# The challenge of NGS data

Guy Cochrane

# Context

# EMBL European Bioinformatics Institute

## Genes, genomes & variation

European Nucleotide Archive
1000 Genomes

Ensembl
Ensembl Genomes

European Genome-phenome Archive
Metagenomics portal

Cross domain resources · Cross domain resources

## Gene, protein & metabolite expression

ArrayExpress

Expression Atlas

Metabolights
PRIDE

## Literature & ontologies

Europe PubMed Central
Gene Ontology
Experimental Factor
Ontology

## Protein sequences, families & motifs

InterPro        Pfam        UniProt

## Molecular structures

Protein Data Bank in Europe
Electron Microscopy Data Bank

## Chemical biology

ChEMBL        ChEBI

## Reactions, interactions & pathways

IntAct        Reactome        MetaboLights

## Systems

BioModels        BioSamples
Enzyme Portal

# European Nucleotide Archive (ENA)



ENA

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

- A **broad platform** for the management, sharing, integration and dissemination of sequence data

- Established in the early 1980s, extended for **new technologies and applications**

- **Globally comprehensive scientific record** and European node of INSDC

- **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**: 1.3 petabase pairs across >1 million taxa, 2,000-5,000 active data providers, global consumer userbase

- Rich submission, discovery and retrieval **software, tools and services**

EMBL-EBI

# Submissions

# Data discovery

*temperature>=10 AND temperature<=25 AND geo_box1(42, 17, 43, 18)*



| Accession | First public | Geographical location | Submitter's sample name | Depth (m) | Environment (Biome) | Temperature (C) | Sampling Site |
|---|---|---|---|---|---|---|---|
| SAMEA2591084 | 2014-07-11 | 42.2038 N 17.715 E | TARA_E500000075 | 5.0 | marine biome (ENVO:00000447) | 17.32198 | TARA_023 |
| SAMEA2591093 | 2014-06-23 | 42.2038 N 17.715 E | TARA_A100000551 | 5.0 | marine biome (ENVO:00000447) | 17.32198 | TARA_023 |
| SAMEA2591094 | 2014-06-23 | 42.2038 N 17.715 E | TARA_A100000553 | 5.0 | marine biome (ENVO:00000447) | 17.32198 | TARA_023 |
| SAMEA2591095 | 2014-06-26 | 42.2038 N 17.715 E | TARA_A100000552 | 5.0 | marine biome (ENVO:00000447) | 17.32198 | TARA_023 |
| SAMEA2591096 | 2014-06-26 | 42.2038 N 17.715 E | TARA_A100000547 | 5.0 | marine biome (ENVO:00000447) | 17.32198 | TARA_023 |
| SAMEA2591097 | 2014-07-18 | 42.2038 N 17.715 E | TARA_E500000056 | 5.0 | marine biome (ENVO:00000447) | 17.32198 | TARA_023 |
| SAMEA2591098 | 2014-07-18 | 42.1735 N 17.7252 E | TARA_E500000081 | 55.0 | marine biome (ENVO:00000447) | 15.194062 | TARA_023 |
| SAMEA2591099 | 2014-07-11 | 42.1735 N 17.7252 E | TARA_E500000080 | 55.0 | marine biome (ENVO:00000447) | 15.194062 | TARA_023 |
| SAMEA2591103 | 2014-06- | 42.1735 N 17.7252 | TARA_A100000559 | 55.0 | marine biome | 15.194062 | TARA_023 |

*tax_tree(10090) AND library_source="GENOMIC" AND instrument_platform="ILLUMINA" AND library_strategy="ChIP-Seq"*



