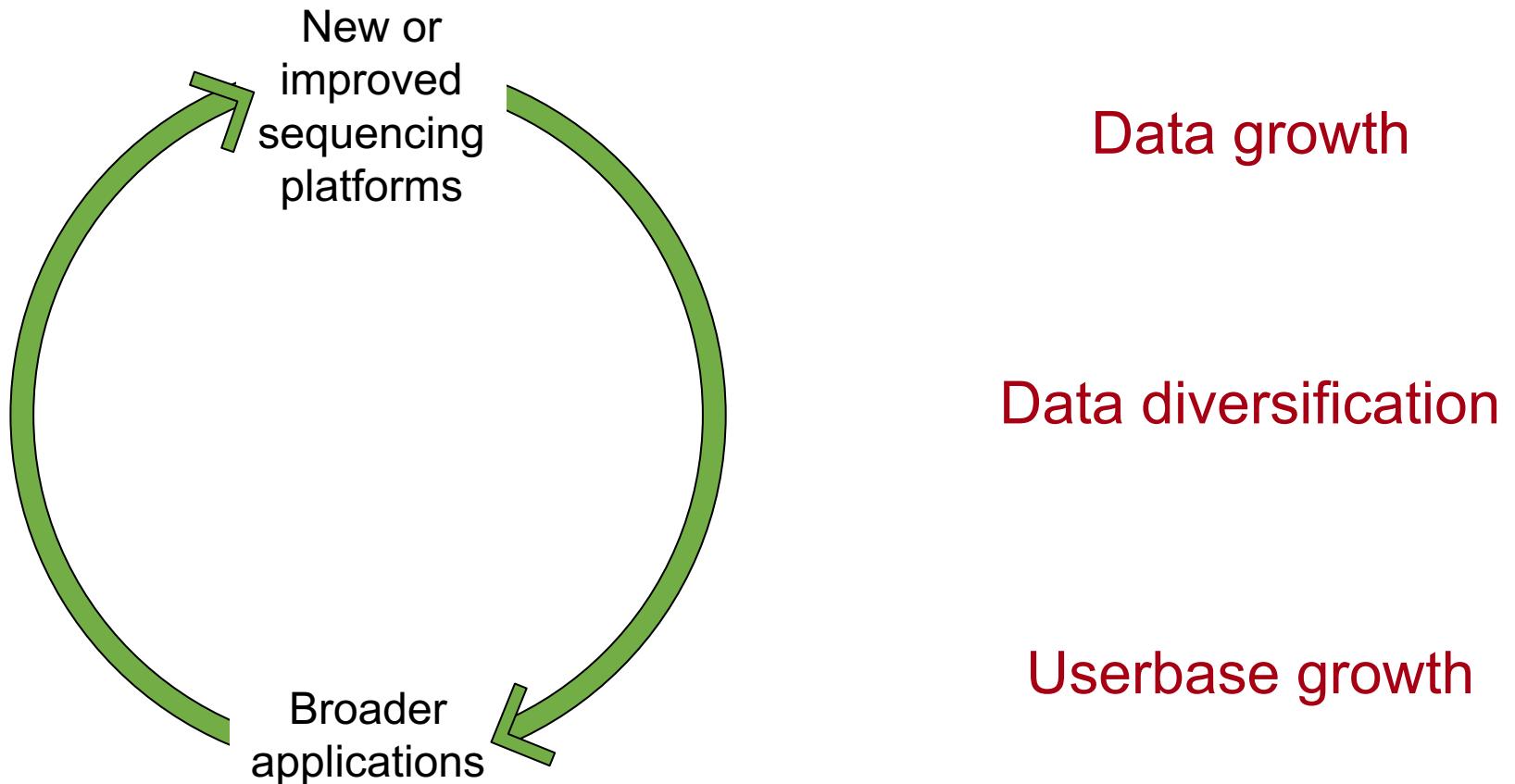


Sequencing communities at scale: dealing with the data

Guy Cochrane





Public environmental sequencing data up to 2002

- Dominant data from small-scale 16S rDNA studies
- Single sequencing platform
- ~35k sequence records
- ~30 kilobase pairs; ~0.5 megabytes

ID AF142985; SV 1; linear; genomic DNA; STD; ENV; 914 BP.
XX
AC AF142985;
XX
DT 06-JUL-1999 (Rel. 60, Created)
DT 10-MAR-2017 (Rel. 132, Last updated, Version 5)
XX
DE Uncultured archaeon ORGANIC1_A 16S ribosomal RNA gene, partial sequence.
XX
KW ENV.
XX
OS uncultured archaeon ORGANIC1_A
OC Archaea; Euryarchaeota; Halobacteria; Haloferacales; Halorubraceae;
OC Halorubrum; environmental samples.
XX
RN [1]
RP 1-914
RX DOI; 10.1046/j.1462-2920.2000.00097.x.
RX PUBMED; 11220308.
RA Bowman J.P., Rea S.M., McCammon S.A., McMeekin T.A.;
RT "Diversity and community structure within anoxic sediment from marine
RT salinity meromictic lakes and a coastal meromictic marine basin, Vestfold
RT Hills, Eastern Antarctica";
RL Environ. Microbiol. 2(2):227-237(2000).
XX
RN [2]
RP 1-914
RA Bowman J.P., McCammon S.A., McMeekin T.A.;
RT ;
RL Submitted (13-APR-1999) to the INSDC.
RL School of Agricultural Science, University of Tasmania, GPO Box 252-54,
RL Hobart, Tasmania 7001, Australia
XX
DR MD5: 489e036d498b21932c168a5dc2c78a5f.
DR SILVA-SSU; AF142985.
XX
FH Key Location/Qualifiers
FH
FT source 1..914
FT /organism="uncultured archaeon ORGANIC1_A"
FT /environmental_sample
FT /mol_type="genomic DNA"
FT /isolation_source="sediment"
FT /clone="ORGANIC1_A"
FT /db_xref="taxon:98847"
FT <1..>914
FT /product="16S ribosomal RNA"
XX
SQ Sequence 914 BP; 233 A; 209 C; 297 G; 175 T; 0 other;
attctggttt atctcgccatg aggacatgtc tattgggatt cgatcttagcc atgtctagtcg
caccgggtta gactctgtgc agatagtcga gttaacacgtg gccaacgtac cttccatggcc
agaataaacct cgggaaaactg aggtcaagac tggataacgtt atgcacgtcg gaatgcaggaa
tattccaaac gtcggcggcg tgaaggatac gggtcgccg gatttagttt aeggttaggt
aacggcccad ctgtggcggta atccgttccg gtcatcgatg tgagggcccg gagacyggat
ctggagacaag attccggggcc ctacggggcc cagcaggccg gaaacccctta cactgcacgc
aagtgcgata ggggaatccc aagtttgatgg gcatatagcc ctgcgttttt tcgactgtaa
ggagggtcaac tggcaagac cggtgcggat ggcggccgat ataccggat
ctcgagttgt gtgttatii attggccatc aagctggccgt aegtggcccg cgaatgtcc
cggggaaatcc acttgcoccaat cgggtggccg tccggcggaaat actgtccggc ttggaaacccg
aaggctcaga ggttgcgtct ggggttggag taaatctt taatccccga cggacgcacccg
atggcgaag cagtctgaga ggacggatcc gacagtgggg gacgaaagct agggtctcga
accggatgg ataccggatgg atggctcgtat gtaaaggatg ctgcgttagat gtggcacccca
ctacggatgg gtgtctgttc gtggaaaga agttaaacgca gcccgttggg aagtactgtcc
gcaaggatgt aacttaaagg aattggccggg ggagcactac aaccggagga gcatcgccgtt
taattttttt caac // 60
120
180
240
300
420
480
540
600
660
720
780
840
900
914

Public environmental sequencing data in 2017

- 16,000 environmental sequencing studies
- 600,000 metagenomics samples
- diverse data
 - raw data
 - 454, Illumina, PacBio, Ion Torrent, ONT, etc.
 - sequenced libraries (600,000)
 - 30 terabase pairs, ~0.5 petabytes
 - shotgun metagenomics
 - metatranscriptomics
 - metagenome/metatranscriptome assemblies (1,100)
 - single cell genomes
 - reference data (e.g. 100,000 bacterial assemblies)
 - identification data (243)
 - gene catalogue (1)

