

PhenoDB on Cavatica

User Guide

Contents

PhenoDB Cavatica Overview	1
Copy PhenoDB Public Apps Into Your Cavatica Project	1
PhenoDB ANNOVAR Annotation Steps	2
Annotation Inputs Description	3
PhenoDB Analysis: Steps	4
Analysis Inputs Description	5
Inheritance Types: Description and Samples Required	6

PhenoDB Cavatica Overview

PhenoDB is a clinical research tool to analyze VCF files from individuals or families with suspected Mendelian disease. It is now available on Cavatica so data can be analyzed while it remains within the platform's security and computation resources.

Additional information about PhenoDB and its creators is available at <https://phenodb.org/>. Not all PhenoDB analysis functionality in the web version is implemented yet in the Cavatica app.

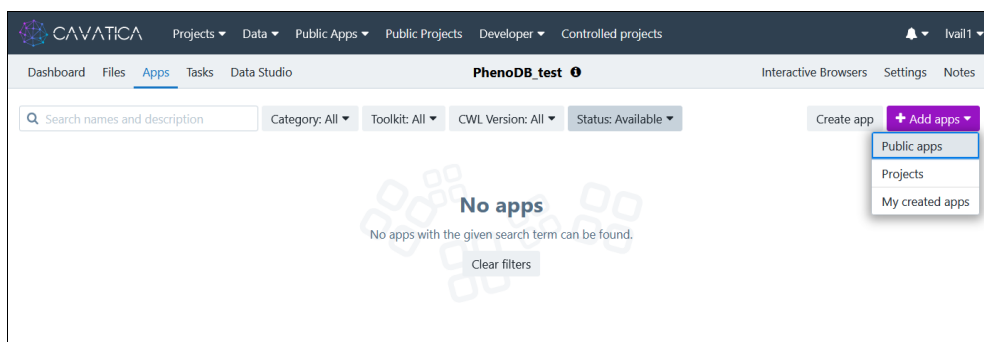
PhenoDB on Cavatica is run in two steps, each in a separate public Cavatica app.

1. PhenoDB ANNOVAR Annotation: annotates VCF files with parameters used for analysis
2. PhenoDB Analysis: reviews the annotated files, and finds variants consistent with user-selected inheritance patterns and other criteria.

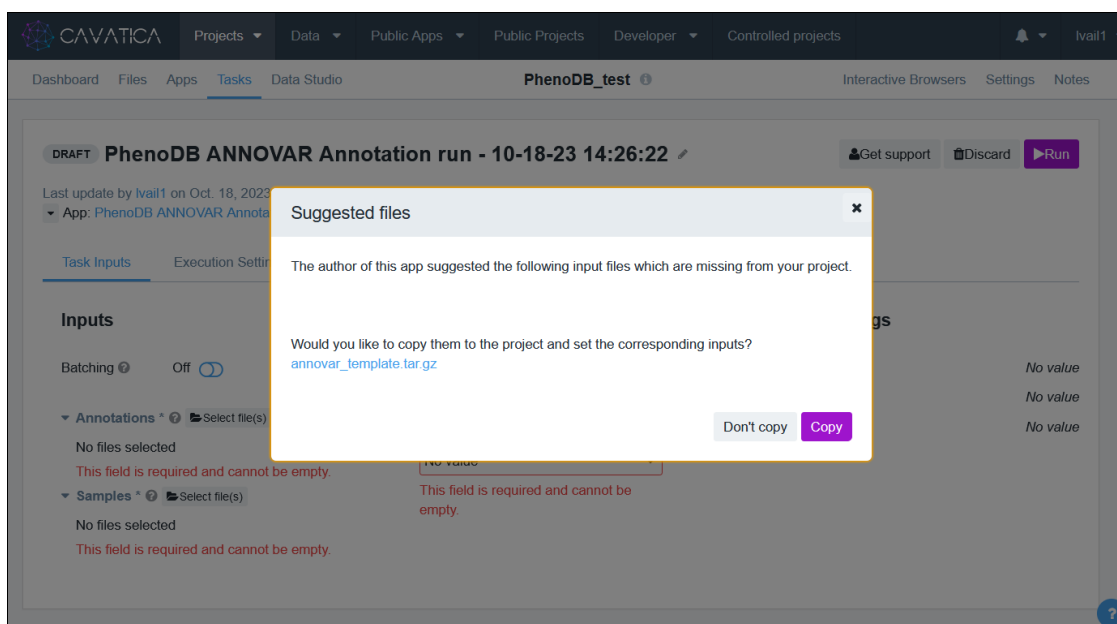
The example inputs used below are available in the public Cavatica project "Commit: PhenoDB App", which also stores the apps. The source code of the apps, its Dockerfile, and Common Workflow Language (CWL) specification are available at <https://github.com/d3b-center/phenodb-cwl-app>. For assistance, email phenodb@jhmi.edu and Laura Vail at lvail1@jhmi.edu.

