

## **Guidance on importing NCBI genomes/metagenomes into GOLD/IMG**

For all public isolate genome and metagenome projects at GenBank, we have a process in place to import those projects into GOLD and annotate them in IMG. Project information and metadata are curated in GOLD. GenBank and SRA data are annotated in IMG.

We **do not accept** such projects manually entered by users.

We regularly import isolate and metagenomes from NCBI. However, we may not have all of them in GOLD/IMG yet.

For any public GenBank genome, these are the possible scenarios:

- 1) The genome is already in IMG.
- 2) Genome analysis project is defined in GOLD but not yet annotated in IMG.
- 3) We don't have the specific genome (sequencing project/analysis project) in GOLD yet.

What should you do if you want to analyze a public genome from GenBank?

First of all, you need to check whether we already have the genome of your interest in IMG or an Analysis Project in GOLD. The best way to do that is to perform a search by NCBI BioSample Accession as shown in Figure 1.

Step 1 - Click on the number of Sequencing Projects in the top left menu table;

Step 2 - Click on the "Select Columns for Table" button;

Step 3 - Expand Project Information section and select NCBI Biosample Accession field;

Step 4 - Click on the "Submit" button;

Step 5 - Perform the search, using the BioSample accession of your interest.

Figure 1

The screenshot shows a web application interface with a navigation bar at the top containing links: Home, Search, Distribution Graphs, Biogeographical Metadata, SRA Explorer, Statistics, GOLD Usage Policy, Team, and Help.

On the left, there is a sidebar with a table of counts:

Studies	49,626
Biosamples	133,305
Sequencing Projects	411,322
Analysis Projects	321,505
Organisms	411,045

Callout 1 points to the 'Analysis Projects' row.

The main content area has a header with 'Current Filters: None Set' and a 'New Search' button. Callout 2 points to a 'Select Columns for Table' button. Callout 5 points to the 'New Search' button.

Below the header is a table with columns: GOLD Project ID, Project Name, NCBI BioSample Accession, and Project Status. The table contains several rows of data, including project names like '8C63\_RNA' and 'riptome -', and accession numbers like 'SAMN18312926'.

A modal window titled 'Select Fields using the Checkboxes' is open over the table. It contains a 'Submit' button (callout 4), a 'Toggle All' checkbox, and a legend: '√ = required/selected column'. Below the legend is an 'Expand All Fields' button. The modal lists various fields with checkboxes:

- PROJECT INFORMATION**
- GOLD Project ID:
- Project Name:
- Other Names:
- Legacy ER Project ID:
- Legacy ER Sample ID:
- Legacy GOLD ID:
- NCBI BioProject Name:
- NCBI BioProject ID:
- NCBI BioProject Accession:
- NCBI Locus Tag:
- NCBI BioSample Accession:  (callout 3)
- Project Comments:
- Project Status:
- Project Description:
- Is Public:
- Add Date:

For example, you are interested in the public NCBI genome *Methylococcus sp. Phi* with BioSample Accession SAMN04957809. The search returns the seq. project with ID Gp0493101.

Click on the number of the Analysis Projects under the "Your current search results are:" menu (see Figure 2).

Figure 2

Home Search Distribution Graphs Biogeographical Metadata SRA Explorer Statistics GOLD Usage Policy

Studies ⓘ	<a href="#">49,626</a>
Biosamples ⓘ	<a href="#">133,305</a>
Sequencing Projects ⓘ	<a href="#">411,322</a>
Analysis Projects ⓘ	<a href="#">321,505</a>
Organisms	<a href="#">411,045</a>

Your current search results are:

Studies	Biosamples	Organisms	Sequencing Projects	Analysis Projects
1	0	1	1	1

Current Filters:  
Project.NCBI BioSample Accession → [SAMN04957809](#) ✕

[Clear All Filters](#) [New Search](#) [Download Results](#) [Refine Search Filters](#)

[Select Columns for Table](#)

GOLD Project ID	Project Name	NCBI BioSample Accession	Project Status
<a href="#">Gp0493101</a>	Methylacidiphilum sp. Phi	<a href="#">SAMN04957809</a>	Permanent Draft

Click on Analysis Project ID that will open the Analysis Project page. Locate the field IMG Taxon ID at the bottom of the page (see Figure 3). The field is populated with a number 2881253745. (Case 1. The genome is already in IMG.)

Figure 3

Gene Count	2360
Estimated Size	2337855
Sequencing Depth (MIGS-31.1) ⓘ	1000.0x
Annotation Pipeline	IMG Annotation Pipeline v.5.0.18
Added By	JGI automated process on 2020-07-04
Last Modified By	JGI automated process on 2020-08-11
ORF Comparison Tools ⓘ	
16S Recovered ⓘ	
16S Recovery Software ⓘ	
<b>PROJECT EXTERNAL REFERENCES</b>	
IMG Taxon ID ⓘ	<a href="#">2881253745</a>
Genbank Accession ⓘ	Genbank ID <a href="#">LXQC00000000</a> (GCA_004421255.1)
Genbank Release Date	2020-03-28
<b>ANALYSIS PROJECT COMPOSITION</b>	
Study	<a href="#">Methylacidiphilum sp. Phi genome sequencing</a>
Number of Studies ⓘ	<a href="#">1</a>
Number of Projects ⓘ	<a href="#">1</a>
Number of Related Analysis Projects ⓘ	0

**Action:** Click on the number, and link will take you to the IMG page for this genome.

Let's say another search returns a sequencing project, but the field IMG Taxon ID of the associated analysis project is blank. It means that though we already have the analysis project in place for this particular genome, we haven't annotated this genome in IMG yet (Case 2).

**Action:** Send us this AP (or APs) via Contact Us form under the Help section of the top menu, and we will prioritize it (them) for you.

There can be a situation where the search does not return a sequencing project. That means that we do not have a sequencing/analysis project in GOLD yet (Case 3).

**Action:** Send us NCBI BioProject and BioSample Accessions for the specific genome (or a list of Accessions for multiple genomes).

If you want to submit **your own assembly** of a public genome/metagenome that is not in IMG, please contact us.