D603.
6

7
8
9
10 [4] Meyer F, Paarmann D, D'Souza M, Olson R , Glass EM, Kubal M, et al. The
11 Metagenomics RAST server – a public resource for the automatic phylogenetic and
12 functional analysis of metagenomes. BMC Bioinformatics. 2008;9:386.
13
14
15

16
17
18
19
20 [5] Field D, Amaral-Zettler L, Cochrane G, Cole JR, Dawyndt P, Garrity GM et al. The
21 Genomic Standards Consortium. PLoS Biol. 2011;9(6):e1001088.
22
23
24

25
26
27
28 [6] Orchard S., Hermjakob H. and Apweiler R. (2003) The Proteomics Standards
29 Initiative. Proteomics. 2003;3(7):1374-6.
30
31
32

33
34
35
36 [7] The Group on Earth Observations Biodiversity Observation Network.
37
38 <http://geobon.org/essential-biodiversity-variables/connect-with-geoss/>. Accessed 19 Jan
39
40
41 2017.
42
43
44

45
46 [8] Pereira HM, Ferrier S, Walters M, Geller GN, Jongman RHG, Scholes RJ et al. Essential
47 Biodiversity Variables. Science. 2013;339(6117):277-8.
48
49
50

51
52
53
54 [9] The Minimum Information about a Metagenome Sequence.
55
56 <http://wiki.genesc.org/index.php?title=MIGS/MIMS>. Accessed 19 Jan 2017.
57
58
59
60
61
62
63
64
65

1 [10] The Minimum Information about a Marker Gene Sequence.

2 <http://wiki.gensc.org/index.php?title=MIMARKS>. Accessed 19 Jan 2017.

3
4
5
6
7 [11] Yilmaz, P, Gilbert, JA, Knight, R, Amaral-Zettler L, Karsch-Mizrachi I, Cochrane G et al.

8 The Genomic Standards Consortium: bringing standards to life for microbial ecology. The
9
10 ISME J. 2011;**5**,1565-7.

11
12
13
14
15
16
17 [12] The Minimum Information about any (x) Sequence, Version 4.0. [http:// wiki. gensc.](http://wiki.gensc.org/index.php?title=MlxS)
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

[org/ index. php? title= MlxS](http://wiki.gensc.org/index.php?title=MlxS). Accessed 19 Jan 2017.

[13] Cochrane G, Karsch-Mizrachi I, Takagi T and International Nucleotide Sequence Database Collaboration. The International Nucleotide Sequence Database Collaboration. Nucleic Acid Res. 2016;**44**:D40-D50.

[14] The ISA framework and tools. <http://isa-tools.org/>. Accessed 7 February 2017.

[15] Ten Hoopen P, Pesant S, Kottman R, Kopf A, Bicak M, Claus S et al. Marine microbial biodiversity, bioinformatics and biotechnology (M2B3) data reporting and service standards. The Standards in Genomic Sciences 2015;**10**:20.

[16] The Common Data Index, Version 3.0. [http:// www. seadatanet. org/ Data-Access/ Common-Data-Index-CDI](http://www.seadatanet.org/Data-Access/Common-Data-Index-CDI). Accessed 19 Jan 2017.

1 [17] The Ocean Biogeographic Information System data standard, Version 1.1.

2 <http://old.iobis.org/node/304>. Accessed 19 Jan 2017.

3
4
5
6
7 [18] The *TARA* Oceans umbrella project record of barcoding and shotgun sequencing.

8 <http://www.ebi.ac.uk/ena/data/view/PRJEB402>. Accessed 7 Feb 2017.

9
10
11
12
13
14
15 [19] The Ocean Sampling Day umbrella project record of amplicon and metagenome
16 sequencing. <http://www.ebi.ac.uk/ena/data/view/PRJEB5129>. Accessed 7 Feb 2017.

17
18
19
20
21
22
23 [20] The record of the *TARA* Oceans shotgun sequencing project of the prokaryotic
24 fraction filtered from seawater. <http://www.ebi.ac.uk/ena/data/view/PRJEB1787>. Accessed
25
26
27
28 7 Feb 2017.

29
30
31
32
33 [21] The record of the Ocean Sampling Day shotgun sequencing project from the year
34 2014. <http://www.ebi.ac.uk/ena/data/view/PRJEB8682>. Accessed 7 Feb 2017.