Context

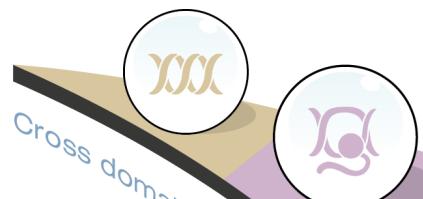
EMBL European Bioinformatics Institute

Genes, genomes & variation

European Nucleotide Archive
1000 Genomes

Ensembl
Ensembl Genomes

European Genome-phenome Archive
Metagenomics portal



Literature & ontologies

Europe PubMed Central
Gene Ontology
Experimental Factor
Ontology



Reactions, interactions & pathways

IntAct

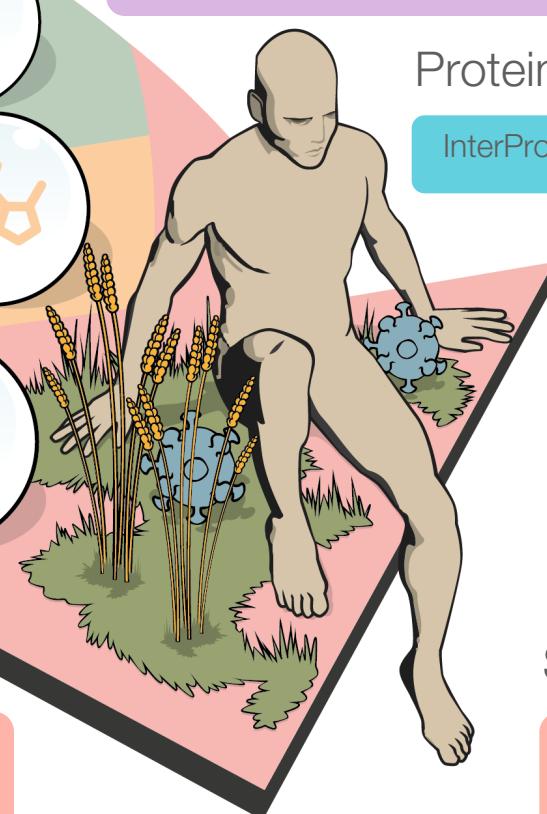
Reactome

MetaboLights

Gene, protein & metabolite expression

ArrayExpress
Expression Atlas

Metabolights
PRIDE



Protein sequences, families & motifs

InterPro

Pfam

UniProt

Molecular structures

Protein Data Bank in Europe
Electron Microscopy Data Bank

Chemical biology

ChEMBL

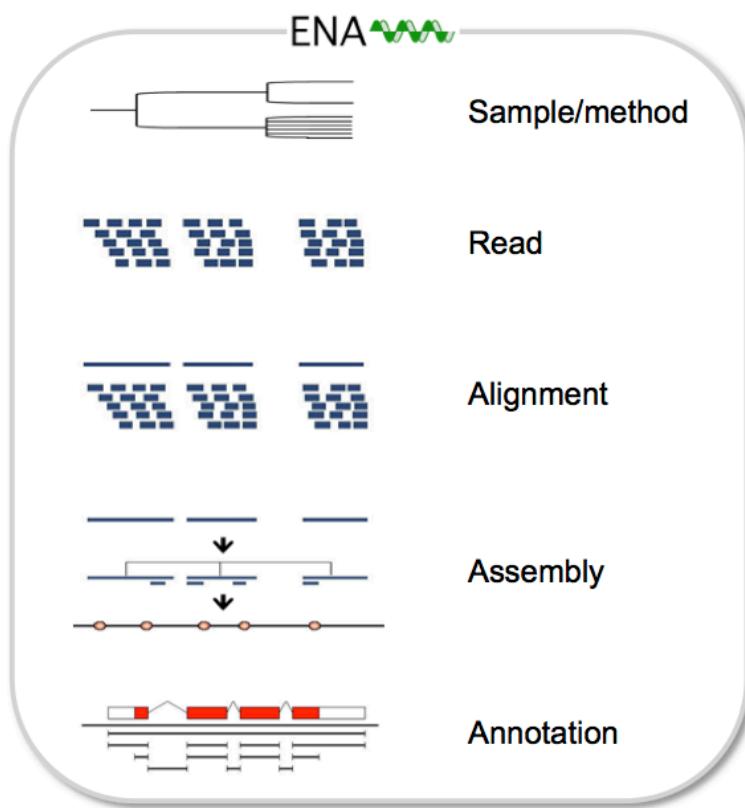
ChEBI

Systems

BioModels
Enzyme Portal

BioSamples

European Nucleotide Archive (ENA)



ENA

<http://www.ebi.ac.uk/ena/>

- Globally comprehensive scientific record and European node of **INSDC**
- A **broad platform** for the management, sharing, integration and dissemination of sequence data
- Established in the early 1980s, extended for **new technologies and applications**
- **Connectivity** with broader EMBL-EBI resources
- Sequence data **foundation**
- **Sustained** within EMBL-EBI under EMBL funding with additional support from EC, UK Research councils, Wellcome Trust, etc.
- **Substantial scale**: 1 submission every 6 minutes, 1.3 petabase pairs across 1.5 million taxa, 2,000-5,000 active data providers, global consumer userbase
- Rich submission, discovery and retrieval **software, tools and services**

Data validation and submissions

Home New Submission Studies Sample Groups Samples Experiments Runs Projects

Start >> Sample >> Finish

Please create new samples by uploading a spreadsheet or by following the instructions below.

Please select the checklist that you wish to use for your sample submission

If you already have a spreadsheet containing your data upload it here.

Default Checklist

Upload Spreadsheet

Book1 - Microsoft Excel

File Home Insert Page Layout Formulas Data Review View

From Access From Web From Text From Other Sources Existing Connections Refresh All Connections

A Z A Z A Z Sort Filter Advanced Text to Columns Remove Duplicates Data Validation Consolidate What-If Analysis Group Ungroup Subtotal Outline

#checklist_accession

	A	B	C	D	E	F	G	H	I	J
1	#checklist_accession	ERC000001								
2	#unique_name_prefix	mouse_dendrocyte_								
3	sample_alias	tax_id	scientific_name	common_name	anonymized_name	sample_title	sample_description	tissue_type	sex	collection
4	#template		10090	Mus musculus	house mouse		mouse dendrocyte	A dendrocyte brain sample was extracted from a mouse	brain	male
5	#units									13/0
6	mouse_dendrocyte_1		10090	Mus musculus	house mouse		mouse dendrocyte	A dendrocyte brain sample was extracted from a mouse	brain	male
7	mouse_dendrocyte_2		10090	Mus musculus	house mouse		mouse dendrocyte	A dendrocyte brain sample was extracted from a mouse	brain	male
8	mouse_dendrocyte_3		10090	Mus musculus	house mouse		mouse dendrocyte	A dendrocyte brain sample was extracted from a mouse	brain	male
9	mouse_dendrocyte_4		10090	Mus musculus	house mouse		mouse dendrocyte	A dendrocyte brain sample was extracted from a mouse	brain	male
10	mouse_dendrocyte_5		10090	Mus musculus	house mouse		mouse dendrocyte	A dendrocyte brain sample was extracted from a mouse	brain	male
11	mouse_dendrocyte_6		10090	Mus musculus	house mouse		mouse dendrocyte	A dendrocyte brain sample was extracted from a mouse	brain	male
12	mouse_dendrocyte_7		10090	Mus musculus	house mouse		mouse dendrocyte	A dendrocyte brain sample was extracted from a mouse	brain	male
										13/0

Sheet1 Sheet2 Sheet3

Ready 85%

Data access

Data discovery

Country

Geographical location

Search by: Bounded box Radius (km)

Map Satellite

Collection date

Select search conditions:

Taxonomy and related

- Taxon name = wastewater metagenome Include subordinate taxa
- Environmental sample = True False
- Strain =
- Sub-strain =

tax_eq(527639)
AND
collection_date>=2016-01-01
AND
geo_box1(55.37, 12.08, 55.79, 13.05)

Search query
tax_eq(527639) AND collection_date>=2016-01-01 AND geo_box1(55.371271003687205, 12.0830078125, 55.790069013291166, 13.052734375)

Search results for tax_eq(527639) AND collection_date>=2016-01-01 AND geo_box1(55.371271003687205, 12.0830078125, 55.790069013291166, 13.052734375)

