

# Lawrence Berkeley National Laboratory

## LBL Publications

### Title

The Genomes OnLine Database (GOLD) v.5: a metadata management system based on a four level (meta)genome project classifications

### Permalink

<https://escholarship.org/uc/item/2gk413jc>

### Authors

Thomas, Alex D.  
Stamatis, Dimitri  
Bertsch, Jon  
et al.

### Publication Date

2014-10-29

# The Genomes OnLine Database (GOLD) v.5: a metadata management system based on a four level (meta)genome project classifications

**AUTHORS:** T. B. K. Reddy<sup>1\*</sup>, Alex D. Thomas<sup>1</sup>, Dimitri Stamatis<sup>1</sup>, Jon Bertsch<sup>1</sup>, Michelle Isbandi<sup>1</sup>, Jakob Jansson<sup>1</sup>, Jyothi Mallajosyula<sup>1</sup>, Ioanna Pagani<sup>1</sup>, Elizabeth A. Lobos<sup>1</sup> and Nikos C. Kyrpides<sup>1,2\*</sup>

1. DOE Joint Genome Institute, 2800 Mitchell Drive, Walnut Creek, CA 94598, USA
2. Department of Biological Sciences, Faculty of Science, King Abdulaziz University, Jeddah, Saudi Arabia
- 3.

\* To whom correspondence may be addressed. [tbreddy@lbl.gov](mailto:tbreddy@lbl.gov), and/or [nckyrpides@lbl.gov](mailto:nckyrpides@lbl.gov)

October 29, 2014

## ACKNOWLEDGMENTS:

We are grateful to all the colleagues who kindly provide information for the more accurate monitoring of the genome projects and their associated metadata. We thank the members of the microbial genomics and metagenomics programs at the JGI for support, useful discussions and exchange of ideas. We would also like to thank Lynn Schriml and Lynette Hirschman for feedback and useful discussion on the EnvO ontology. This work was performed under the auspices of the US Department of Energy Office of Science, Biological and Environmental Research Program, and by the University of California, Lawrence Berkeley National Laboratory under contract no. DE-AC02-05CH11231.

## DISCLAIMER:

This document was prepared as an account of work sponsored by the United States Government. While this document is believed to contain correct information, neither the United States Government nor any agency thereof, nor The Regents of the University of California, nor any of their employees, makes any warranty, express or implied, or assumes any legal responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by its trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof, or The Regents of the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof or The Regents of the University of California.

**Some supplementary files may need to be viewed online via your Referee Centre at <http://mc.manuscriptcentral.com/nar>.**

# The Genomes OnLine Database (GOLD) v.5: a metadata management system based on a four level (meta)genome project classification

T. B. K. Reddy<sup>1\*</sup>, Alex D. Thomas<sup>1</sup>, Dimitri Stamatis<sup>1</sup>, Jon Bertsch<sup>1</sup>, Michelle Isbandi<sup>1</sup>, Jakob Jansson<sup>1</sup>, Jyothi Mallajosyula<sup>1</sup>, Ioanna Pagani<sup>1</sup>, Elizabeth A. Lobos<sup>1</sup> and Nikos C. Kyrpides<sup>1,2\*</sup>

<sup>1</sup> Prokaryotic Super Program, DOE Joint Genome Institute, Walnut Creek, CA 94598 USA

<sup>2</sup> Department of Biological Sciences, Faculty of Science, King Abdulaziz University, Jeddah, Saudi Arabia

\* To whom correspondence should be addressed. Tel: 1-925-927-2580; Fax: 1-925-296-5666; Email: [tbreddy@lbl.gov](mailto:tbreddy@lbl.gov), [nckyrpides@lbl.gov](mailto:nckyrpides@lbl.gov)

## ABSTRACT

The Genomes OnLine Database (GOLD, <http://www.genomesonline.org>) is a comprehensive online resource to catalogue and monitor genetic studies worldwide. GOLD provides up-to-date status on complete and ongoing sequencing projects along with a broad array of curated metadata. Here we report version 5 (v.5) of the database. The newly designed database schema and web user interface supports several new features including the implementation of a four level (meta)genome project classification system and a simplified intuitive web interface to access reports and launch search tools. The database currently hosts information for about 19,200 studies, 56,000 Biosamples, 56,000 sequencing projects, and 39,400 analysis projects. More than just a catalogue of worldwide genome projects, GOLD is a manually curated, quality controlled metadata warehouse. The problems encountered in integrating disparate and varying quality data into GOLD are briefly highlighted. GOLD fully supports and follows the Genomic Standards Consortium (GSC) Minimum Information standards.

## INTRODUCTION

The Genomes OnLine Database (GOLD) is a data management system for cataloguing and continuous monitoring of sequencing projects worldwide. GOLD collects, curates, and disseminates metadata associated with those projects. GOLD is currently in its fifth version (1–6). With rapidly decreasing costs for sequencing, the number of sequencing projects and the amount of sequence data generated is increasing at an exponential rate. As these data are submitted to various public resources like GenBank (7) and EMBL (8) or analysis platforms like Integrated Microbial Genomes (IMG) (9) and MG-RAST (10), it becomes increasingly important to document the associated metadata in order to facilitate comparative analysis and hypothesis generation. The Genomics Standards Consortium (GSC) mandates the Minimum Information about any (x) Sequence (MIxS) specifications to be used when making sequence data available in public repositories (11, 12). GOLD is fully compliant with the GSC's MIxS standards in capturing metadata and provides a platform to query projects based on various metadata features.