| Study accession | Sample accession | Run accession | Scientific name | Fastq files (ftp) | Fastq files (galaxy) | Submitter's sample name |
|---|---|---|---|---|---|---|
| PRJEB6568 | SAMEA2604495 | ERR537823 | Mus musculus | File 1 File 2 | File 1 File 2 | E-MTAB-2661:Exp2-Irf5-WT-120m-R2 |
| PRJEB6568 | SAMEA2604492 | ERR537824 | Mus musculus | File 1 File 2 | File 1 File 2 | E-MTAB-2661:Exp2-Irf5-KO-000m-R1 |
| PRJEB6568 | SAMEA2604494 | ERR537825 | Mus musculus | File 1 File 2 | File 1 File 2 | E-MTAB-2661:Exp2-Irf5-input-120m |
| PRJEB6568 | SAMEA2604493 | ERR537826 | Mus musculus | File 1 File 2 | File 1 File 2 | E-MTAB-2661:Exp2-Irf5-WT-000m-R1 |
| PRJEB6568 | SAMEA2604496 | ERR537827 | Mus musculus | File 1 File 2 | File 1 File 2 | E-MTAB-2661:Exp2-Irf5-WT-120m-R1 |
| PRJEB6568 | SAMEA2604491 | ERR537828 | Mus musculus | File 1 File 2 | File 1 File 2 | E-MTAB-2661:Exp2-Irf5-KO-120m-R2 |
| PRJEB6568 | SAMEA2604500 | ERR537829 | Mus musculus | File 1 File 2 | File 1 File 2 | E-MTAB-2661:Exp2-Irf5-input-000m |
| PRJEB6568 | SAMEA2604498 | ERR537830 | Mus musculus | File 1 File 2 | File 1 File 2 | E-MTAB-2661:Exp2-Irf5-KO-120m-R1 |
| PRJEB6568 | SAMEA2604497 | ERR537831 | Mus musculus | File 1 | File 1 | E-MTAB-2661:Exp2-Irf5-WT-000m- |

# Cross-references and tagging
GUI: http://www.ebi.ac.uk/ena/data/xref/search

# A global platform for the sequence-based rapid identification of pathogens



COMPARE: the enabling system for rapid identification, containment and mitigation of emerging infectious diseases and foodborne outbreaks by generation and comparison of genomic information on samples and pathogens across sectors, time and locations, with additional contextual data.

# Ocean Samping Day and Tara Oceans

# Technology

# Data accumulation



Cook CE, Bergman MT, Finn RD, Cochrane G, Birney E, Apweiler R. The European Bioinformatics Institute in 2016: Data growth and integration. Nucleic Acids Res. 2016 Jan;44(D1) D20-6. doi:10.1093/nar/gkv1352. PMID: 26673705; PMCID: PMC4702932.

EMBL-EBI

# Starting point

TGAGCTCTAAGTACC
329183050298757

# Sequence compression



- Encoding of read starts and differences

- 3.5x–100x compression over existing formats

- Scales favourably with increasing read length and density

Fritz, M.H. Leinonen, R., et al. (2011) Efficient storage of high throughput DNA sequencing data using reference-based compression. *Genome Res.* 21 (5), 734-40.

EMBL-EBI

# Quality compression

TGAGCTCTAAGTACC
329183050298757

Horizontal reduction

Vertical reduction

002020010022212

-2---30---9---7

EMBL-EBI

# Quality compression: simple, horizontal reduction



Photograph from MichaelMaggs, http://en.wikipedia.org/wiki/File:Amanita_muscaria_(fly_agaric).JPG

# Models for data reduction



Jong-Seok Lee et al. (2009), http://mmspg.epfl.ch/files/content/sites/mmspl/files/shared/lee_icme.pdf

# CRAM performance



100                            10                          1                       0.1