Sample		Records		Reports			
Sample (47)		<input type="checkbox"/> Accession <input type="checkbox"/> Cell type <input type="checkbox"/> Cultivar <input type="checkbox"/> Ecotype <input type="checkbox"/> Identified by <input type="checkbox"/> Mating type <input type="checkbox"/> Sex (as submitted) <input type="checkbox"/> Sub-strain <input type="checkbox"/> Tax ID <input type="checkbox"/> Center name <input type="checkbox"/> Environment (Biome) <input type="checkbox"/> Salinity (PSU) <input type="checkbox"/> Protocol Label <input type="checkbox"/> Host status <input type="checkbox"/> Host gravity <input type="checkbox"/> Environmental package <input type="checkbox"/> Sequencing method		<input type="checkbox"/> Secondary sample accession <input type="checkbox"/> Collected by <input type="checkbox"/> Culture collection <input type="checkbox"/> Environmental sample <input type="checkbox"/> Isolate <input type="checkbox"/> Serotype <input type="checkbox"/> Specimen voucher <input type="checkbox"/> Tissue library <input checked="" type="checkbox"/> Scientific name <input type="checkbox"/> Depth (m) <input type="checkbox"/> Environment (Feature) <input type="checkbox"/> Sampling Campaign <input type="checkbox"/> Project Name <input type="checkbox"/> Host sex <input type="checkbox"/> Host phenotype <input type="checkbox"/> Investigation type <input type="checkbox"/> Target gene		<input type="checkbox"/> Bio material <input type="checkbox"/> Collection date <input type="checkbox"/> Description <input type="checkbox"/> First public <input type="checkbox"/> Isolation source <input type="checkbox"/> Strain <input type="checkbox"/> Tissue type <input type="checkbox"/> Submitter's sample name <input type="checkbox"/> Elevation (m) <input type="checkbox"/> Environment (Material) <input type="checkbox"/> Sampling Site <input type="checkbox"/> Host <input type="checkbox"/> Host sex (as submitted) <input type="checkbox"/> Host genotype <input type="checkbox"/> Experimental factor <input type="checkbox"/> pH	

Showing results 1 - 10 of 47 results

Accession	Collection date	Country	First public	Geographical location	Scientific name	Investigation type
SAMEA404463	2016-04-07	Denmark	2016-06-22	55.608536 N 12.450516 E	wastewater metagenome	metagenome
SAMEA404464	2016-04-11	Denmark	2016-06-22	55.694984 N 12.615026 E	wastewater metagenome	metagenome
SAMEA404465	2016-04-07	Denmark	2016-06-22	55.694984 N 12.615026 E	wastewater metagenome	metagenome
SAMEA404466	2016-04-13	Denmark	2016-06-22	55.694984 N 12.615026 E	wastewater metagenome	metagenome
SAMEA404467	2016-04-17	Denmark	2016-06-22	55.694984 N 12.615026 E	wastewater metagenome	metagenome
SAMEA404468	2016-03-29	Denmark	2016-06-22	55.640495 N 12.505821 E	wastewater metagenome	metagenome
SAMEA404469	2016-04-02	Denmark	2016-06-22	55.640495 N 12.505821 E	wastewater metagenome	metagenome

Showing results 1 - 10 of 47 results

Accession	Collection date	Country	First public	Geographical location	Scientific name	Investigation type
SAMEA404463	2016-04-07	Denmark	2016-06-22	55.608536 N 12.450516 E	wastewater metagenome	metagenome
SAMEA404464	2016-04-11	Denmark	2016-06-22	55.694984 N 12.615026 E	wastewater metagenome	metagenome
SAMEA404465	2016-04-07	Denmark	2016-06-22	55.694984 N 12.615026 E	wastewater metagenome	metagenome
SAMEA404466	2016-04-13	Denmark	2016-06-22	55.694984 N 12.615026 E	wastewater metagenome	metagenome
SAMEA404467	2016-04-17	Denmark	2016-06-22	55.694984 N 12.615026 E	wastewater metagenome	metagenome
SAMEA404468	2016-03-29	Denmark	2016-06-22	55.640495 N 12.505821 E	wastewater metagenome	metagenome
SAMEA404469	2016-04-02	Denmark	2016-06-22	55.640495 N 12.505821 E	wastewater metagenome	metagenome

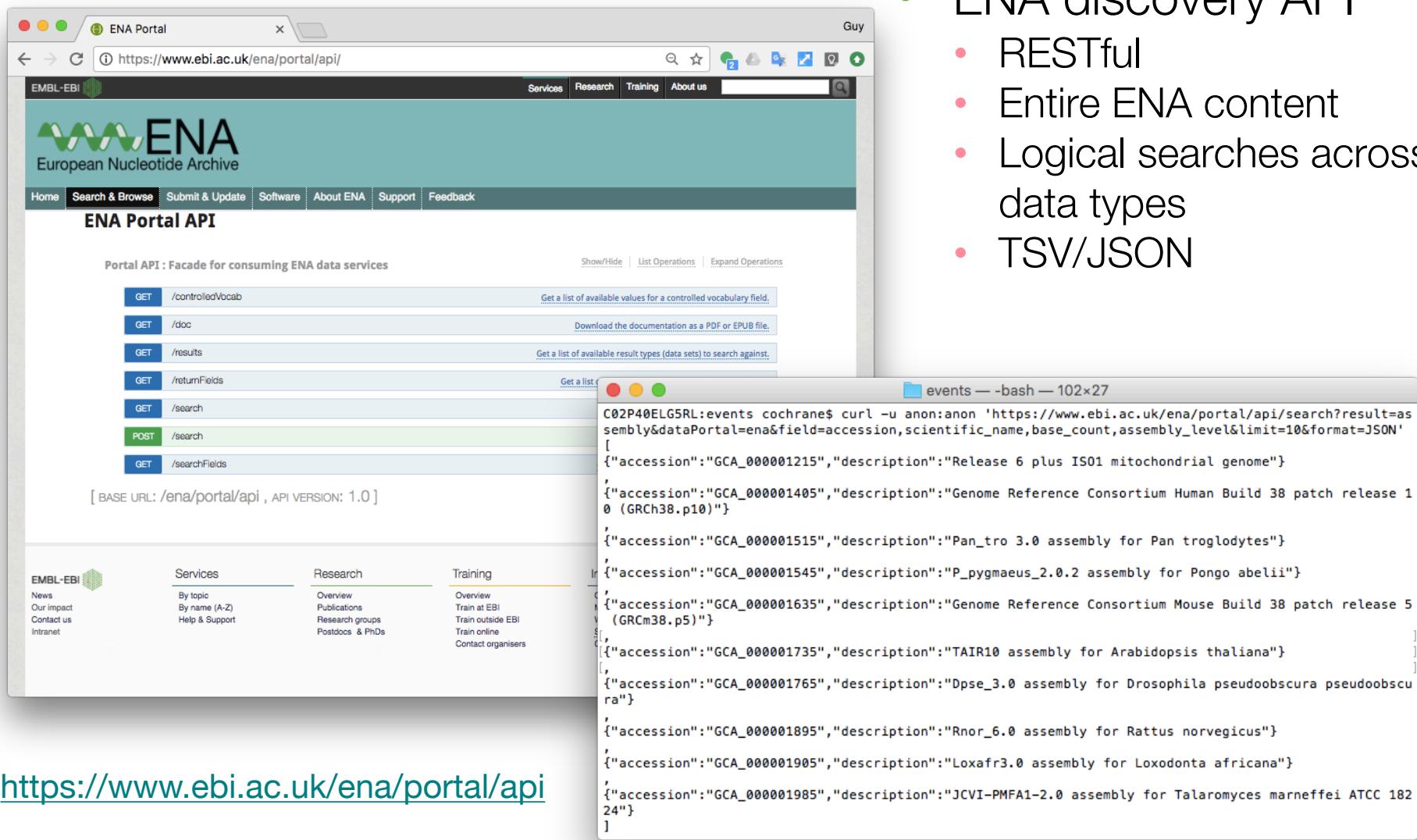
Navigation Attributes

Location on map

Map Satellite

Map data ©2016 Google Report a map error

Programmatic access



Public environmental sequencing data in 2017

- 16,000 environmental sequencing studies
- 600,000 metagenomics samples
- diverse data
 - raw data
 - 454, Illumina, PacBio, Ion Torrent, ONT, etc.
 - sequenced libraries (600,000)
 - 30 terabase pairs, ~0.5 petabytes
 - shotgun metagenomics
 - metatranscriptomic
 - metagenome/metatranscriptome assemblies (1,100)
 - single cell genomes
 - reference data (e.g. 100,000 bacterial assemblies)
 - identification data (243)
 - gene catalogue (1)

ENA as data management platform for metagenomics

- In-project data sharing
- Data validation
- Dissemination and publishing
- Data coordination
- for customers small....,
- medium...