1  
2  
3 GOLD supports the IMG family of data management systems (9, 13–15) as a gatekeeper  
4 of projects and metadata, and requires that projects are annotated with at least minimal  
5 metadata. In fact an entry in GOLD and compliance with required metadata is a prerequisite to  
6 submit a project to the IMG systems for annotation. The main steps in the process include  
7 project registration in GOLD, project submission to IMG for annotation, and finally publication of  
8 results in the GSC's journal, Standards in Genomic Sciences (SIGS)  
9 (<http://www.standardsingenomics.com/>), or other journals of your choice. Since GOLD complies  
10 with MIxS, all available required metadata is already in place to publish in SIGS.  
11  
12

13  
14  
15 In the past, when sequencing was still expensive and only a limited number of high-  
16 interest organism genomes were sequenced, maintaining the associated information in a  
17 catalogue format was sufficient. With lower sequencing costs, many more genomes are now  
18 being sequenced as part of a single study. Initiatives such as the Human Microbiome Project  
19 (HMP) (16) and Genomic Encyclopedia of Bacteria and Archaea (GEBA) (17, 18) are a couple of  
20 examples where several thousands of genomes were sequenced as part of a single initiative. The  
21 emergence of high-throughput sequencing technologies and the development of analysis tools for  
22 studying metagenomes has facilitated the rapid growth in metagenome studies as well. It is also  
23 becoming more common to use multiple sequencing approaches on the same sample(s), for  
24 example the Functional Encyclopedia of Bacteria and Archaea (FEBA) (19). In such cases it is  
25 important to collect common metadata pertaining to these samples and organize all of the  
26 samples under one or more relevant studies.  
27  
28

29  
30  
31 The increasing variety of sequencing and analysis projects need to be linked and tracked  
32 in a seamlessly integrated system. One of the major limitations of the previous versions of the  
33 database has been the assumption of a one-to-one relationship between related components.  
34 For example, the previous versions could not correlate multiple sequencing projects to a single  
35 sample. In the event an isolate genome and metagenome were derived from a single sample a  
36 separate record for each sequence would need to be created. Similarly the previous versions  
37 could not capture the multiple sequencing projects of a combined assembly nor was it possible to  
38 connect multiple analyses to a single sequence project. Another limitation was that all genome  
39 projects were designated as isolates, an incorrect assignment for a genome assembled from a  
40 metagenome. These issues necessitated a new mechanism to organize various components of  
41 sequencing studies.  
42  
43

#### 44 45 46 47 48 49 50 **NEW TO THIS RELEASE**

51  
52 Version 5 of the database is founded on a fundamentally redesigned schema to accommodate a  
53 four level project classification system (Fig. 1). The new classification system is comprised of  
54 Studies, Biosamples, Sequencing Projects (SPs), and Analysis Projects (APs). Studies constitute  
55 the highest level of classification in the system, containing Biosamples, SPs, and APs that are  
56  
57  
58  
59  
60

1  
2  
3 part of a single initiative. GOLD's Biosamples represent the physical isolate or environmental  
4 material from which genetic material is extracted for sequencing. GOLD's Biosamples have no  
5 relation to NCBI BioSamples. GOLD's SPs represent sequencing protocols such as whole  
6 genome sequencing, transcriptomes, metagenomes, metatranscriptomes, methylation  
7 sequencing etc. applied to Biosamples. APs are the analytical processes applied to the SPs.  
8 Multiple different assemblies or annotations of the same SPs, would result to multiple different  
9 APs with varying metadata that need to be captured. These four components are described in  
10 more detail below.  
11  
12  
13  
14

15 In addition to the four levels described above, one more entity has been introduced in the  
16 new schema to provide metadata information for the individual organisms. In the previous  
17 versions of the database, each sequencing project of an isolate organism, included both the  
18 metadata for the sequencing information and the organism in a single record. Increasingly, the  
19 genome of a single organism is being sequenced more than once, by different groups, making it  
20 inefficient to associate the same organism metadata individually with every different project.  
21 GOLD v.5 defines and curates the organism records with core taxonomy, environmental, and  
22 other metadata independently of their associated SPs. As a result, this entity can be used by all  
23 SPs without the need for curating and propagating redundant metadata. By doing so, v.5 now  
24 enables the identification of all the organisms with different but synonymous names.  
25  
26  
27  
28  
29

30 Historically, the focus of the database was to provide a comprehensive coverage to all  
31 prokaryotic genomes and metagenomes. We are in the process of systematically integrating  
32 eukaryotic SPs into GOLD. Projects are introduced in the database from three main streams: (1)  
33 projects sequenced at the JGI are automatically added following a number of Quality Control  
34 (QC) checks; (2) projects submitted to the database from individual researchers around the world;  
35 and (3) projects available at the NCBI's BioProject portal.  
36  
37  
38

39 The previous versions of the database provided read-only project reporting system. This  
40 served user needs for accessing project information and searching for projects based on specific  
41 metadata. However, the user interface for project creation and curation was provided through a  
42 separate system called IMG-GOLD. The new version has enabled the seamless integration of  
43 these two formerly separated functions into a single resource.  
44  
45

46 Isolate genomes via their associated Biosamples are now classified using the same five-  
47 tier hierarchical classification system previously developed and implemented for metagenomes  
48 (20). Over 10,000 public isolate genomes have been classified accordingly. Over 9,000 isolate  
49 genomes have also been curated to add strain habitat classifications. This field refers to the  
50 specific habitat of the strain according to the strain isolation information, as opposed to the  
51 previous general habitat in the database which corresponds to the species. The controlled  
52 vocabulary of the strain habitat has been mapped to the hierarchical ecosystem classification. For  
53 example, there are 161 genomes for organisms with the classification path Host-associated  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 (ecosystem) -> mammals (ecosystem category)-> digestive system (ecosystem type) -> foregut  
4 (ecosystem subtype)-> rumen (specific ecosystem). The strain habitats within this group include  
5 'sheep rumen', 'cattle rumen' and 'goat rumen' with 37, 53, and 1 genome respectively. Thus,  
6 there is manual curation of organism Biosamples with specific habitat terms.  
7  
8

9 A newly designed web interface provides access to data through various pre-selected  
10 reports, project distribution graphs, statistics and an intuitive search interface that allows a user to  
11 search based on an array of metadata fields. The new implementation also provides access for  
12 public users to search for APs submitted to the IMG systems.  
13  
14

## 15 **GOLD DATABASE ORGANIZATION AND DATA OVERVIEW**

### 16 **The Four-level classification system:**

17  
18  
19  
20 The current release organizes genome, metagenome, and other sequencing projects into a  
21 system of four levels which are described below.  
22  
23

### 24 **GOLD Study:**

25  
26  
27 A Study represents the highest-level organization. Studies include one or more Biosamples and  
28 their associated SPs and Aps that have been grouped to investigate a related research topic of  
29 interest. For example, the HMP (16), GEBA (17, 18) and KMG (21) studies represent typical  
30 cases where researchers set out to explore a specific topic by sequencing thousands of samples.  
31 Studies like GEBA-MDM (22) and FEBA (19) applied several different sequencing strategies  
32 (e.g., isolate genomes, single-cell genomes, metagenomes and transcriptomes etc.) as part of a  
33 single study. Studies may be composed of one to hundreds of Biosamples from a wide range of  
34 ecological settings (Fig. 2). Each Biosample may also yield several different SPs, each of which  
35 may yield multiple APs (Fig. 3, Table 1). Study IDs are referred to as "Gs" IDs in the new system.  
36 A GOLD Study is analogous to the NCBI's umbrella BioProject, and may contain one or more  
37 NCBI BioSamples.  
38  
39  
40  
41  
42  
43