Copy PhenoDB Public Apps Into Your Cavatica Project

1. In the project where you want to run PhenoDB analyses, open the Apps tab. In the upper left, click "Add Apps," and select "Public Apps" as the source.



2. In the search bar, enter “PhenoDB” and two apps should show up. Click “Copy” for both of them. Close this search window.
3. Refresh your “Apps” page so the copied apps appear, and click “Run” on the PhenoDB ANNOVAR Annotation. A popup will ask whether to copy the ANNOVAR template file into your project: this is a required input, so click “Copy.” This will automatically populate the field “Annotations” for all runs.



PhenoDB ANNOVAR Annotation Steps

1. Enter your inputs, descriptions below. The field “Annotations” should be automatically filled in when you open the app with a file from the public project. Each VCF can only contain data for a single sample.

CAVATICA Projects Data Public Apps Public Projects Developer Controlled projects

Dashboard Files Apps **Tasks** Data Studio **PhenoDB_test** Interactive Browsers Settings Notes

DRAFT PhenoDB ANNOVAR Annotation run - 10-18-23 14:26:22 [Get support](#) [Discard](#) [Run](#)

Last update by lvail1 on Oct. 18, 2023 10:26
App: PhenoDB ANNOVAR Annotation - Revision: 0

Task Inputs Execution Settings

Inputs

Batching [?](#) Off ☐

Annotations * [Change selection](#)

annovar_template.tar.gz

Samples * [Change selection](#)

NA12878-0196534405_proband.vcf.gz
NA12891-0123958385_father.vcf.gz
NA12892-1109184861_mother.vcf.gz

App Settings

[Edit parameters](#) [Show editable](#)

Human_Assembly * [?](#)

Hg38

Output Settings

annotated_vcf No value
debug [?](#) No value
errorFile No value

2. Completing the task may take up to one hour or longer, depending on the size and number of the input files. When completed, there should be one new file for each VCF you entered, with “hg_38_multianno.txt” added to the end of the file's original name

CAVATICA Projects Data Public Apps Public Projects Developer Controlled projects

Dashboard Files Apps **Tasks** Data Studio **PhenoDB_test** Interactive Browsers Settings Notes

COMPLETED PhenoDB ANNOVAR Annotation run - 10-18-23 14:26:22 [Get support](#) [View stats & logs](#) [Edit and rerun](#)

Executed on Oct. 18, 2023 10:31 by lvail1
Spot Instances: On [?](#) Memoization (WorkReuse): Off [?](#) Price: \$0.33 [?](#) Duration: 50 minutes [?](#)
App: PhenoDB ANNOVAR Annotation - Revision: 0

Task Inputs Execution Settings

Inputs

Batching [?](#) On ☒

Annotations [Change selection](#)

annovar_template.tar.gz

Samples [Change selection](#)

NA12878-0196534405_proband.vcf.gz
NA12891-0123958385_father.vcf.gz
NA12892-1109184861_mother.vcf.gz

App Settings

[Show non-default](#)

Human_Assembly [?](#)

Hg38

Output Settings

annotated_vcf

NA12878-0196534405_proband.hg38_multianno.txt
NA12891-0123958385_father.hg38_multianno.txt
NA12892-1109184861_mother.hg38_multianno.txt

debug [?](#)

debug.txt
errorFile No value

3. These output files can then be analyzed with the PhenoDB Analysis app

Annotation Inputs Description

- VCF File(s): Individual VCF files for all people included in the analysis
 - This step can be done for multiple families at once
 - The files can be either .gz compressed, or uncompressed
 - Set “Batching” to On for the VCFs, and select “Batch by: File.” This helps the process run faster
- Human_Assembly: Select either Hg19 or Hg38
- Annotations: ANNOVAR reference data for the annotation and analysis. Only use the file “annovar_template.tar.gz” from the public project

PhenoDB Analysis: Steps

1. Begin with annotated VCF files generated by the 'PhenoDB ANNOVAR Annotation' app
2. Enter your inputs, and descriptions of each item are below. Please note that there will be red text when you first click “+” to add a new Sample, and it will go away once all the fields for that Sample have been filled in. You can enter any number of Samples, but all should be related to the proband.