UniEuk

... and large

The screenshot shows the EBI Metagenomics homepage. At the top, there's a navigation bar with links for Services, Research, Training, and About us. Below the header, a main banner features a magnifying glass over a globe icon, with the text "Submit, analyse, visualize and compare your data." and a "SUBMIT DATA" button. To the right of the banner, there are four data summary boxes:

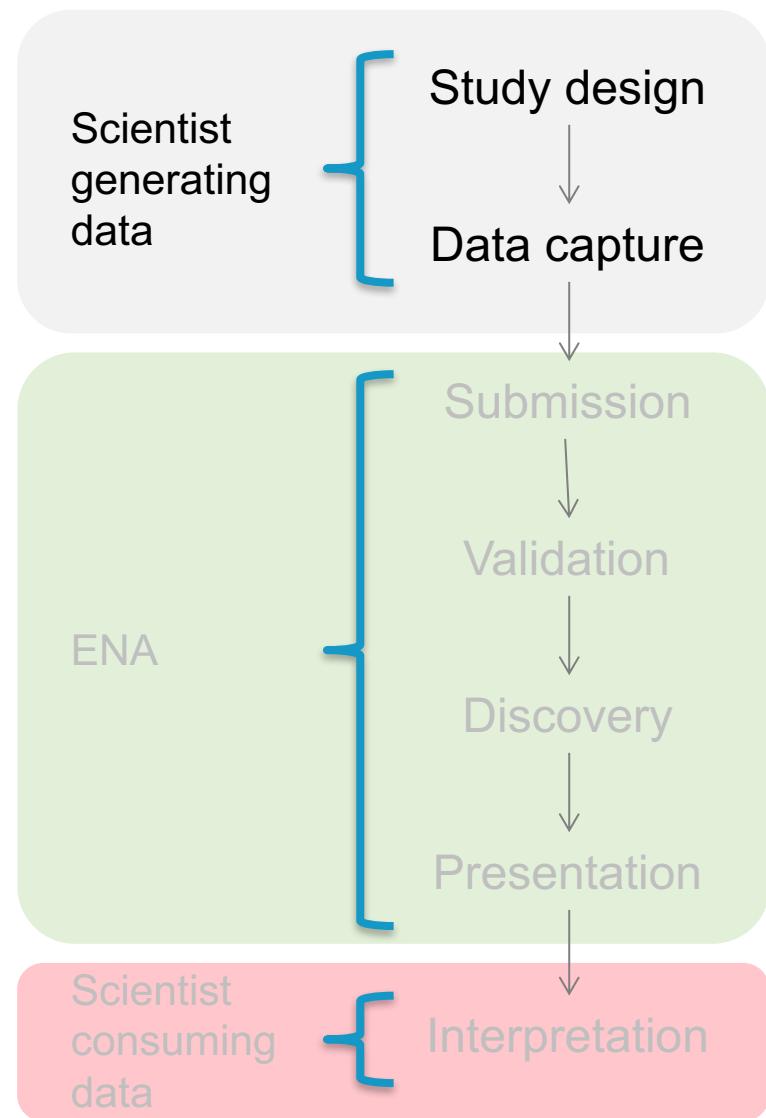
Category	Count	Type
data sets	107123	
amplicons	75946	
assemblies	102	
metabarcoding	1258	
metagenomes	16973	
metatranscriptomes	1235	
runs	93874	
samples	72211	
projects	1148	
runs	6288	
samples	6060	
projects	166	

Below the summary boxes, there's a section titled "Browse projects" with a "By selected biomes" grid. The grid includes icons for Soil, Engineered, Marine, Host-associated human, and Host-associated plant environments, each with a count of samples. To the right, there's a "Latest projects" section with two entries:

- Characterization of tissue microbiota using the Illumina MiSeq sequencing technology after pipeline validation on Mock communities**
Characterization of tissue microbiota using 16S metagenomic sequencing by Illumina MiSeq of: - feces - ileum - adipose tissue - heart - liver - brain - muscle The validation of the ...
View more - 181 samples
- Substrate variations triggered the emergent of different active bacterial and archaeal assemblages during biomethane production**
Biomethane has been regarded as one of the promising renewable energy supplies. However, the microbial communities involved in methane synthesis have not been fully characterized. By ...
View more - 8 samples

Standards

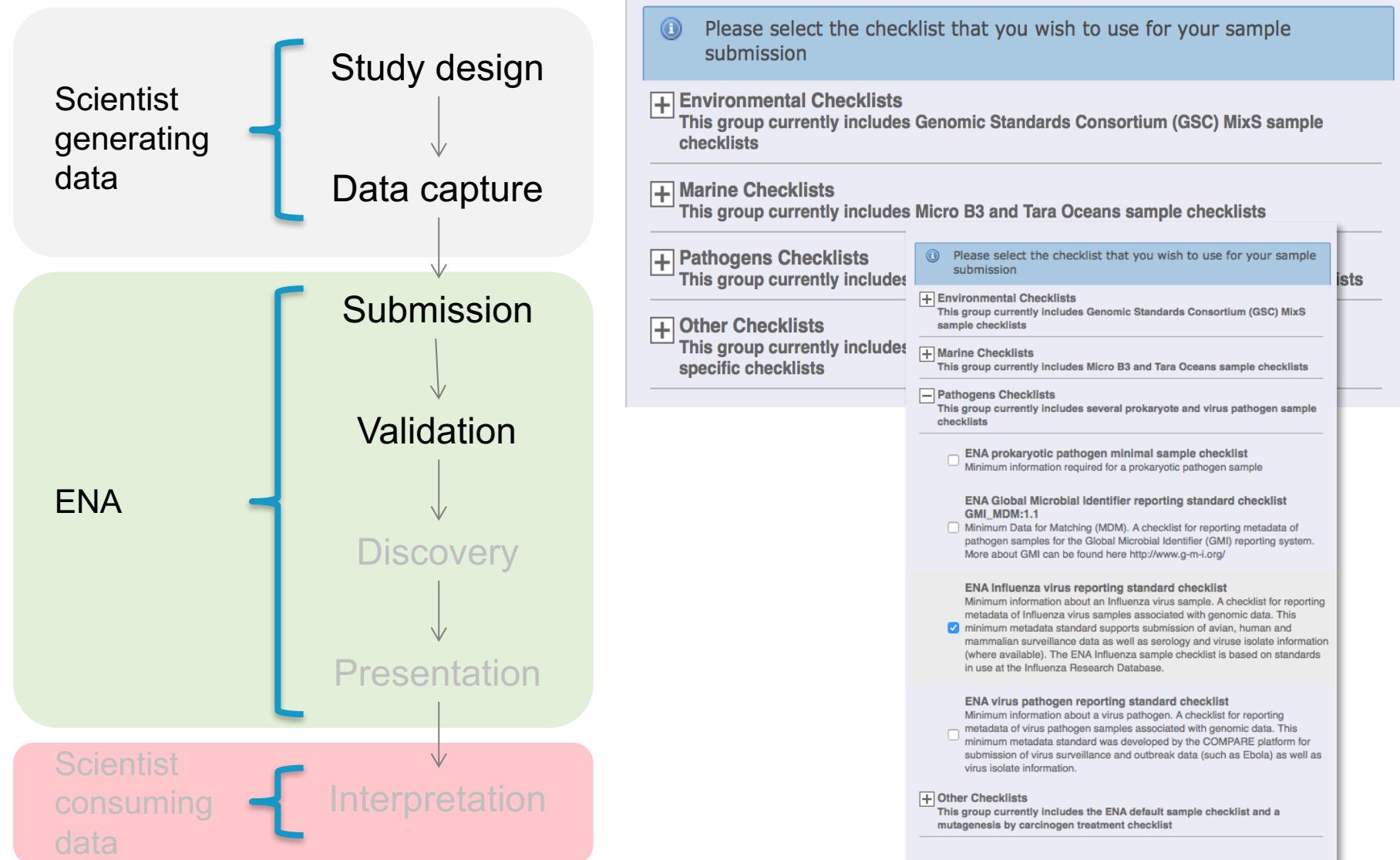
Standards through the data journey



The screenshot shows the ENA website interface with the following sections:

- Header:** EMBL-EBI, ENA European Nucleotide Archive, Services, Research, Training, About us.
- Search Bar:** Examples: BN000065, hostone, Search, Advanced Sequence.
- Checklist:** ERC000012 (GSC MIxS air)
- Description:** Genomic Standards Consortium package extension for reporting of measurements and observations obtained from the environment where the sample was obtained. By choosing the environmental package, a selection of fields can be made from a relevant subsets of the GSC terms.
- geography:** A table showing fields for environmental package, geographic location (latitude), and geographic location (longitude).
- local environment conditions:** A table showing fields for barometric pressure, humidity, pollutants, and solar irradiance.
- Associated host information schema:** A detailed diagram showing the schema for "host_associated ('yes' or 'no')". If YES: Associated host information includes host_disease_stage, host_disease_outcome, host_health_state, isolation_source, host_description, and other descriptive information relating to the host. If NO: Environmental information includes isolation_source (describes the physical, environmental and/or local where the organism was derived) and lat_lon OR country (geographical coordinates of the location where the organism indicated in terms of political names for).