Bits/base

EMBL-EBI

# CRAM performance

# Human aspects

# Standards

event
sampling — sample
measurement — environment
organism

data/time, longitude, latitude
site, platform, campaign
depth, sample title, protocol label
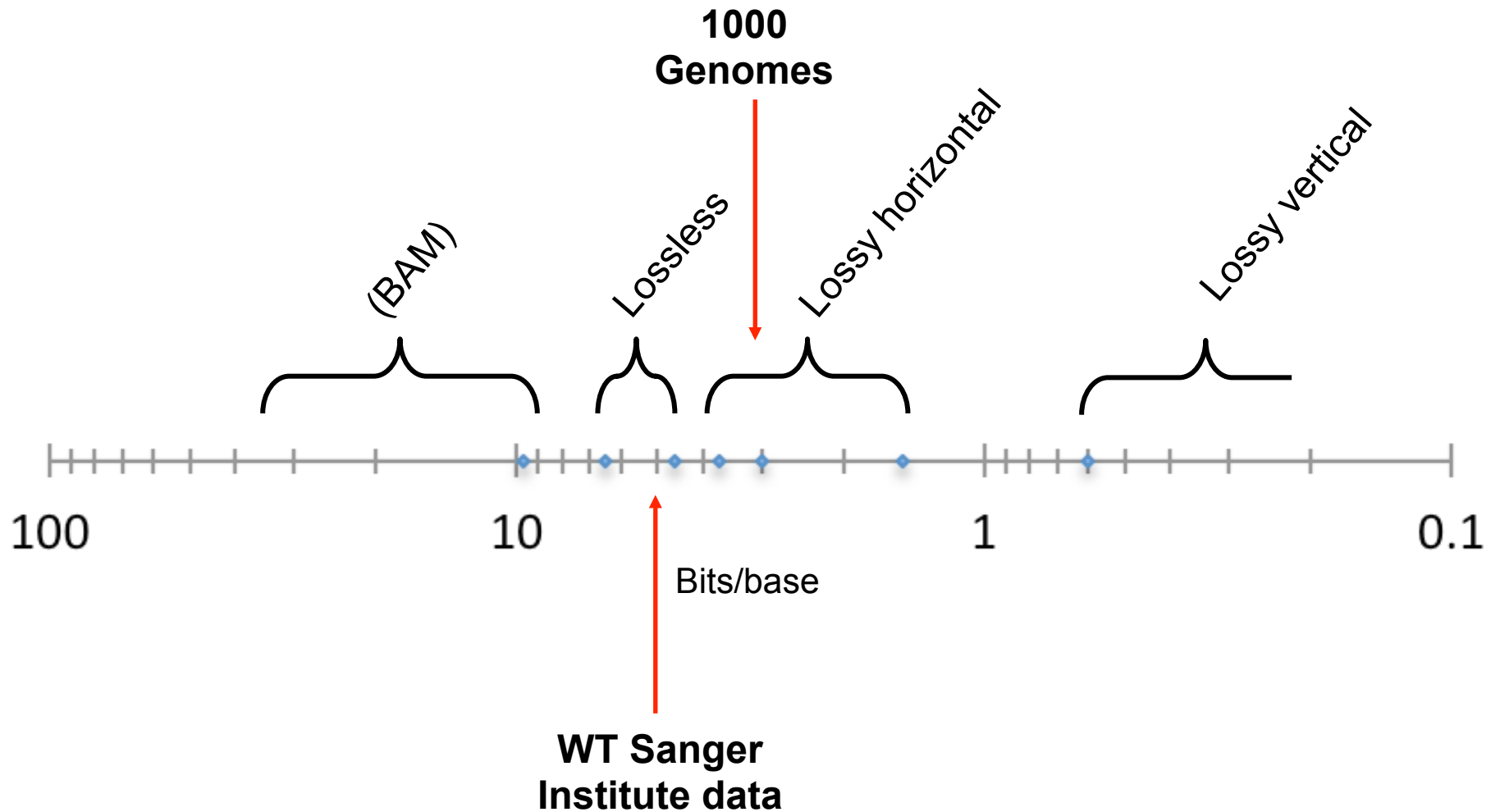parameter ID
temperature, salinity, material, feature, biome
taxon ID
scientific name

event
sampling — sample
measurement — environment
organism

event ID, device, method, comment
authors, project, objective
content, container, quantity, size fraction >, size fraction <, treatment storage, treatment chemical
parameter name, dimensions, quantity, currency, comment, method, units
environ. parameters, marine region
sex, size, biomass, biovolume, life stage

genomic
oceanographic
biodiversity
optional
recommended
mandatory
M2B3 checklist

**Biodiversity Bioinformatics Biotechnology**

**COMMENTARY**    **Open Access**

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

Petra ten Hoopen[1], Stéphane Pesant[2], Renzo Kottmann[3], Anna Kopf[3,9], Mesude Bicak[4], Simon Claus[5], Klaas Deneudt[5], Catherine Borremans[6], Peter Thijsse[7], Stefanie Dekeyzer[5], Dick MA Schaap[7], Chris Bowler[8], 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 compos...

...future marine survey projects will add value to these explorations and will continue to build a powerful marine data infrastructure from which ecosystems biology and

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 dissemin...

...a lack of interoperability between standards and the 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
[1]European 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

BioMed Central

EMBL-EBI

# Metagenomics

Workflow diagram:

Raw reads → Sequence file preparation → Initial reads → QC → Processed reads → RNA identification & masking

RNA identification & masking → rRNAs → 16S → Taxonomic analysis → Taxonomic assignments

RNA identification & masking → Reads with rRNA masked → ORF predictions → Predicted CDS → Functional analysis → Functional assignments

Legend:
- Download EBI metagenomics (box)
- Download ENA (green box)
- Process/component (circle)
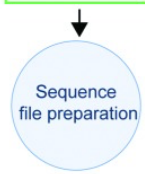- Multiple files (stacked box)
- Not available for download (dashed box)

