

FAANG metadata - sample specification

This document describes the specification for all sample metadata. You can find an overview of our metadata and archival plans in the overview document. The experiment and analysis documents are also in this git repo.

In the sample context, we consider donor animals, tissue samples, primary cells or other biological material to be samples. All Samples must be registered in BioSamples at EMBL-EBI as this samples archive has the best support for related to and derived from sample relationships. The NCBI BioSample database is a peer of the EMBL-EBI BioSamples, and they exchange data regularly. FAANG samples should be registered in the EMBL-EBI BioSamples prior to data submission. This document describes the attributes which must be associated with any BioSamples submission.

Sample metadata requirements

Most requirements are laid out like this:

- **attribute name** (*data type*) a brief description

The data types are described later in this document.

Common

These attributes should be present on every sample record.

Required:

- **Sample name** (*text*) sample names should follow the naming rules listed below. Each name must be unique.
- **Material** (*ontology term*) the type of material being described. This will be used to decide what metadata are required and must be one of the expected terms:
 - organism
 - tissue specimen
 - cell specimen
 - cell culture
 - pool of specimens
- **project** (*text*) project name - this should always be 'FAANG'. This will allow the DCC to identify FAANG samples

Optional:

- **description** (*text*) a brief description of the sample
- **availability** (*URL*) either a link to a web page giving information on sample availability (who to contact and if the sample is available), or a e-mail address to contact about availability. E-mail addresses should be prefixed with 'mailto:', e.g. 'mailto:samples@example.ac.uk'. In either case, long term support of the web page or e-mail address is necessary. Group e-mail addressees are preferable to individual.

Animal

An animal sampled for FAANG. The following attributes are in addition to the attributes listed in the 'Common' section above. The **material** should be reported as organism.

Required:

- **Organism** (*NCBI taxon ID*)
- **sex** (*ontology term*) animal sex, described using any child term of PATO_0000047
- **birth date** (*date*) birth date, in the format YYYY-MM-DD, or YYYY-MM where only the month is known
- **breed** (*ontology term*) animal breed, described using a term from the Livestock Breed Ontology
- **health status** (*ontology term*) Healthy animals should have the term normal, otherwise use the as many disease terms as necessary from EFO

Optional:

- **birth location** (*text*) name of the birth location
- **birth location latitude** (*number*) latitude of the birth location in decimal degrees. Units should be specified as 'decimal degrees'
- **birth location longitude** (*number*) longitude of the birth location in decimal degrees. Units should be specified as 'decimal degrees'
- **birth weight** (*number*) weight, in kilograms or grams. Units must be specified
- **placental weight** (*number*) weight, in kilograms or grams. Units must be specified.
- **pregnancy length** (*number*) length of time, in days, weeks or months
- **delivery timing** (*ontology term*) possible values
 - early parturition
 - full-term parturition
 - delayed parturition
- **delivery ease** (*text*) possible values
 - normal autonomous delivery

- c-section
- veterinarian assisted
- **physiological conditions**(*ontology term*) use as many terms as necessary from ATOL)
- **environmental conditions**(*ontology term*) as many terms as necessary from EOL)
- **phenotype** (*ontology term*) as many terms as required from the VT, ATOL or MP) ontologies
- **pedigree** (*URL*) a link to pedigree information for the animal

Links to other records:

- **child of** (*sample*) sample name or Biosample ID for sire/dam. Required if related animals are part of FAANG, e.g. quads.

Specimen

A piece of tissue taken from an animal. The following attributes are in addition to the attributes listed in the ‘Common’ section above. The **material** should be reported as tissue specimen.

Required:

- **specimen collection date**(*date*) date at which the specimen was collected
- **animal age at collection** (*number*) animal age at the point of collection, in years, months, weeks or days. Units must be specified. An estimate is acceptable where the age is not precisely known.
- **developmental stage** (*ontology term*) a child term of life cycle stage
- Animal Disease / health status at point of collection
- **tissue** (UBERON term preferred)
- **specimen collection protocol** (*protocol*) a link to the protocol followed when taking the specimen
- **fasted status** - (*text*) One of the following values, for which the criteria *must* be specified in the protocol:
 - fed
 - fasted
 - unknown
- **health status at collection** (*ontology term*) Healthy animals should have the term normal, otherwise use the as many disease terms as necessary from EFO

Optional:

- **physiological_conditions**(*ontology term*) as many terms as necessary from ATOL
- **number_of_pieces** (*number*) Units must be specified as ‘count’
- **specimen volume** (*number*) Units must be specified as either ‘square centimeters’, ‘liters’ or ‘milliliters’
- **specimen size**(*number*) Units must be specified as either ‘meters’, ‘centimeters’ or ‘millimeters’
- **specimen weight** (*number*) Units must be specified as either ‘grams’, ‘kilograms’
- **specimen picture url** (*URL*) Link to a picture of the specimen
- **gestational age at sample collection** (*number*) If the animal was pregnant when the specimen was taken, state how long had it been pregnant for. Units must be specified as ‘days’ or ‘weeks’.