### 44 **GOLD Biosample:**

45  
46  
47 Biosamples provides a description of the individual environmental sample, from which the  
48 organism, or genetic material (DNA or RNA) was isolated for downstream SPs. There are two  
49 types of Biosamples, organisms and biomes (environmental samples). Historically, samples were  
50 either isolated organisms for whole genome sequencing or environmental samples for  
51 metagenomics. However, it is becoming increasingly common to apply multiple sequencing  
52 techniques to a single sample and thus initiating several different SPs from the same starting  
53 material. For example, from a single biosample, whole DNA can be extracted for a metagenome  
54 and a metatranscriptome SP, as well as cells for single-cell genome project (Fig. 2) (19). The  
55  
56  
57  
58  
59  
60

1  
2  
3 need to manage and organize this type of complexity has led to the creation of GOLD  
4 Biosamples, which are quite distinct from NCBI's Biosamples. While GOLD Biosamples are  
5 organized above the sequencing projects in order to provide linkage of multiple sequencing  
6 projects originating from the same physical sample, NCBI's Biosamples are associated with  
7 individual sequencing projects, providing metadata only for that sequencing project. NCBI's  
8 Biosamples are also used *in lieu* of BioProjects to represent individual sequencing projects as in  
9 the case of multi-isolate projects. GOLD Biosample IDs are represented as "Gb" IDs.

#### 14 **GOLD Sequencing Project (SP):**

16 A number of technological advances have enabled an increasing diversity of SP types (Fig. 3,  
17 Table 2). SPs represent individual sequencing deliverables such as metagenomes,  
18 metatranscriptomes, 16S sequences, single-cell genome sequences, isolate transcriptomes, or  
19 isolate whole genome sequences. As mentioned above, material from one Biosample can be the  
20 basis for more than one SP. GOLD SP's are often connected to a single NCBI BioProject, which  
21 could lead to the misconception that there is a one to one analogy between them. NCBI's  
22 BioProjects represent a mixture of project types that include the umbrella or multi-isolate types  
23 that are more analogous to the GOLD's Studies. This lack of standardization in NCBI BioProjects  
24 is one of the data management challenges the new GOLD classifications aims to address. GOLD  
25 Sequencing Project IDs are represented as "Gp" IDs. Each sequencing project can contain one or  
26 more APs.

#### 34 **GOLD Analysis Project (AP):**

36 APs represent individual data processing methodologies or approaches that are undertaken for a  
37 given SP. As the diversity of data processing and analysis (eg assembly, structural, and  
38 functional annotation) methods has increased, so has the diversity of APs (Fig. 3). More  
39 specifically, the data generated from a single SP may be processed through multiple different  
40 approaches, as researchers exploring various different assembly methods or the same assembly  
41 with different annotation parameters. As shown in Figure 1, a researcher may also generate a  
42 combined assembly from multiple SPs and submit the data for annotation as one AP. This is  
43 more common in the case of single-cell genome projects where sparse sequence data from two  
44 related single cells can result in a better assembly and thereby more coverage of the genome of  
45 the organism being studied. One of the major limitation of the previous systems was the inability  
46 to represent these complex APs with their parent SPs. The current release fills this unmet need in  
47 representing different APs. AP IDs are represented as "Ga" IDs, and there are currently six  
48 different types:

55 (a) **Default AP:** this represents the standard assembly and annotation process applied for any  
56 sequencing project.  
57  
58  
59  
60



1  
2  
3 (b) **Default-screened AP and default - unscreened AP**: these are applicable only for single cell  
4 genome projects where contamination is a major issue due to extraneous DNA or due to errors  
5 during cell sorting/isolation events. Accordingly, there is a need to distinguish between APs that  
6 have gone through a decontamination round (screened) versus those that have not (unscreened).

7  
8  
9 (c) **Combined assembly AP**: these APs use data from multiple SPs which are combined into a  
10 single assembly, which is then submitted for annotation. For example whole genome shotgun  
11 sequencing may be applied to a set of single cell genomes from the same Biosample and the  
12 data from each single cell genome can be used to generate a combined assembly for a better  
13 genome reconstruction. Alternatively, metagenomic sequences from multiple different Biosamples  
14 may be combined into a single assembly. Tracking these many relationships between  
15 Biosamples, SPs, and APs within a Study is a key feature of new GOLD.

16  
17 (d) **Genome from metagenome AP**: these APs represent individual genomes extracted from  
18 metagenomics data. Advances in metagenomic assembly and binning (  
19 <http://ggkbase.berkeley.edu/>, 23 ) have enabled the reconstruction of partial or entire genomes  
20 directly from metagenomic sequencing project.

21  
22 (e) **Reassembly AP**: represent the APs created when an already processed genome is subjected  
23 to different assembly methods to generate a new assembly.

24  
25 (f) **Reannotation AP**: represent the AP created for annotating a genomes that has been  
26 annotated before.

27  
28 (g) **Metatranscriptome Mapping AP**: these APs represent the mapping of the  
29 metatranscriptomic data on the metagenomic sequences in order to connect functional processes  
30 to genes.

## 31 32 33 34 35 36 37 **GOLD BY NUMBERS**

### 38 39 **Studies:**

40  
41 As of September 2014, there are 19,242 Studies in GOLD. These include 472 metagenomic  
42 studies (i.e. have at least one metagenome sequencing project) and 18,770 non-metagenomic  
43 studies. Studies have been generally growing in size and complexity and are increasingly  
44 composed of Biosamples from more diverse environments (Fig. 2). There is also an increasing  
45 number of sequencing strategies applied to each Biosamples (Fig. 2, Table 1) as well a growing  
46 number of APs used within a study (Fig. 3).