35
36
37
38
39
40
41 [22] The record of a *TARA* Oceans sample from depth 5m and salinity 38psu.
42
43
44 <https://www.ebi.ac.uk/metagenomics/projects/ERP001736/samples/ERS477979>. Accessed
45
46
47 7 Feb 2017.

48
49
50
51 [23] The record of an Ocean Sampling Day sample from depth 5m.
52
53
54 <https://www.ebi.ac.uk/metagenomics/projects/ERP009703/samples/ERS667511>. Accessed
55
56
57 7 Feb 2017.

1 [24] The record of an Ocean Sampling Day sample with salinity 38psu.
2 <https://www.ebi.ac.uk/metagenomics/projects/ERP009703/samples/ERS667548>. Accessed
3
4
5 7 Feb 2017.
6

7
8
9
10 [25] The record of an oil spill water sample from Gulfport.
11
12 <http://www.ebi.ac.uk/ena/data/view/SAMN00194025>. Accessed 7 Feb 2017.
13
14
15

16
17
18 [26] An example of a study XML. <http://www.ebi.ac.uk/ena/submit/preparing-xmls#study>. Accessed 19 Jan 2017.
19
20
21
22

23
24
25 [27] An example of a sample XML. <http://www.ebi.ac.uk/ena/submit/preparing-xmls#sample>. Accessed 19 Jan 2017.
26
27
28
29

30
31
32 [28] The validating XMLs document. <http://www.ebi.ac.uk/ena/submit/validating-xmls>.
33
34
35
36 Accessed 19 Jan 2017.
37

38
39
40 [29] An example of an experiment XML. <http://www.ebi.ac.uk/ena/submit/preparing-xmls#experiment>. Accessed 19 Jan 2017.
41
42
43
44
45

46
47
48 [30] An example of a run XML. <http://www.ebi.ac.uk/ena/submit/preparing-xmls#run>.
49
50
51
52 Accessed 19 Jan 2017.
53

54
55 [31] The document on the ENA-supported read file formats.
56
57
58
59 <http://www.ebi.ac.uk/ena/submit/read-file-formats>. Accessed 19 Jan 2017.
60
61
62

1
2 [32] The document on the Picard set of command line tools.
3

4
5 <http://broadinstitute.github.io/picard/>. Accessed 19 Jan 2017.
6
7

8
9
10 [33] The document on recommended usage of FASTQ files.
11

12 <http://www.ebi.ac.uk/ena/submit/read-data-format>. Accessed 7 February 2017.
13
14
15

16
17
18 [34] The document on the NGSUtils tools for next-generation sequencing analysis.
19

20 <http://ngsutils.org/>. Accessed 19 Jan 2017.
21
22
23

24
25 [35] Markowitz VM, Chen IM, Chu K, Szeto E, Palaniappan K, Pillay M, Ratner A et al.
26

27
28 IMG/M 4 version of the integrated metagenome comparative analysis system. *Nucleic Acid*
29
30 *Res.* 2015;42:D568-73.
31
32
33

34
35
36 [36] Robertsen EM, Kahlke T, Raknes IA, Pedersen E, Semb EK, Ernsten M et al. META-
37

38 pipe – pipeline annotation, analysis and visualisation of marine metagenomic sequence
39

40
41 data. <https://arxiv.org/abs/1604.04103v1>. Accessed 19 Jan 2017.
42
43
44

45
46 [37] Leipzig J. A review of bioinformatic pipeline frameworks. *Briefings in Bioinformatics*.
47

48
49 2016; doi:[10.1093/bib/bbw020](https://doi.org/10.1093/bib/bbw020).
50
51
52

53
54 [38] The Greengenes Database. <http://greengenes.secondgenome.com/>. Accessed 19 Jan
55