# EBI metagenomics in 2016 - an expanding and evolving resource for the analysis and archiving of metagenomic data

Alex Mitchell, Francois Bucchini, Guy Cochrane, Hubert Denise, Petra ten Hoopen, Matthew Fraser, Sebastien Pesseat, Simon Potter, Maxim Scheremetjew, Peter Sterk, and Robert D. Finn*

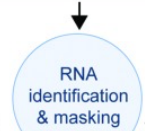Jump to: Biological process | Molecular function | Cellular component

Most frequent GO terms (molecular function)

GO terms list
Click to hide
- aminoacyl-trna ligase activity
- atpase activity
- cofactor binding
- dna polymerase activity
- gtpase activity
- hydrolase activity
- hydrolase activity, acting on ester bonds
- hydrolase activity, acting on glycosyl bonds
- isomerase activity
- ligase activity
- lyase activity
- metal ion binding
- nucleic acid binding
- nucleotide binding
- oxidoreductase activity
- peptidase activity
- rna polymerase activity
- transferase activity
- transporter activity
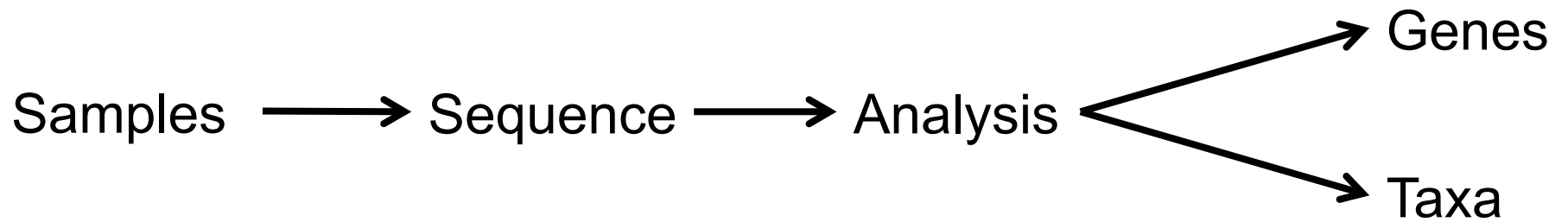- other (less than 2%)

BMC
Bioinformatics

# RIEMS: a software pipeline for sensitive and comprehensive taxonomic classification of reads from metagenomics datasets

Matthias Scheuch[†], Dirk Höper[*†] and Martin Beer

Scheuch *et al*. (2015) BMC Bioinformatics 16:69
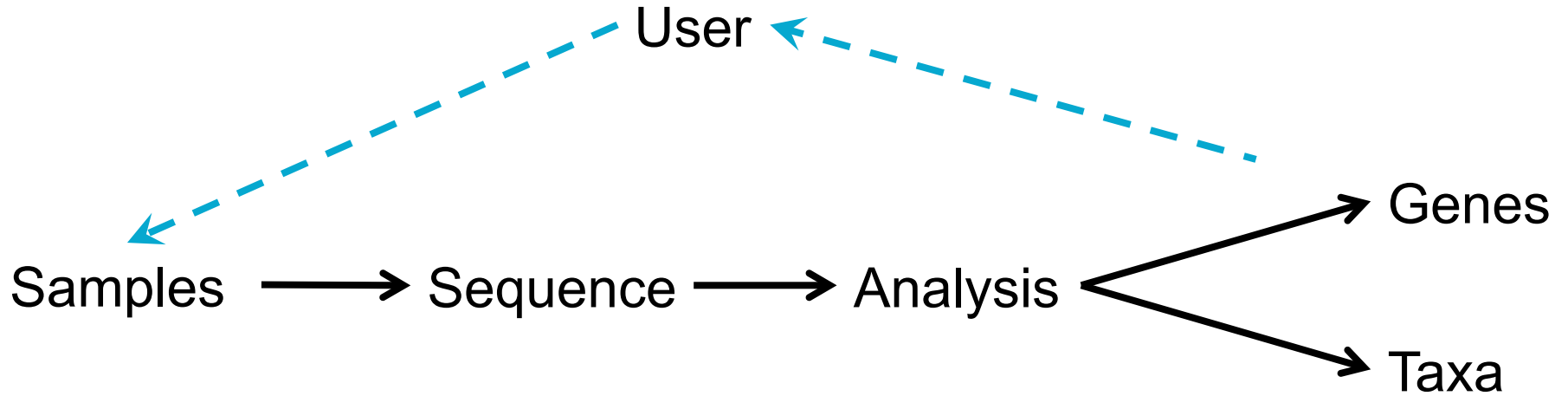
EMBL-EBI