### 47 48 49 50 **Biosamples:**

51  
52 There are currently 56,403 Biosamples in the database which are classified as host-associated  
53 (9,922 samples), engineered (1,270 samples), environmental (3,637 samples), and unclassified  
54 (36,435 samples). Organism Biosamples represent more than 150 GOLD phylogenetic  
55  
56  
57  
58  
59  
60

1  
2  
3 classifications. Biome Biosamples represent more than 200 unique GOLD ecosystem  
4 classifications.  
5

### 6 7 **Sequencing Projects:**

8  
9 There are currently 56,458 Sequencing Projects reported in the database. These include 47,932  
10 whole genome sequencing (WGS) projects distributed across 36,824 bacteria, 5,822 eukaryal,  
11 and 851 archaeal projects. There are also 4,351 metagenomic SPs, distributed across 1,567  
12 host-associated, 239 engineered, and 2,545 environmental projects. In addition to the genomic  
13 and metagenomic SP, the database provides information on 1,200 transcriptomic and 797  
14 metatranscriptomic SPs. While there are only 50 targeted gene survey SPs, all of these are part  
15 of Studies that include metagenomic data and most include metatranscriptomic data (Table 1).  
16 The database also provides information on 13 transposon mutagenesis SPs. As this technique is  
17 becoming more high-throughput more projects of this type can be expected (19). A similar growth  
18 is expected for the methylation SPs, only 15 of which are currently available in the database.  
19  
20  
21  
22  
23

### 24 25 **Analysis Projects:**

26  
27 38,573 APs are currently reported of which 36,755 are default APs. For single-cell SPs there are  
28 856 default - screened and 1082 default-unscreened APs. There are also 107 transcriptome  
29 mapping and 80 metatranscriptome mapping APs. Finally, 30 combined assembly APs from 310  
30 SPs in 11 Studies are available in GOLD. All of the Sequencing Projects used for combined  
31 assembly were also used for 'default' APs.  
32  
33  
34

### 35 36 **ACCESSING GOLD**

37  
38 GOLD provides free access to all publicly available data, project status reports and other  
39 statistical information. Data can be accessed by various pre-computed reports or by querying the  
40 database using search functions. Menu tabs to allow users to choose Search, Distribution  
41 Graphs, Biogeographical Metadata and Statistics options to access data are also available from  
42 the front page. A list of all the public projects in the database is also available for download.  
43  
44  
45

### 46 47 **Distribution Graphs:**

48  
49 Automatically generated pie charts that describe the different types of projects in the database  
50 are now available. These include data organized by SP type, sequencing status, phylogenetic  
51 table, phylogenetic tree and Biosample classification in separate tabs.  
52

### 53 54 **Biogeographical Metadata:**

55  
56 The geographic distribution of Biosamples can be visualized via the Google Map and Google  
57 Earth options. These can also be used to select Biosamples based on their geographic location.  
58  
59  
60

1  
2  
3 The Google Map feature aggregates Biosamples by geographic location into circles noting the  
4 number of Biosamples in a group when viewing larger spatial extents. These groupings are  
5 ungrouped as the map view is focused using the zoom feature. The map view can be focused on  
6 the location of a biosample when it is selected from a list next to the map. The Google Earth  
7 feature provides a similar tool but with a 3-dimensional global perspective.  
8  
9

#### 10 11 **Statistics:**

12  
13  
14 The GOLD statistics page provides several precomputed user friendly, easy to interpret graphs,  
15 bar charts and pie charts about various sequencing projects. Refer to supplementary material for  
16 more details about various precomputed charts.  
17

#### 18 19 **SEARCHING THE GOLD DATABASE**

20  
21 The Search function can be used to query the database based on various search criteria that  
22 encompass all four levels of the project classification system or based on various metadata  
23 features. A drop-down menu allows the choice of three search options, Quick Search, Advanced  
24 Search, and Metadata Search.  
25  
26

#### 27 28 **Quick Search:**

29  
30 Quick Search allows a user to search through the most frequently used fields/identifiers across  
31 the four levels in the database (Studies, Biosamples, Sequencing Projects and APs).  
32  
33

#### 34 35 **Advanced Search:**

36  
37 Advanced Search provides options to query metadata fields in each level of the new classification  
38 system. Results are provided as a list according to the search criteria, with fields used displayed  
39 in separate columns. Search result can be redefined by removing any search term by clicking  
40 'remove' next to the search term in the column header. Search results may also be refined  
41 directly in the results table by modifying the search term to any field by clicking the "+" under the  
42 column header. There is also a 'Select Fields' button on the left, which allows the user to add  
43 additional fields.  
44  
45  
46

#### 47 48 **Metadata Search:**

49  
50 Metadata search is designed to query the database using various metadata identifiers. These  
51 include the classification by the domains of the project organism, Archaea, Bacteria, Eukarya or  
52 all. The various search tabs contain graphical and tabular representation of the numbers of  
53 projects or organisms. This approach serves to obtain an overall picture of projects and samples  
54 according to chosen criteria, and produces a sortable table and also plots these lists in a pie-chart  
55 for easy reference.  
56  
57  
58  
59  
60

## CREATING AND EDITING PROJECTS IN GOLD

Registered users can submit new projects or edit their existing entries.

**Editing:** Existing projects can be updated using a new inline-editing user interface. For editing existing entries a user needs to login and select the entry of their interest. When clicking on a field, an edit box is launched with existing values in it. One can update the value and save. The inline-edit feature seamlessly integrates the edit functionality with user interface without the need for launching a separate edit form.

**Creating New Projects:** Registered users can create new SPs using the new project entry interface. Creating a new SP also requires defining all related database entities i.e. Study, Biosample, and Organism when applicable for isolate genome projects. As shown in supplementary material the new project entry landing page provides the following options (a) create a new SP (b) create new AP (c) review your Studies, Biosamples, and Sequencing Projects.

The new SP creation interface will walk a user through a series of steps to define new projects or select existing projects. For example, launching 'Create a new Sequencing Project' will first ask if this is metagenome (biome) or isolate (organism) project. This information is used to launch appropriate forms and guide users through the process. Next, a user will be asked to enter a Study for the SP. If this is a returning user adding additional SPs to an existing Study, the user will be able to choose the existing Study. Otherwise the user will be asked to define a new one. Once the Study is created a Biosample must be defined. Again the user may define a new Biosample or select an existing Biosample. If the SP is for an isolate organism, the user must select an existing organism from the database or define a new organism. After the Study, Biosample and/or organism are created, the user will be able to define a new SP. All the required fields are marked with an asterisk and tool tips are provided with appropriate examples to guide a user in defining new projects. Help pages are available to provide explanation on specific database terminology. If a SP is defined a user can select "Create a new Analysis Project for submission to IMG" to define an AP. A single SP can have multiple APs to represent different assemblies and/or gene calling methodologies applied. Study, Biosample and Organism entries created but not yet associated with a sequencing project are saved as drafts. Users can access these from the "My Data" table as well as select these from the pulldown list as part of new SP creation interface.

