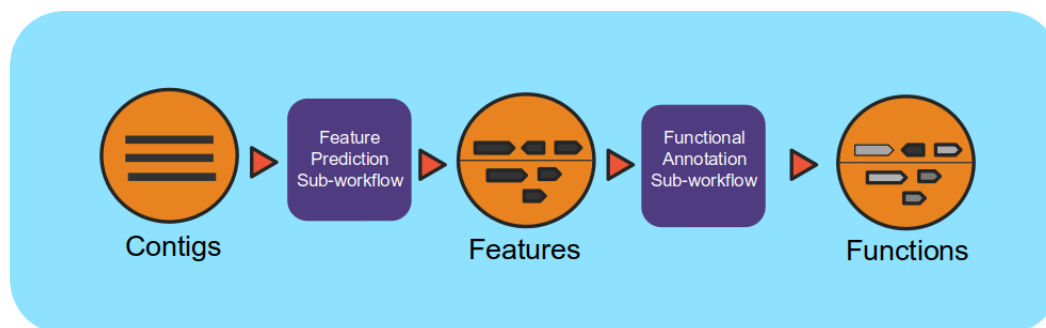


# Metagenome Annotation Workflow (v1.0.0)



## Overview

This workflow takes assembled metagenomes and generates structural and functional annotations.

## Running the Workflow

Currently, this workflow can be run in [NMDC EDGE](#) or from the command line. (CLI instructions and requirements are found [here](#).)

Tutorial videos on how to run each workflow in NMDC EDGE are found [here](#).

## Input

Metagenome Annotation requires assembled contigs in a FASTA file. This input can be the output from the Metagenome Assembly workflow, and this is recommended.

- **Acceptable file formats:** .fasta, .fa, .fna, .fasta.gz, .fa.gz, .fna.gz

## Details

The workflow uses a number of open-source tools and databases to generate the structural and functional annotations. The input assembly is first split into 10MB splits to be processed in parallel. Depending on the workflow engine configuration, the split can be processed in parallel. Each split is first structurally annotated, then those results are used for the functional annotation. The structural annotation uses tRNAscan\_se, RFAM, CRT, Prodigal and GeneMarkS. These results are merged to create a consensus structural annotation. The resulting GFF is the input for functional annotation which uses multiple protein family databases (SMART, COG, TIGRFAM, SUPERFAMILY, Pfam and Cath-FunFam) along with custom HMM models. The functional predictions are created using Last and HMM. These annotations are also merged into a consensus GFF file. Finally, the respective split annotations are merged to generate a single structural annotation file and single functional annotation file. In addition, several summary files are generated in TSV format.

## Software Versions

- Conda
- tRNAscan-SE >= 2.0
- Infernal 1.1.2
- CRT-CLI 1.8
- Prodigal 2.6.3
- GeneMarkS-2 >= 1.07
- Last >= 983
- HMMER 3.1b2
- TMHMM 2.0

## Output

The main outputs are the structural annotation file and the functional annotation file. The functional annotation file can be an input for the MAGs Generation workflow.

Primary Output Files	Description
Structural Annotation	Consensus structural annotation file from multiple tools (.gff)
Functional Annotation	Consensus functional annotation file from multiple tools (.gff)
KEGG summary	KEGG gene function tabular summary (.tsv)
EC summary	Enzyme Commission tabular summary (.tsv)
Gene phylogeny summary	Gene phylogeny tabular summary (.tsv)

## Running the Metagenome Annotation Workflow in NMDC EDGE

### Select a workflow

1. From the Metagenomics category in the left menu bar, select 'Run a Single Workflow'.
2. Enter a **unique** project name with no spaces (underscores are fine).
3. A description is optional, but helpful.
4. Select 'Metagenome Annotation' from the dropdown menu under Workflow.

## Input


This workflow accepts Illumina data in FASTA format as the input; the file can be compressed. It is highly recommended to input the assembled contigs from the Metagenome Assembly workflow. **Acceptable file formats:** .fasta, .fa, .fasta.gz, .fa.gz, fna.gz.

5. Click the button to the right of the input blank for data to select the data file for analysis (if there are separate files, there will be two input blanks). A box called 'Select a File' will open to allow the user to find the desired file(s) from previously run projects, the public data folder, or files uploaded by the user.
6. Then click 'Submit'.

## Output

The General section of the output shows which workflow and which tools were run and the run time information.



File	Size	Last Modified
 MetagenomeAnnotation		
Annotation_Test.faa	20.53 MB	20 days ago
Annotation_Test_cath_funfam.gff	11.89 MB	20 days ago
Annotation_Test_cog.gff	7.92 MB	20 days ago
Annotation_Test_contigs.fna	51.30 MB	20 days ago
Annotation_Test_crt.crisprs	11 kB	20 days ago
Annotation_Test_ec.tsv	1.27 MB	20 days ago
Annotation_Test_functional_annotation.gff	17.43 MB	20 days ago
Annotation_Test_gene_phylogeny.tsv	10.45 MB	20 days ago
Annotation_Test_ko.tsv	2.36 MB	20 days ago
Annotation_Test_ko_ec.gff	44.29 MB	20 days ago
Annotation_Test_pfam.gff	9.71 MB	20 days ago
Annotation_Test_product_names.tsv	5.21 MB	20 days ago
Annotation_Test_proteins.cath_funfam.domtblout	151.86 MB	20 days ago
Annotation_Test_proteins.cog.domtblout	51.46 MB	20 days ago
Annotation_Test_proteins.pfam.domtblout	15.08 MB	20 days ago
Annotation_Test_proteins.smart.domtblout	7.59 MB	20 days ago
Annotation_Test_proteins.supfam.domtblout	339.68 MB	20 days ago
Annotation_Test_proteins.tigrfam.domtblout	3.00 MB	20 days ago
Annotation_Test_smart.gff	3.33 MB	20 days ago
Annotation_Test_structural_annotation.gff	9.99 MB	20 days ago
Annotation_Test_structural_annotation_stats.json	6 kB	20 days ago
Annotation_Test_structural_annotation_stats.tsv	3 kB	20 days ago
Annotation_Test_supfam.gff	12.60 MB	20 days ago
Annotation_Test_tigrfam.gff	1.79 MB	20 days ago
rc	2 B	20 days ago
script	35 kB	20 days ago