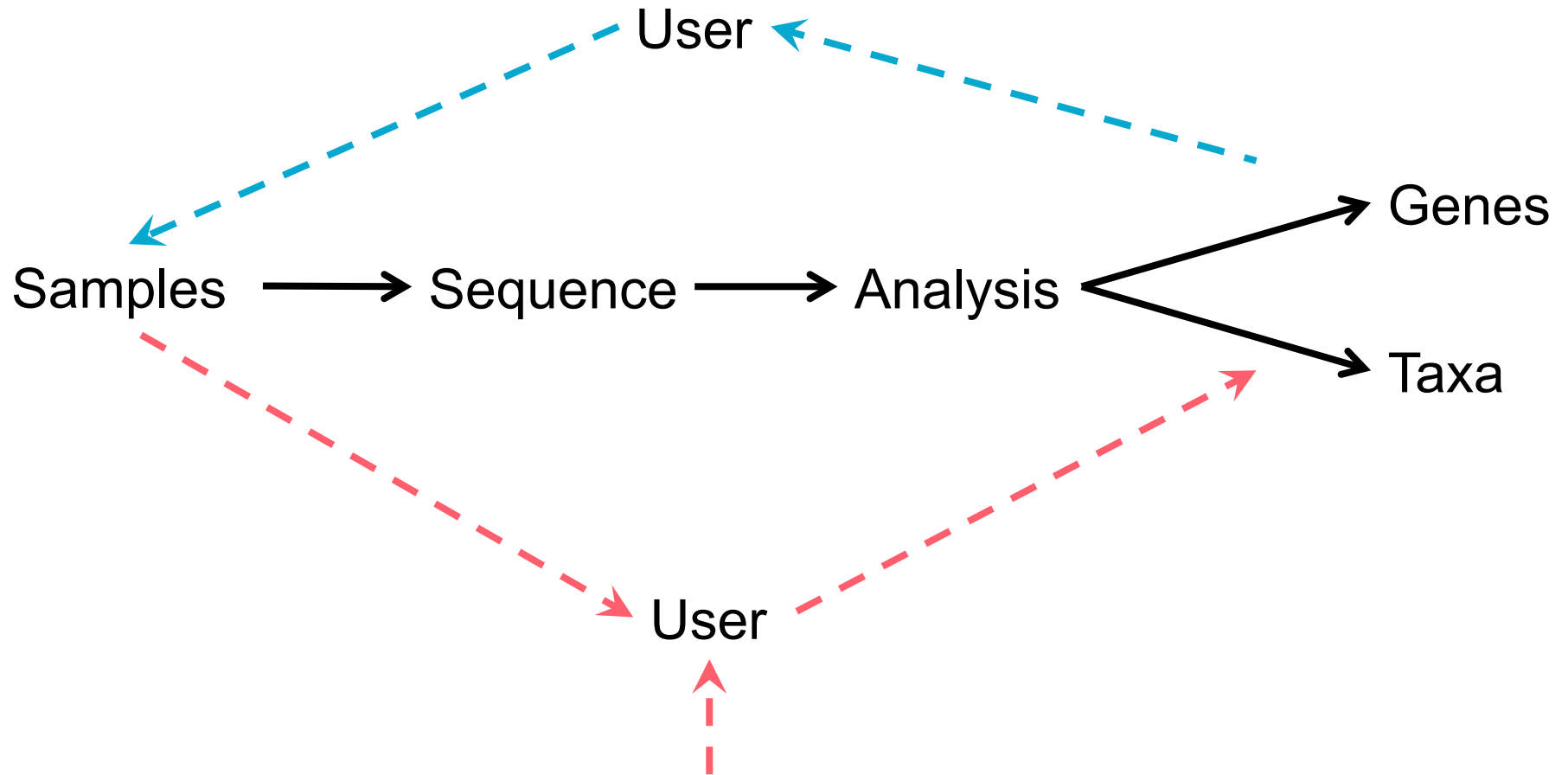# Identification data

# Identification data accumulation

- Identification against context is informative

- (Molecular observations may be tentative or high confidence)

- Coincidental observations (time, place, virulence, host phenotypes)

- 'Capture' identifications to 'connect' coincidental observations

EMBL-EBI

# Identification data

# Identification data

# Acknowledgements