## DATA IMPORT AND CURATION CHALLENGES

GOLD continuously monitors sequencing projects around the world both through direct submissions from users and through data imports from major public resources, such as NCBI (7). A series of cross checks have been implemented to ensure high data quality, manually verify data

1  
2  
3 conflicts, and curate metadata during and after import into the database. Due to the nature of  
4 data organization and data quality enforcement standards at different public resources it is  
5 challenging and curation intensive to keep the import processes working. For a list of examples  
6 see the supplementary material.  
7  
8

9 The aim of listing these issues is two fold (i) to express the difficulties that any integrated  
10 public database resource like GOLD is facing in representing disparate information and (ii) to  
11 highlight the need for more manual data curation and quality control checks at major public  
12 resources like NCBI. If the data are corrected at the source it saves time and effort for several  
13 groups around the world. For example correctly representing the sequencing center names and  
14 geographic coordinates at the source would eliminate the need for all other databases who import  
15 data from NCBI to come up with their own procedures for finding and resolving these issues.  
16 NCBI systems serve as a large democratizing force providing unrestricted access for users  
17 around the world to submit their data and freely share with the rest of the world. With such a  
18 broad mandate and unhindered access, it is difficult to enforce strict standards, but at least some  
19 the above listed issues can be mitigated with more manual curation and quality control processes  
20 in place. These challenges are not unique to this database, but to all who rely on public database  
21 resources. Thus there is a strong case for the stakeholders and funding agencies to support data  
22 curation efforts at public resources (23).  
23  
24  
25  
26  
27  
28  
29

### 30 **FUTURE DIRECTIONS:**

31  
32  
33 Future developments will focus on data integration, expanding metadata fields and providing  
34 sophisticated search options across the metadata fields at different classification levels.  
35  
36

37 Data Integration: we will continue importing public metagenome sequencing projects from NCBI  
38 and EBI. We will expand our semi-automatic NCBI isolate genome import process to include  
39 multi-isolate NCBI BioProjects, where more than one isolate genome is listed under a single  
40 NCBI BioProject with different NCBI BioSamples as opposed to represented by individual NCBI  
41 BioProjects.  
42  
43  
44

45  
46 Expanding metadata fields: The growing complexity of the SPs and the diversity of the GOLD  
47 Biosamples collected from specific locations and conditions necessitate GOLD to constantly  
48 expand metadata fields. We plan to incorporate all of the MlxS environmental packages and  
49 include metadata fields that are not currently available in the database.  
50  
51  
52

53 Metadata Miner: The advanced search feature in the current release provides an option to search  
54 among a multitude of metadata fields within each of the four project classification levels. For data  
55 mining and hypothesis generation often it is important to search across different levels using  
56  
57  
58  
59  
60

1  
2  
3 different metadata fields at the same time. For example the search for “aerobic bacterial whole  
4 genome sequencing projects that have a project relevance of medical, with human as a host and  
5 project status of complete” in the current implementation would need to be executed in multiple  
6 steps at different GOLD classification levels. We plan to implement an integrated Metadata Miner  
7 that would facilitate complex searches across all four levels of GOLD. Such an advanced  
8 metadata mining tool will make it easy for users to execute searches similar to the above  
9 example.

## 14 CONCLUSION

16 The steady increase in the number of sequencing studies carried out around the world coupled  
17 with the complexity of the samples, diversity of sequencing strategies, and expanding analysis  
18 methods necessitates an integrated metadata warehouse like GOLD. As outlined above both  
19 through our current release and proposed feature enhancements like Metadata Miner, GOLD is  
20 uniquely positioned to organize sequence metadata and provide unhindered access both for  
21 hypothesis generation and testing. GOLD’s rich metadata coupled with seamless integration with  
22 the IMG analysis systems provides users with the ability to look at their data and analyze results  
23 from a whole different perspective with associated metadata. This helps in understanding the  
24 observations as well as asking questions to find answers hitherto impossible without curated  
25 metadata. Towards this goal GOLD will continue expanding in terms of metadata fields well as  
26 the numbers of projects integrated from various sources around the world.

## 34 ACKNOWLEDGEMENTS

36 We are grateful to all the colleagues who kindly provide information for the more accurate  
37 monitoring of the genome projects and their associated metadata. We thank the members of the  
38 microbial genomics and metagenomics programs at the JGI for support, useful discussions and  
39 exchange of ideas. We would also like to thank Lynn Schriml and Lynette Hirschman for feedback  
40 and useful discussion on the EnvO ontology.

## 45 FUNDING

47 This work was performed under the auspices of the US Department of Energy Office of Science,  
48 Biological and Environmental Research Program, and by the University of California, Lawrence  
49 Berkeley National Laboratory under contract no. DE-AC02-05CH11231.

## 52 REFERENCES

- 54 1. Kyrpides, N.C. (1999) Genomes OnLine Database (GOLD 1.0): a monitor of complete and  
55 ongoing genome projects world-wide. *Bioinformatics*, **15**, 773–774.