Standards through the data journey



Standards through the data journey

Scientist generating data

Study design

Data capture

ENA

Submission

Validation

Discovery

Presentation

Scientist consuming data

Interpretation

ENEA European Nucleotide Archive

Please subscribe to ena-announce mailing list here: listserver.ebi.ac.uk/mailman/listinfo/ena-announce, to receive alerts about ENA services.

Advanced Upload accession

Search query

tax_tree(11320) AND geo_circ(52.069881131843665, 5.097687499999893, 269.7631590474306) AND country="The Netherlands" AND host_status="diseased" AND checklist="ERC000032"

Select domain:

- Assembly
- Sequence
- Contig set
- Coding
- Non-coding
- Read
- Analysis
- Trace
- Study
- Taxon
- Sample**
- Environmental
- Marker

Select search conditions:

Taxonomy and related

Taxon name = Influenza A virus
Include subordinate taxa
NCBI Catalogue of Life

Environmental sample = True False

Strain =

Collected by =

Identified by =

Country = The Netherlands

Geographical location

Search by: Bounded box Radius (km)

Map Satellite

Centre point Latitude 52.069881131Longitude 5.0976874999 Radius (km) 269.7631590

Google

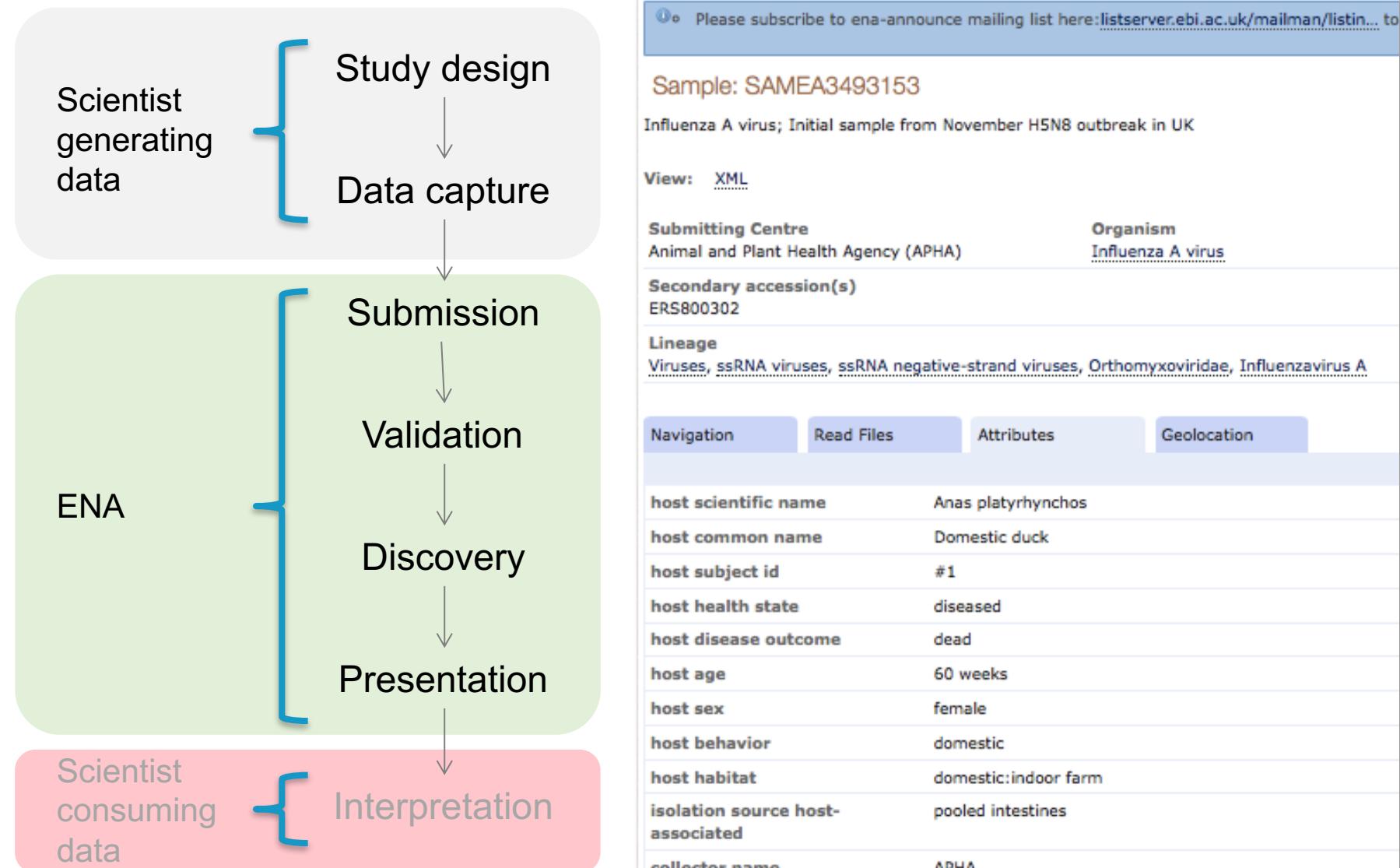
Depth (m) =

Elevation (m) =

Altitude (m) =

Environment (Biome) =

Standards through the data journey



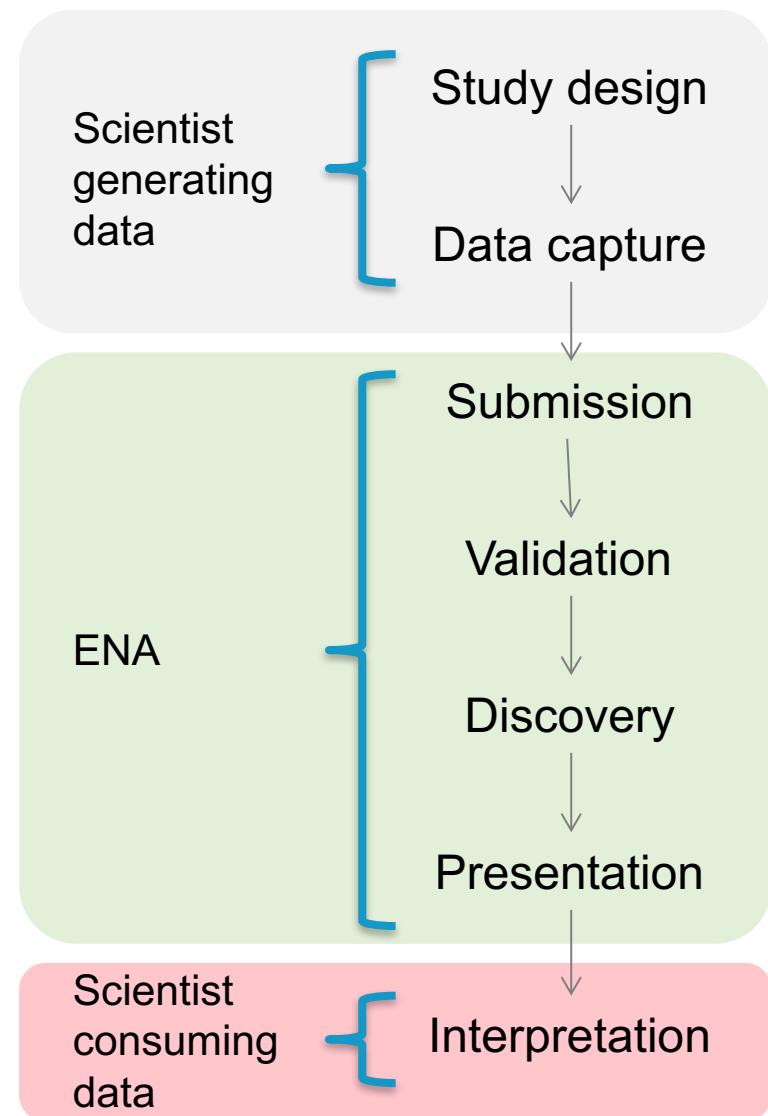
Download: 1 - 2 of 2 results in TEXT

[Hide](#) [Select columns](#)