56
57 2017.
58
59
60
61
62
63
64
65

- 1
2
3
4
5 [39] The Silva Database. <https://www.arb-silva.de/>. Accessed 19 Jan 2017.
6
7
8
9
10 [40] The NCBI nr Database. <https://www.ncbi.nlm.nih.gov/refseq/>. Accessed 19 Jan 2017.
11
12
13
14
15 [41] Huntemann M, Ivanova NN, Mavromatis K, Tripp HJ, Paez-Espino D, Tennessen K et
16 al. The standard operating procedure of the DOE-JGI Metagenome Annotation Pipeline
17 (MAPv v4.). Standards in Genomic Sciences. 2016; doi:10.1186/s40793-016-0138-x.
18
19
20 [42] Sboner A, Mu XJ, Greenbaum D, Auerbach RK and Gerstein MB. The real costs of
21 sequencing: higher than you think! Genome Biology. 2011;12:125.
22
23
24
25
26
27
28 [43] The record of the TARA Oceans Ocean Microbiome Project.
29
30 <http://www.ebi.ac.uk/ena/data/view/PRJEB7988>. Accessed 7 Feb 2017.
31
32
33
34
35
36 [44] The document on a genome assembly submission to the ENA.
37
38 <http://www.ebi.ac.uk/ena/submit/genomes-sequence-submission>. Accessed 19 Jan 2017.
39
40
41
42
43
44 [45] The National Science Foundation National Ecological Observatory Network.
45
46 <http://www.neonscience.org/>. Accessed 19 Jan 2017.
47
48
49
50
51 [46] The National Science Foundation Critical Zone Observatory.
52
53 <http://criticalzone.org/national/>. Accessed 19 Jan 2017.
54
55
56
57
58
59
60
61
62
63
64
65

1 [47] Meinicke P. UProC: tools for ultra-fast protein domain classification. Bioinformatics.
2 2014;31(9):1382-8.
3

4
5
6
7 [48] Finn RD, Coggill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL et al. The Pfam
8 protein families database: towards a more sustainable future. Nucleic Acid Res.
9 2016;44:D279-85.
10

11
12
13 [49] Wilkinson MD, Dumontier M, Aalbersberg IJJ, Appleton G, Axton M, Baak A et al. The
14 FAIR Guiding Principles for scientific data management and stewardship. Sci. Data 2016;
15 doi:10.1038/sdata.2016.18 (2016).
16
17

18 [50] The controlled vocabulary for the INSDC country qualifier.
19 <http://www.insdc.org/country.html>. Accessed 19 Jan 2017.
20
21

22 [51] The Environment Ontology browser. [http://www.environmentontology.org/Browse-
23 EnvO](http://www.environmentontology.org/Browse-EnvO). Accessed 19 Jan 2017.
24
25

26 [52] The SeaDataNet L06 controlled vocabulary of platform categories.
27 http://seadatanet.maris2.nl/v_bodc_vocab_v2/search.asp?lib=L06. Accessed 19 Jan 2017.
28
29

30 [53] The SeaDataNet P02 controlled vocabulary of parameters.
31 http://seadatanet.maris2.nl/v_bodc_vocab_v2/search.asp?lib=P02. Accessed 19 Jan 2017.
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 [54] The controlled vocabulary of BODC data storage units.

2 http://seadatanet.maris2.nl/v_bodc_vocab_v2/search.asp?lib=P06. Accessed 19 Jan 2017.

3
4
5
6
7 [55] The CHEBI ontological classification of small chemical compounds.

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

9
10
11
12
13
14
15 [56] The National Center for Biotechnology Information taxonomy index.

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

17
18
19
20
21
22 [57] The controlled vocabulary for sequencing instrument models.

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

24
25
26
27
28
29
30 [58] The controlled vocabulary for the library source.

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

32
33
34
35
36
37
38 [59] The controlled vocabulary for the library strategy.

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

40
41
42
43
44
45
46 [60] The controlled vocabulary for the library selection.

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

Figure Legends

Figure 1: A generalised metagenomics data analysis workflow in the context of other 'omics' approaches.

Figure 2: A common data model for read data and associated metadata.

Figure 3: Schematic overview of best practice for analysis metadata collection with example fields.

Tables

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 [50]
collection date	ISO8601 date and time
environment (biome)	ENVO class [51]
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.