- 1  
2  
3 2. Bernal, a, Ear, U. and Kyrpides, N. (2001) Genomes OnLine Database (GOLD): a monitor of  
4 genome projects world-wide. *Nucleic Acids Res.*, **29**, 126–7.
- 5  
6 3. Liolios, K., Tavernarakis, N., Hugenholtz, P. and Kyrpides, N.C. (2006) The Genomes On Line  
7 Database (GOLD) v.2: a monitor of genome projects worldwide. *Nucleic Acids Res.*, **34**,  
8 D332–4.
- 9  
10 4. Liolios, K., Mavromatis, K., Tavernarakis, N. and Kyrpides, N.C. (2008) The Genomes On Line  
11 Database (GOLD) in 2007: status of genomic and metagenomic projects and their  
12 associated metadata. *Nucleic Acids Res.*, **36**, D475–9.
- 13  
14 5. Liolios, K., Chen, I.-M. a, Mavromatis, K., Tavernarakis, N., Hugenholtz, P., Markowitz, V.M.  
15 and Kyrpides, N.C. (2010) The Genomes On Line Database (GOLD) in 2009: status of  
16 genomic and metagenomic projects and their associated metadata. *Nucleic Acids Res.*, **38**,  
17 D346–54.
- 18  
19 6. Pagani, I., Liolios, K., Jansson, J., Chen, I.-M. a, Smirnova, T., Nosrat, B., Markowitz, V.M. and  
20 Kyrpides, N.C. (2012) The Genomes OnLine Database (GOLD) v.4: status of genomic and  
21 metagenomic projects and their associated metadata. *Nucleic Acids Res.*, **40**, D571–9.
- 22  
23 7. Benson, D. a, Cavanaugh, M., Clark, K., Karsch-Mizrachi, I., Lipman, D.J., Ostell, J. and  
24 Sayers, E.W. (2013) GenBank. *Nucleic Acids Res.*, **41**, D36–42.
- 25  
26 8. Leinonen, R., Akhtar, R., Birney, E., Bower, L., Cerdeno-Tárraga, A., Cheng, Y., Cleland, I.,  
27 Faruque, N., Goodgame, N., Gibson, R., et al. (2011) The European Nucleotide Archive.  
28 *Nucleic Acids Res.*, **39**, D28–31.
- 29  
30 9. Markowitz, V.M., Chen, I.-M. a, Palaniappan, K., Chu, K., Szeto, E., Pillay, M., Ratner, A.,  
31 Huang, J., Woyke, T., Huntemann, M., et al. (2014) IMG 4 version of the integrated  
32 microbial genomes comparative analysis system. *Nucleic Acids Res.*, **42**, D560–7.
- 33  
34 10. Meyer, F., Paarmann, D., D'Souza, M., Olson, R., Glass, E.M., Kubal, M., Paczian, T.,  
35 Rodriguez, a, Stevens, R., Wilke, a, et al. (2008) The metagenomics RAST server - a public  
36 resource for the automatic phylogenetic and functional analysis of metagenomes. *BMC*  
37 *Bioinformatics*, **9**, 386.
- 38  
39 11. Field, D., Amaral-Zettler, L., Cochrane, G., Cole, J.R., Dawyndt, P., Garrity, G.M., Gilbert, J.,  
40 Glöckner, F.O., Hirschman, L., Karsch-Mizrachi, I., et al. (2011) The Genomic Standards  
41 Consortium. *PLoS Biol.*, **9**, e1001088.
- 42  
43 12. Yilmaz, P., Kottmann, R., Field, D., Knight, R., Cole, J.R., Amaral-Zettler, L., Gilbert, J. a,  
44 Karsch-Mizrachi, I., Johnston, A., Cochrane, G., et al. (2011) Minimum information about a  
45 marker gene sequence (MIMARKS) and minimum information about any (x) sequence  
46 (MIxS) specifications. *Nat. Biotechnol.*, **29**, 415–20.
- 47  
48 13. Markowitz, V.M., Chen, I.-M. a, Chu, K., Szeto, E., Palaniappan, K., Pillay, M., Ratner, A.,  
49 Huang, J., Pagani, I., Tringe, S., et al. (2014) IMG/M 4 version of the integrated  
50 metagenome comparative analysis system. *Nucleic Acids Res.*, **42**, D568–73.
- 51  
52 14. Markowitz, V.M., Mavromatis, K., Ivanova, N.N., Chen, I.-M. a, Chu, K. and Kyrpides, N.C.  
53 (2009) IMG ER: a system for microbial genome annotation expert review and curation.  
54 *Bioinformatics*, **25**, 2271–8.
- 55  
56  
57  
58  
59  
60

- 1  
2  
3 15. Markowitz, V.M., Chen, I.-M. a, Chu, K., Szeto, E., Palaniappan, K., Jacob, B., Ratner, A.,  
4 Liolios, K., Pagani, I., Huntemann, M., et al. (2012) IMG/M-HMP: a metagenome  
5 comparative analysis system for the Human Microbiome Project. *PLoS One*, **7**, e40151.  
6  
7  
8 16. Nelson, K.E., Weinstock, G.M., Highlander, S.K., Worley, K.C., Creasy, H.H., Wortman, J.R.,  
9 Rusch, D.B., Mitreva, M., Sodergren, E., Chinwalla, A.T., et al. (2010) A catalog of reference  
10 genomes from the human microbiome. *Science*, **328**, 994–9.  
11  
12 17. Wu, D., Hugenholtz, P., Mavromatis, K., Pukall, R., Dalin, E., Ivanova, N.N., Kunin, V.,  
13 Goodwin, L., Wu, M., Tindall, B.J., et al. (2009) A phylogeny-driven genomic encyclopaedia  
14 of Bacteria and Archaea. *Nature*, **462**, 1056–60.  
15  
16 18. Kyrpides, N.C., Hugenholtz, P., Eisen, J. a, Woyke, T., Göker, M., Parker, C.T., Amann, R.,  
17 Beck, B.J., Chain, P.S.G., Chun, J., et al. (2014) Genomic Encyclopedia of Bacteria and  
18 Archaea: Sequencing a Myriad of Type Strains. *PLoS Biol.*, **12**, e1001920.  
19  
20 19. Blow, M.J., Deutschbauer, A.M., Hoover, C.A., Lamson, J., Pennacchio, L.A., Price, M.N.,  
21 Waters, J., Wetmore, K.M., Bristow, J. and Arkin, A.P. (2013) Functional Encyclopedia of  
22 Bacteria and Archaea. Poster session presented at: Genomics of Energy & Environment  
23 User Meeting; Walnut Creek, CA.  
24  
25 20. Ivanova, N., Tringe, S.G., Liolios, K., Liu, W.-T., Morrison, N., Hugenholtz, P. and Kyrpides,  
26 N.C. (2010) A call for standardized classification of metagenome projects. *Environ.*  
27 *Microbiol.*, **12**, 1803–5.  
28  
29 21. Kyrpides, N.C., Woyke, T., Eisen, J. a, Garrity, G., Lilburn, T.G., Beck, B.J., Whitman, W.B.,  
30 Hugenholtz, P. and Klenk, H.-P. (2014) Genomic Encyclopedia of Type Strains, Phase I:  
31 The one thousand microbial genomes (KMG-I) project. *Stand. Genomic Sci.*, **9**, 1278–84.  
32  
33 22. Rinke, C., Schwientek, P., Sczyrba, A., Ivanova, N.N., Anderson, I.J., Cheng, J.-F., Darling,  
34 A., Malfatti, S., Swan, B.K., Gies, E. a, et al. (2013) Insights into the phylogeny and coding  
35 potential of microbial darkmatter- Supplementary Information. *Nature*, **499**, 431–7.  
36  
37 23. Kyrpides, N.C. (2009) Fifteen years of microbial genomics: meeting the challenges and  
38 fulfilling the dream. *Nat. Biotechnol.*, **27**, 627–32.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