<input checked="" type="checkbox"/> Accession	<input type="checkbox"/> Secondary sample accession	<input type="checkbox"/> Bio material	<input type="checkbox"/> Cell line
<input type="checkbox"/> Cell type	<input type="checkbox"/> Collected by	<input type="checkbox"/> Collection date	<input checked="" type="checkbox"/> Country
<input type="checkbox"/> Cultivar	<input type="checkbox"/> Culture collection	<input type="checkbox"/> Description	<input type="checkbox"/> Developmental stage
<input type="checkbox"/> Ecotype	<input type="checkbox"/> Environmental sample	<input checked="" type="checkbox"/> First public	<input type="checkbox"/> Germline
<input type="checkbox"/> Identified by	<input type="checkbox"/> Isolate	<input type="checkbox"/> Isolation source	<input checked="" type="checkbox"/> Geographical location
<input type="checkbox"/> Mating type	<input type="checkbox"/> Serotype	<input type="checkbox"/> Serovar	<input type="checkbox"/> Sex
<input type="checkbox"/> Sex (as submitted)	<input type="checkbox"/> Specimen voucher	<input type="checkbox"/> Strain	<input type="checkbox"/> Sub-species
<input type="checkbox"/> Sub-strain	<input type="checkbox"/> Tissue library	<input type="checkbox"/> Tissue type	<input type="checkbox"/> Variety
<input type="checkbox"/> Tax ID	<input checked="" type="checkbox"/> Scientific name	<input type="checkbox"/> CoL tax ID	<input type="checkbox"/> CoL scientific name
<input type="checkbox"/> Submitter's sample name	<input type="checkbox"/> Checklist	<input type="checkbox"/> Center name	<input type="checkbox"/> Depth (m)
<input type="checkbox"/> Elevation (m)	<input type="checkbox"/> Altitude (m)	<input type="checkbox"/> Environment (Biome)	<input type="checkbox"/> Environment (Feature)
<input type="checkbox"/> Environment (Material)	<input type="checkbox"/> Temperature (C)	<input type="checkbox"/> Salinity (PSU)	<input type="checkbox"/> Sampling Campaign
<input type="checkbox"/> Sampling Site	<input type="checkbox"/> Sampling Platform	<input type="checkbox"/> Protocol Label	<input type="checkbox"/> Project Name
<input checked="" type="checkbox"/> Host	<input type="checkbox"/> Host tax id	<input checked="" type="checkbox"/> Host Status	<input checked="" type="checkbox"/> Host sex
<input type="checkbox"/> Host sex (as submitted)	<input type="checkbox"/> Host body site	<input type="checkbox"/> Host gravidity	<input type="checkbox"/> Host phenotype
<input type="checkbox"/> Host genotype	<input type="checkbox"/> Host growth conditions	<input type="checkbox"/> Environmental package	<input type="checkbox"/> Ingestion type

Accession	Collection date	Country	First public	Geographical location	Scientific name	Submitter's sample name	Host	Host Status	Host sex
SAMEA3493153	2014-11-14	United Kingdom	2016-02-17	54.001 N 0.434 W	Influenza A virus	AVP-14-036255	Anas platyrhynchos	diseased	female
SAMEA3860663	2014-11-04	Germany	2016-02-17		Influenza A virus	A/turkey/Germany/AR2485-L01478/2014	Meleagris gallopavo	diseased	

Standards through the data journey



Please subscribe to ena-announce mailing list here: listserver.ebi.ac.uk/mailman/listinfo/ena-announce	
Sample: SAMEA3493153	
Influenza A virus; Initial sample from November H5N8 outbreak in UK	
View:	XML
Submitting Centre	Animal and Plant Health Agency (APHA)
Organism	Influenza A virus
Secondary accession(s)	ERS800302
Lineage	Viruses, ssRNA viruses, ssRNA negative-strand viruses, Orthomyxoviridae, Influenzavirus A
Navigation Read Files Attributes Geolocation	
host scientific name	Anas platyrhynchos
host common name	Domestic duck
host subject id	#1
host health state	diseased
host disease outcome	dead
host age	60 weeks
host sex	female
host behavior	domestic
host habitat	domestic:indoor farm
isolation source host-associated	pooled intestines
collector name	APHA

Look out for...

ACCEPTED MANUSCRIPT

The metagenomic data life-cycle: standards and best practices

Petra ten Hoopen, Robert D. Finn, Lars Ailo Bongo, Erwan Corre, Bruno Fosso, Folker Meyer, Alex Mitchell, Eric Pelletier, Graziano Pesole, Monica Santamaria, Nils Peder Willassen, Guy Cochrane 

Gigascience gix047. DOI: <https://doi.org/10.1093/gigascience/gix047>

Published: 16 June 2017

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

quencing (3) data analysis and (4) data archiving

OXFORD
ACADEMIC

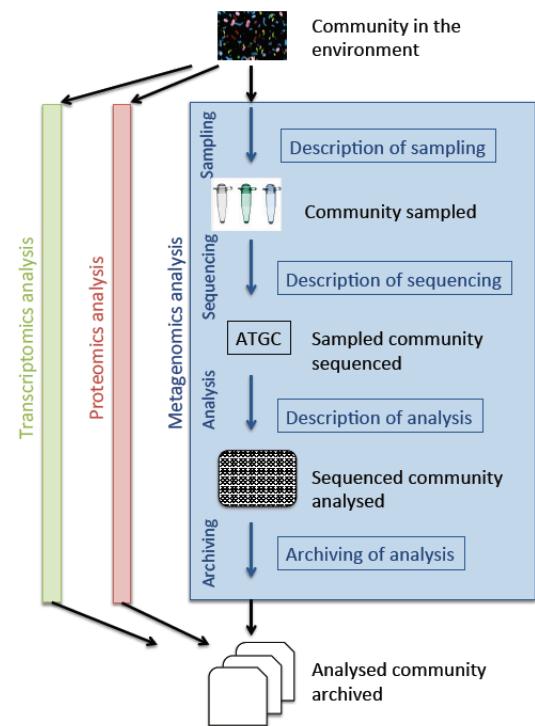
(GIGA)ⁿ
SCIENCE

generation have been to some extent addressed by the scientific

marine research, we summarise essential variables sampling processes and sequencing procedures iment. These aspects of metagenomics dataset

analysis have been to some extent addressed by the scientific

- Authors from 9 institutions, lead author Petra ten Hoopen, input from marine-community@elixir-europe.org
- Covers full workflow from sampling, through sequencing, analysis and result publication
- Promotes systematic archiving of analysis outputs



Scaling: internal tools

Sustainable biocuration

...to

From...

Checklist Editor Checklist Groups Checklists Field Groups Fields Logout

Create Edit Delete

Minimum information about a Tara Oceans sample. A checklist for reporting metadata of oceanic plankton samples associated with genomic data from t...

ENa virus pathogen reporting standard checklist 35 fields

Minimum information about a virus pathogen. A checklist for reporting metadata of virus pathogen samples associated with genomic data. This minimum...

ENa Influenza virus reporting standard checklist 69 fields

Minimum information about an Influenza virus sample. A checklist for reporting metadata of Influenza virus samples associated with genomic ...

ENa multi carcinogen checklist 1 templates

<templates>
<template id="ERT800802" version="3">
 <token name="RNA gene"></forms_name>
 <description>For ribosomal RNA genes from prokaryotic, nuclear or organelar DNA. All rRNAs are considered partial.</description>
 <token name="LAB_ORGANISM_NAME" type="TAXON_FIELD" mandatory="true" display_name="Organism" description="Full name of organism (generally Genus+species). Forma...