M2B3-mandatory saline water sample provenance descriptors	descriptor format
INVESTIGATION_campaign	text
INVESTIGATION_site	text
INVESTIGATION_platform	SDN:L06 controlled vocabulary [52]
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 [53], SDN:P06 [54] 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 [55]	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 [56]	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 [57]
library source	controlled vocabulary [58]
library strategy	controlled vocabulary [59]
library selection	controlled vocabulary [60]
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

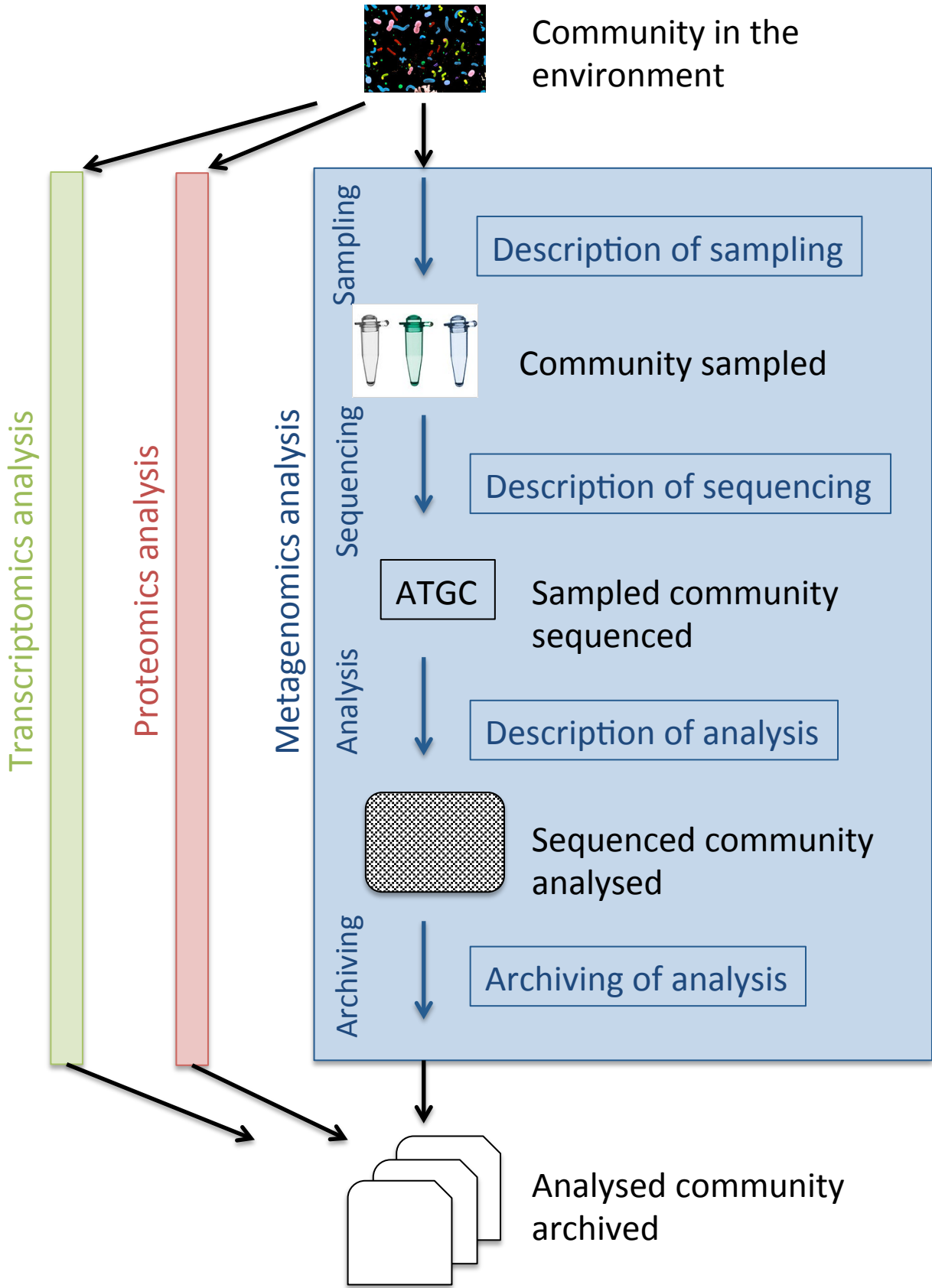
Table 5: Non-mandatory sequencing descriptors (column A) and formats (column B).