## TABLE AND FIGURE LEGENDS

**Table 1.** GOLD Sequencing strategy combinations used within a study.

**Figure 1.** Four level project classification system implemented in v.5 to describe studies, Biosamples, sequencing projects, and analysis projects. Studies group one or more related Biosamples. Biosamples describe an individual sample of genetic material. Sequencing projects are the sequencing deliverables from the Biosamples. Analysis projects are the data processing methods applied to sequencing projects. A) Biosamples may be merged prior to sequencing projects (ex. 16S amplicon data combined prior to sequencing). B) Sequencing Projects may be merged prior to analysis (ex. multiple single-cell genomes combined for assembly).

**Figure 2.** Study Biosamples, ecosystem categories, and sequencing strategies. Each point is a GOLD study. The size of the point represents the number of ecosystem categories within a study. The position on the y-axis notes the number of Biosamples within a study. The color of each point notes the number of unique sequencing strategies used within a study.

**Figure 3.** Sequencing and analysis projects per study over time. Color notes the types of sequencing strategies used within a study. Size indicates the number of analysis projects within a study.

**Table 1.** GOLD Sequencing strategy combinations used within a study.

<b>Sequencing Strategy Combinations</b>	<b>No. Studies</b>	<b>No. Sequencing Projects</b>
Whole Genome Sequencing	18211	46905
Metagenome	403	3315
Transcriptome, Whole Genome Sequencing	76	1989
Metagenome, Metatranscriptome	39	927
Metagenome, Metatranscriptome, Targeted Gene Survey	6	682
Transcriptome	402	596
Metagenome, Whole Genome Sequencing	1	217
Metatranscriptome	9	54
Metagenome, Targeted Gene Survey	4	53
smRNA, Transcription Start Site, Transcriptome, Transposon Mutagenesis Sequencing, Whole Genome Sequencing	1	34
Metatranscriptome, Targeted Gene Survey	1	21
Methylation	4	15
smRNA, Transcriptome	1	14
Plasmid	2	2
smRNA	1	1
Transposon Mutagenesis Sequencing	1	1

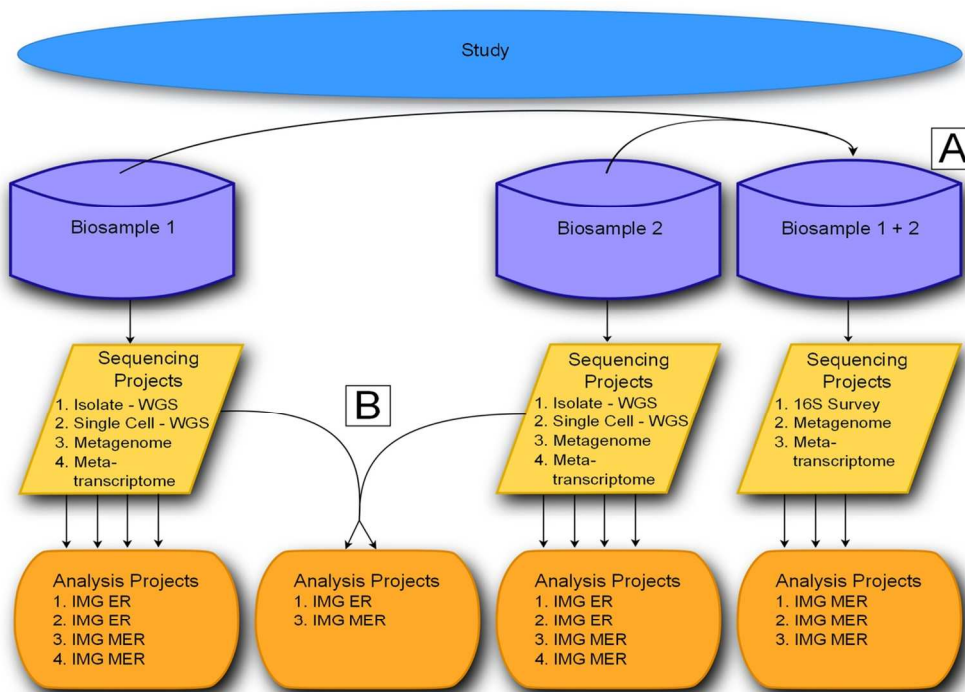


Figure 1  
117x84mm (300 x 300 DPI)