All Fields Included Fields Filter: NONE

Associated host information

ON	host scientific name	Mandatory	Recommended	Single	Edit
		Optional		Multiple	

ON	host common name	Mandatory	Recommended	Single	Edit
		Optional		Multiple	

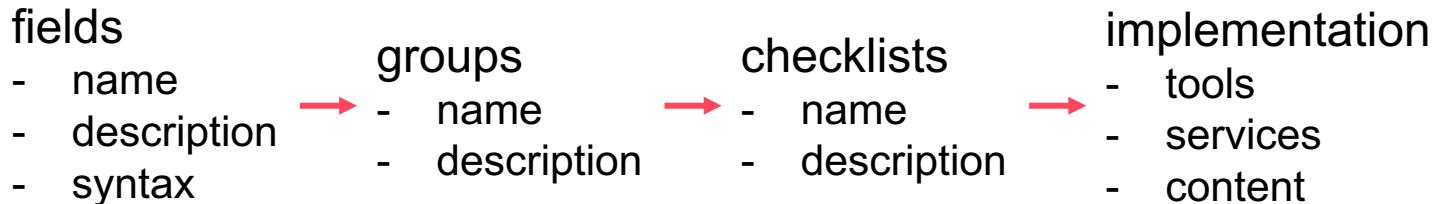
ON	host health state	Mandatory	Recommended	Single	Edit
		Optional		Multiple	

ON	host subject id	Mandatory	Recommended	Single	Edit
----	-----------------	-----------	-------------	--------	------

ena_sequence_templates.xml (no function selected) Not registered

- Submission-level processing
 - Class-level insertion of scientific organisation

Checklist Editor



The screenshot shows the Checklist Editor application interface. At the top, there is a navigation bar with tabs: Checklist Editor, Checklist Groups, Checklists (which is highlighted with a red box), Field Groups, and Fields. On the far right of the navigation bar is a Logout button.

The main content area is divided into two main sections:

- Checklists Section:** This section displays a list of pre-defined checklists:
 - ENAvirus pathogen reporting standard checklist (35 fields)
 - ENAIfluenza virus reporting standard checklist (59 fields)
 - ENAmutagenesis by carcinogen treatment checklist (11 fields)
- Field Groups Section:** This section displays a list of field groups under the heading "Associated host information". Each group has settings for "ON" (checkbox), "host scientific name" (text input), "Mandatory" (radio button), "Recommended" (radio button), "Single" (radio button), and "Edit" (button). Similar groups are listed for host common name, host health state, host subject id, host disease outcome, and host age.

Scaling: community engagement

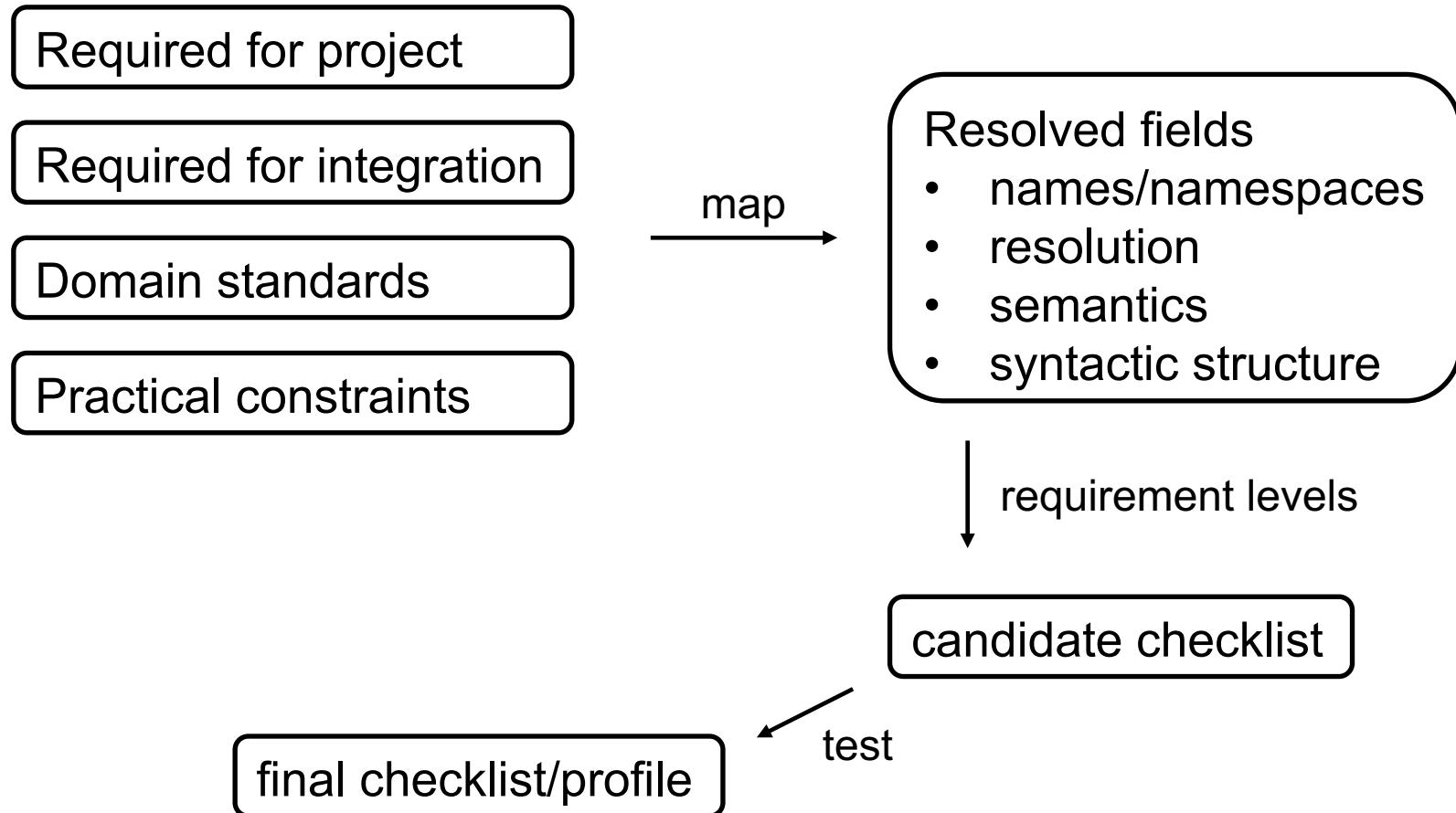
MIxS family of standards



<http://gensc.org>

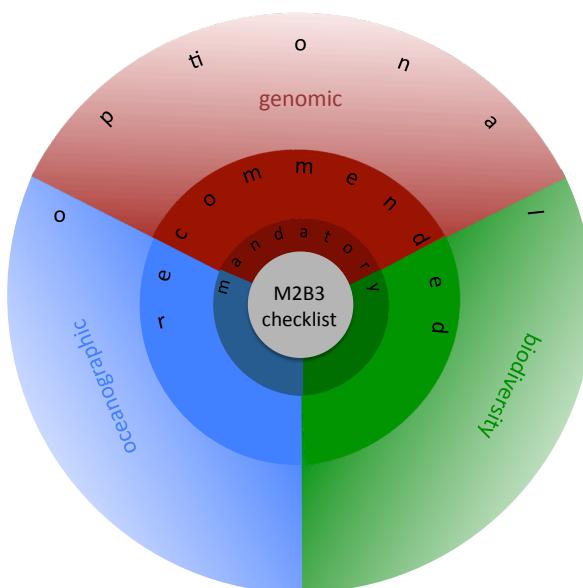
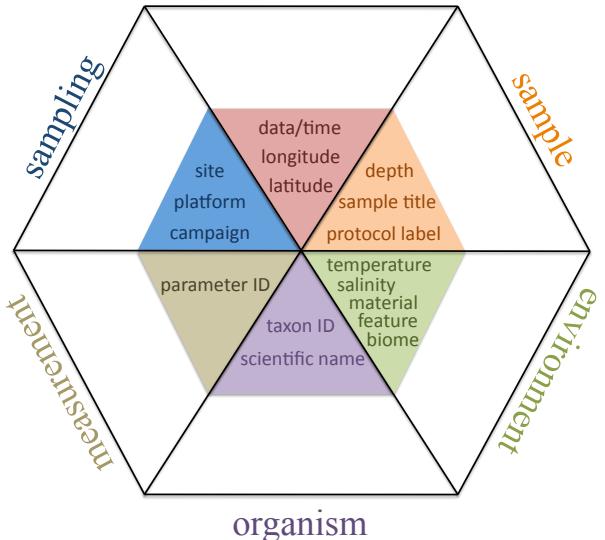
Specification projects	MIGS		MIMS	MIMARKS		New checklists
Checklists	EU	BA	PL	VI	ORG	metagenomes survey specimen e.g., pan-genomes
Shared descriptors	collection date, environmental package, environment (biome), environment (feature), environment (material), geographic location (country and/or sea, region), geographic location (latitude and longitude), investigation type, project name, sequencing method, submitted to INSDC					
Checklist-specific descriptors	assembly, estimated size, finishing strategy, isolation and growth condition, number of replicons, ploidy, propagation, reference for biomaterial				target gene	
Applicable environmental packages (measurements and observations)	Air Host-associated Human-associated Human-oral Human-gut Human-skin Human-vaginal Microbial mat/biofilm Miscellaneous natural or artificial environment Plant-associated Sediment Soil Wastewater/sludge Water				Yilmaz et al. Nature Biotech. 2011; 29:415-420	