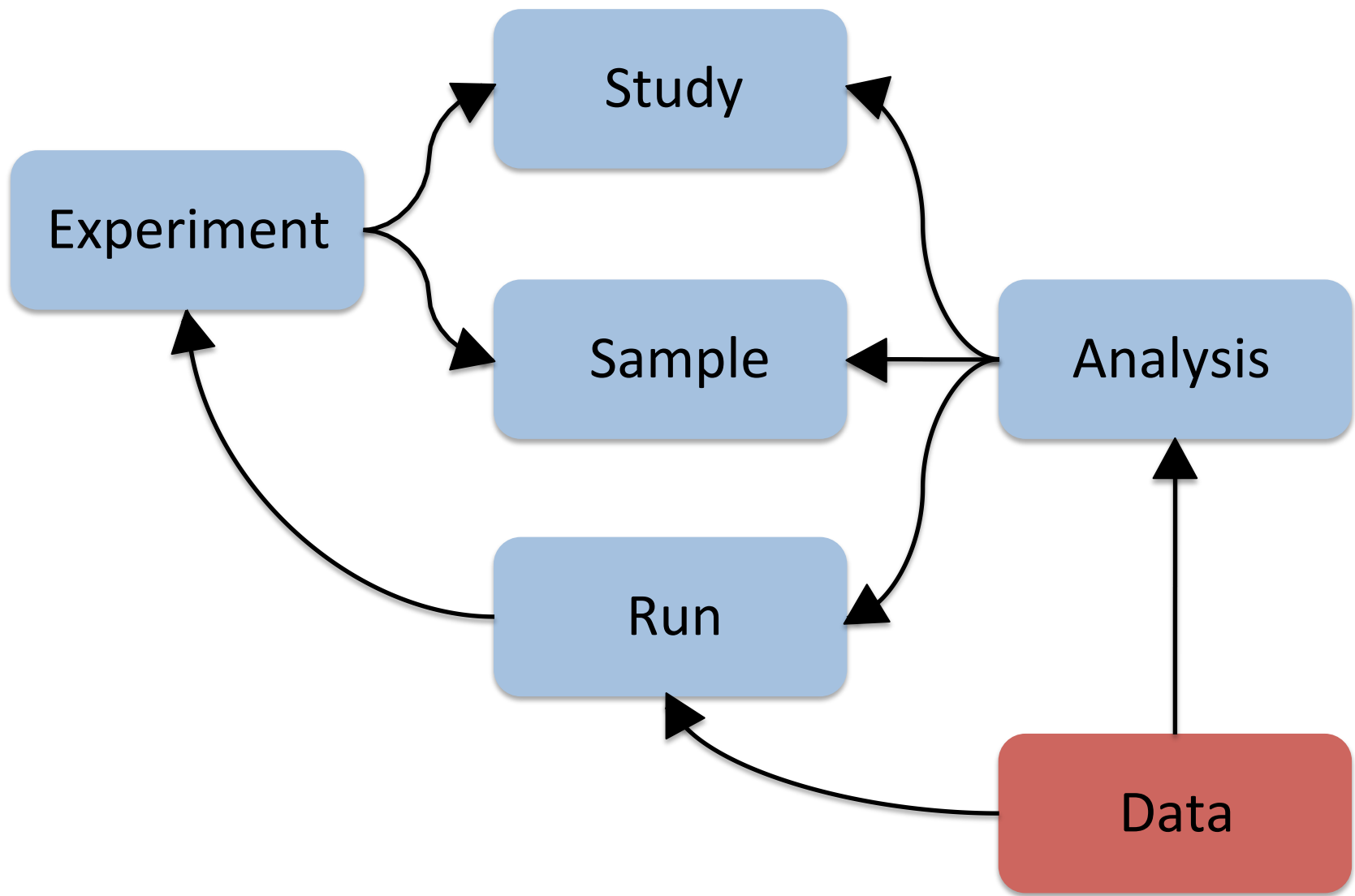
Column C suggests the descriptor’s potential importance for metagenomic data analysis (H – high relevance, M – medium relevance, L – low relevance).

A non-mandatory descriptors of sequencing provenance	B descriptor format	C 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	
Name: InterProScan	
Input(s)	Output 1 Component: FragGeneScan
	Reference DBs: InterPro 58.0 signature lib.
Tools	Name: InterPro
	Version: 5.19-58.0
	Source Code: https://www.ebi.ac.uk/interpro/interproscan.html
	Parameters
Output(s)	<ol style="list-style-type: none">1. Hit XML2. GO slim

C

Component Dependencies