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

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

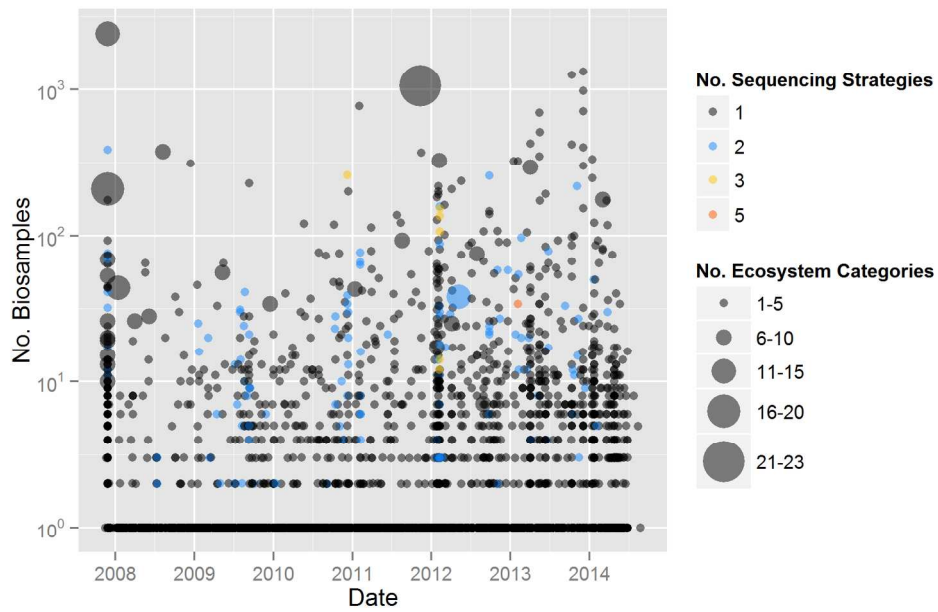


Figure 2

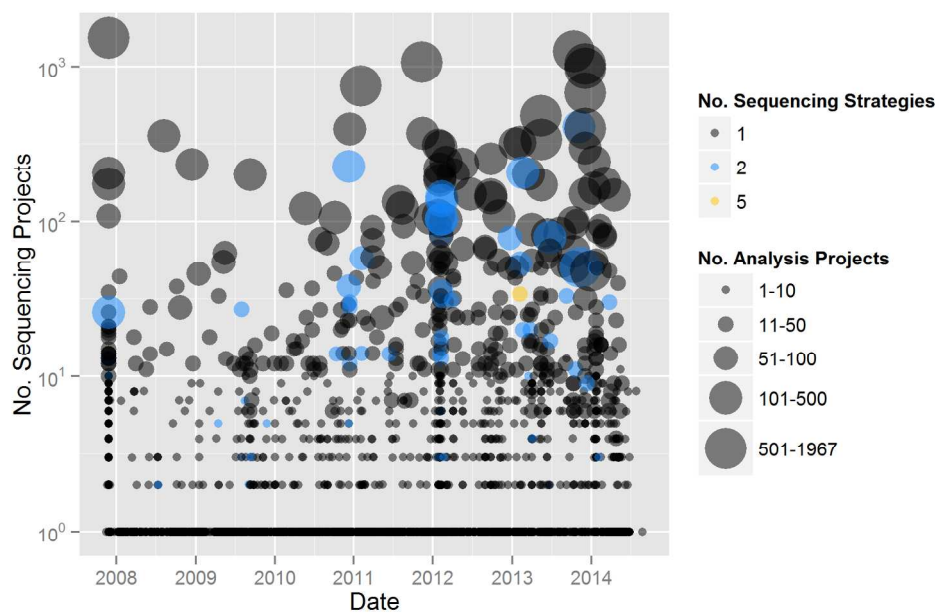


Figure 3

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

## Supplementary Materials

### S1. STATISTICS:

- i) Complete and Permanent Draft Genome Totals (by year and status): This bar chart lists Complete and Permanent Draft genome projects that were completed over the last 9 years.
- ii) Genome Totals (by year and Status): This line diagram tracks the total number of genome sequencing projects registered for each year since 1995. These are further divided into complete and incomplete projects.
- iii) Project Totals (by year and Phylogenetic Group): This graph represents Archaeal, Bacterial, Eukaryotic and Metagenomic projects registered for every year since 2007.
- iv) Phylogenetic distribution of Bacterial Genome Projects: This pie chart breaks down bacterial genome sequencing projects according to major bacterial phyla like Actinobacteria, Firmicutes and Proteobacteria. The number of genomes sequenced in each category are represented in this pie chart.
- v) Projects By Major Sequencing Centers: This pie chart displays the number of genomes sequenced by major sequencing centers worldwide. It may be noted that the top seven major sequencing centers represent 73% of the total projects sequenced.
- vi) Project Relevance of Bacterial Genome Projects: This pie chart shows the breakdown of the number of projects annotated with different relevance terms such as Medical, Environmental, Biotechnological etc.
- vii) Major Sequencing Centers for Archaeal and Bacterial Genomes: This pie chart shows the number of Archaeal and Bacterial genomes sequenced by major sequencing centers.

## S2. CREATING AND EDITING PROJECTS IN GOLD



New project entry landing page for creating new projects in GOLD

### S3. DATA IMPORT AND CURATION CHALLENGES

Here is a select list of issues we face while importing projects into GOLD.

- Frequent data format changes. This necessitates investigating the changes vis-a-vis their overall impact on data import processes in GOLD as well as updating the necessary scripts.

- Ambiguous representation of Genome Projects at BioProject site. Eg: NCBI BioProject ID: 189730. This is a whole genome sequencing project. But it is annotated as Material = Other and Capture = Other instead of correctly representing it as Material = Genome; Capture = Whole. So we can't import such projects automatically. We have to manually review and determine if they are sequencing projects or not.

- Multi-isolate projects. Instead of creating a separate BioProject entry for each isolate genome, NCBI often lists multiple isolate genomes under a single BioProject. Though they are labelled as Scope = "Multiisolate", it require special handling to represent them as in GOLD as individual genome projects. Eg: NCBI BioProject ID: 238952 consists of 6 different individual isolate genomes. In GOLD we represent these as 6 different sequencing projects. As such there is no error on the part of NCBI. This is a classic example of how different resources organize and represent data in their systems. Though such differences look subtle, it often require both engineering and curation efforts to make data imports work.

- We often encounter projects that were listed as multi-isolates, but they are not actually multi-isolate projects. We manually check all multi-isolate projects before importing them into GOLD. It is possible BioProject users specify these erroneously as multi-isolates during submission. Eg: PRJNA238854 appears to be reporting an isolate genome but it is scoped as multi-isolate.

- Cryptic or uninformative sequencing center names on BioProjects. Sequencing center name is a required field in GOLD. We use this info to track sequencing projects carried out at major sequencing centers and provide related statistics on GOLD.

Eg: UH, UB, TDC, 'wqedd wed we' etc. BioProject ID: 34783 lists 'wqedd wed we' as sequencing center name. Again these values are what users provided at the time of project submission to BioProject. We look at the GenBank record or associated publications to infer what UH, UB, TDC etc. mean and curate the same in GOLD.



1  
2  
3 - Inaccurate geolocation info due geolocation name and coordinates mismatch.

4 Eg: BioSamples under PRJNA252784 and PRJNA252785 list China, Liaoning as geographic  
5 location. Where as the coordinates map to Northern California, USA.  
6  
7

8  
9 - Phage genome sequence included as part of host genome sequence. Eg: PRJNA247.

10 Chlamydomonada reinhardtii AR39 genome project also listing Chlamydia phage phiCPAR39  
11 genome sequencing.  
12  
13

14  
15 Duplicate BioProjects: PRJNA243100 is a duplicate entry for PRJNA243545. Both entries are for  
16 the same organism from same institute Shanghai JiaoTong University. In GOLD we choose to list  
17 PRJNA243545 which is associated with GenBank record.  
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