Standards development workflow



Example: M2B3

event



COMMENTARY

Open Access

Marine microbial biodiversity, bioinformatics and biotechnology (M2B3) data reporting and service standards

Petra ten Hoopen¹, Stéphane Pesant², Renzo Kottmann³, Anna Kopf^{3,9}, Mesude Bicak⁴, Simon Claus⁵, Klaas Deneudt⁵, Catherine Borremans⁶, Peter Thijssen⁷, Stefanie Dekeyzer⁵, Dick MA Schaap⁷, Chris Bowler⁸, Frank Oliver Glöckner^{3,9} and Guy Cochrane^{1*}

Abstract

Contextual data collected concurrently with molecular samples are critical to the use of metagenomics in the fields of marine biodiversity, bioinformatics and biotechnology. We present here Marine Microbial Biodiversity, Bioinformatics and Biotechnology (M2B3) standards for "Reporting" and "Serving" data. The M2B3 Reporting Standard (1) describes minimal mandatory and recommended contextual information for a marine microbial sample obtained in the epipelagic zone, (2) includes meaningful information for researchers in the oceanographic, biodiversity and molecular disciplines, and (3) can easily be adopted by any marine laboratory with minimum sampling resources. The M2B3 Service Standard defines a software interface through which these data can be discovered and explored in data repositories. The M2B3 Standards were developed by the European project Micro B3, funded under 7th Framework Programme "Ocean of Tomorrow", and were first used with the Ocean Sampling Day initiative. We believe that these standards have value in broader marine science.

Keywords: Data standard, Marine, Molecular, Biodiversity, Microbial, Bioinformatics, Reporting, Interoperability

Background

An immense wealth of genetic, functional and morphological diversity in marine ecosystems remains unexplored, offering the potential for substantial scientific and biotechnological discoveries. Indeed, significant interest in this area has led to large-scale initiatives, such as Tara Oceans [1], the Global Ocean Survey [2] and Malaspina [3], that target the exploration of marine biodiversity on planetary scales. While the shared goal of such initiatives is the development of an understanding of the compon-

biotechnology will derive benefit. Prerequisite for the successful exploitation of acquired data are standards that enable interoperability in the data infrastructure.

Just as marine studies span many disciplines (e.g. biological, oceanographic, molecular), use of data from marine studies requires approaches that traverse the many disciplines, asking questions, for example, of species distribution, physical oceanographic parameters, molecular biology and data licensing. Each discipline has established infrastructure and best practice for the dissemination,

ten Hoopen et al. *Standards in Genomic Sciences* (2015) 10:20

positories, id analysis when data a lack of a consistent environment for the discovery and retrieval of data.

The Marine Microbial Biodiversity, Bioinformatics, Biotechnology Project (Micro B3) [4] unites intensive oceanographic monitoring, thorough biodiversity studies

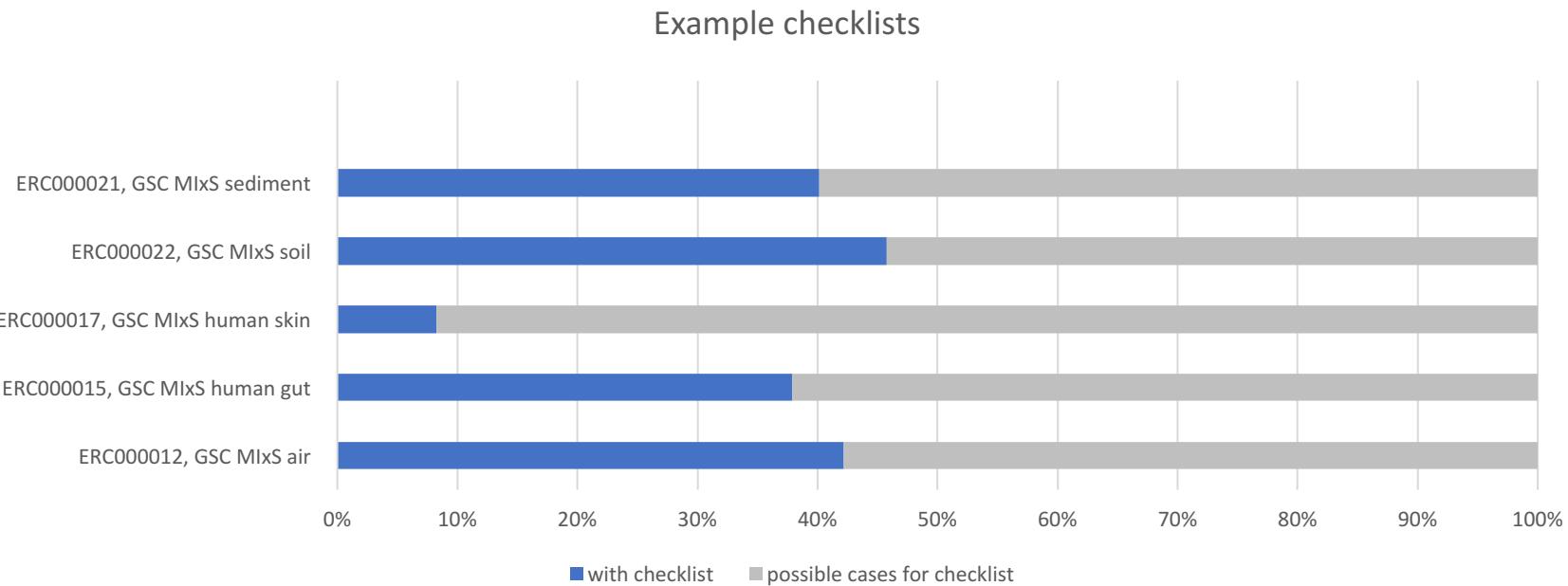
*Correspondence: cochrane@ebi.ac.uk
¹European Nucleotide Archive, EMBL-EBI, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD, UK
Full list of author information is available at the end of the article



© 2015 ten Hoopen et al. *bioRxiv preprint doi: https://doi.org/10.1101/031222; this version posted March 10, 2015. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>).*

Uptake

- **51%** (100,366/195,702) of metagenomics-related samples use appropriate checklists
- **(18%)** (188,213/1,059,594) overall usage of MIxS and MIxS-informed checklists)



- Uptake of individual fields may be greater

*includes pre-publication data

Acknowledgements

- **Content & support:** Ana Cerdeño-Tárraga, Ana Luisa Toribio, Petra ten Hoopen, Marc Rosello, Jeena Rajan, Isabel Santos Magalhaes, **Clara Amid**; Ian Streeter, Susan Fairley, David Richardson, Peter Harrison, **Laura Clarke**
- **Technology (presentation & collaborative tools):** Blaise Alako, Simon Kay, **Nima Pakseresht**; Xin Liu, Suran Jayathilaka, **Nicole Silvester**
- **Technology (submissions & data back-end):** Daniel Vaughan, Neil Goodgame, Iain Cleland, Josué Martinez Villacorta, Dmitriy Smirnov, Kethi Reddy, Vadim Zalunin, Rasko Leinonen, **Thomas Keane**
- **Infrastructure:** EMBL-EBI Technical Services Cluster
- **Gigascience publication co-authors**
- **INSDC Partners:** NCBI, DDBJ
- **Funders**