Links to other records:

- **derived from** (*sample*) sample name or BioSample ID for an *animal* record (required).

Purified cells

Cells purified from a specimen. The following attributes are in addition to the attributes listed in the ‘Common’ section above. The **material** should be reported as cell specimen.

Required:

- **markers** (*text*) markers used to isolate and identify the cell type
- **cell type** (*ontology term*) a term from the CL ontology
- **purification protocol** (*protocol*) protocol describing how the cells were purified

Links to other records:

- **derived from** (*sample*) sample name or BioSample ID for a *specimen* record (required).

Cell culture

Cells cultured from a specimen or purified cells. The following attributes are in addition to the attributes listed in the ‘Common’ section above.

Required:

- **culture type**(*ontology term*) a child term of BTO_0000214
- **cell type** (*ontology term*) a term from the CL ontology
- **cell culture protocol** (*protocol*) protocol describing how the cells were purified
- **culture conditions** (*text*) brief description of culture conditions (e.g. 'on feeder cells', 'E8 media')
- **number of passages** (*number*) number of times the cell line has been re-plated and allowed to grow back to confluency or to some maximum density if using suspension cultures

Links to other records:

- **derived from** (*sample*) sample name or BioSample ID for a *specimen* or *purified cell* record (required).

Pooled samples

Where samples are pooled, a new sample record should be created, containing

- **pooling protocol** (*protocol*)

Links to other records:

- **derived from** (*sample*) sample name or BioSample ID for a *specimen*, *purified cells* or *cell culture* record (required).

Data types for sample attributes

BioSamples takes sample records with a set of attributes. Each attribute has a name and a value. It can also have 'Units', or a 'Term Source' and a 'Term Source ID'. The Term Source and ID allow us to refer to entries in other databases or ontologies. This is fully described on the BioSamples help pages. The following section describe the expectations for each data type within FAANG.

date

Dates should be reported in an ISO 8601 format, YYYY-MM-DD for dates or YYYY-MM for months. To ensure clarity, the format must be reported as the 'units'.

NCBI taxon ID

A species name and identifier from the NCBI Taxonomy database. For example, a human would be described with a value of 'Homo sapiens', a term source of 'NCBI Taxonomy' and a term source ID of 9606.

number

A number, with units specified. BioSamples recommends that units are given without abbreviations. For example, a birth weight could have a value of 1.3 and the units specified as 'kilograms'.

protocol

A URL link to a protocol document on the FAANG FTP site. Please contact the FAANG data coordination centre to have your protocol documents added to the FTP site.

text

Text, using US English spellings.

URL

A URL, such as 'http://faang.org/'. Depending on the context, http, ftp, mailto links may be appropriate. Examples:

- ftp, ftp://ftp.faaang.ebi.ac.uk/ftp/README
- http, http://faang.org/
- mailto, mailto:bob@example.org

ontology term

A reference to an ontology term. The attribute value should be the term label. The term source should be the ontology used, and the term source ID should be an ID from that ontology. For example, cerebral cortex could be described with a term source of 'UBERON', a term source ID of 'UBERON:0000956' and a value of 'cerebral cortex'.

location

A location should be reported as using three attributes:

- **location** (*text*) name of the location
- **location latitude** (*number*) latitude in decimal degrees. Units should be reported as ‘decimal degrees’
- **location longitude** (*number*) longitude in decimal degrees. Units should be reported as ‘decimal degrees’

sample

Samples can be referred to in two ways. If the sample you need to reference is in the submission, use the sample name. If the sample was already submitted, use the BioSample ID (e.g. SAMEA2821491).

Missing data

Where data cannot be included in a submission, submit one of these text values instead

- ‘not applicable’
- ‘not collected’ (i.e. will always be missing)
- ‘not provided’ (i.e. may be added later)
- ‘restricted access’ (i.e. it isn’t missing, we just can’t include it in a public document)

The use of these values will interact with the metadata validation system as follows:

- attribute is required
 - not applicable, not collected, not provided - validation will regard these as an error
 - restricted access - validation will generate a warning
- attribute is recommended
 - not applicable, not collected, not provided - validation will generate a warning
 - restricted access - pass
- attribute is optional
 - validation will pass with any of missing values terms

Sample naming

We propose a sample naming scheme comprising the following elements:

- short species code
- lab or institute short name
- alpha numeric sample ID from LIMS

The purpose is to ensure that samples are uniquely and clearly identified, with reasonably short names.

Short species codes:

- *Bos taurus* BTA
- *Sus scrofa* SSC
- *Ovis aries* OAR
- *Gallus gallus* GGA
- *Equus caballus* ECA
- *Capra hircus* CHR

Submission

Samples should be submitted to BioSamples@EBI. All samples tagged with a project of 'FAANG' will be added to the FAANG BioSamples group. Samples in this group will be synced to BioSample@NCBI periodically. Samples in BioSamples@EBI/BioSample@NCBI can be referenced in submissions to SRA at EBI and NCBI.

Validation

The DCC team at EBI will check the submitted metadata against the specification. Samples that do not meet the minimum requirements will be not be included in FAANG data releases.