The screenshot shows the CAVATICA web interface for the 'PhenoDB Analysis' app. The top navigation bar includes 'Dashboard', 'Files', 'Apps', 'Tasks', and 'Data Studio'. The app title is 'PhenoDB_test'. Below the title, there's a 'DRAFT' status and a timestamp 'PhenoDB Analysis run - 10-18-23 15:32:31'. A 'Last update by' field shows 'Ivalit on Oct. 18, 2023 11:32'. The 'App' is 'PhenoDB Analysis - Revision: 0'. The interface is divided into three main sections: 'Inputs', 'App Settings', and 'Output Settings'.

Inputs:

- Batching:** Off (toggle)
- Samples:** A list of samples with columns for 'Affected_Status', 'Relationship', 'Sex', and 'VCF'. The first sample is 'NA12878-0196534405_proband hg38_multiturno.txt'.
- App Settings:**
 - Analysis Type(s):** Autosomal recessive - Compound, Autosomal recessive - Homozygote, Autosomal dominant - Variants.
 - Exclude minor allele frequency greater than:** 0.01.
 - RefGene_Gene_Location:** exonic, exonic:splicing, splicing.
- Output Settings:**
 - Analysis result:** No value
 - analysis_summary:** No value
 - debug:** No value
 - exceptions:** No value

3. When complete, there should be two files for each selected analysis type: one 'Analysis result' tsv containing the variants of interest for that inheritance pattern, and one 'Analysis summary' file describing the filtering steps that led to the result.

CAVATICA Projects Data Public Apps Public Projects Developer Controlled projects

Dashboard Files Apps **Tasks** Data Studio **PhenoDB_test** Interactive Browsers Settings Notes

COMPLETED PhenoDB Analysis run - 10-18-23 15:32:31 [Get support](#) [View stats & logs](#) [Edit and rerun](#)

Executed on Oct. 18, 2023 11:36 by hval1

Spot Instances: On [M](#) Memoization (Work/Reuse): Off [M](#) Price: \$0.01 [M](#) Duration: 2 minutes [M](#)

App: PhenoDB Analysis - Revision 0

Inputs

- Samples**
 - Affected_Status: Affected
 - Relationship: Proband
 - Sex: Male
 - VCF: NA12878-0196534405_proband.hg38_multianno.txt
- Affected_Status: Unaffected
- Relationship: Mother
- Sex: Female
- VCF: NA12892-1109184861_mother.hg38_multianno.txt
- Affected_Status: Affected
- Relationship: Father
- Sex: Male
- VCF: NA12891-0123956385_father.hg38_multianno.txt

App Settings Show non-default

- Analysis Type(s)**
 - Autosomal recessive - Compound heterozygous
 - Autosomal recessive - Homozygous
 - Autosomal dominant - Variants
- Exclude minor allele frequency greater than: 0.01
- RefGene_Gene_Location**
 - exonic
 - exonic splicing
 - splicing

Output Settings

- Analysis result**
 - PhenoDB_Analysis_AD_V_2023_10_18_11-38-28.tsv
 - PhenoDB_Analysis_AR_CH_2023_10_18_11-38-28.tsv
 - PhenoDB_Analysis_AR_H_2023_10_18_11-38-28.tsv
- analysis_summary**
 - Log_AD_V_2023_10_18_11-38-28.txt
 - Log_AR_CH_2023_10_18_11-38-28.txt
 - Log_AR_H_2023_10_18_11-38-28.txt
- debug**
 - _1_debug.txt
 - exceptions: No value

4. Cavatica has built-in preview for text files, including column sort in the .tsv analysis results. Click on the names of the files, and select “Preview” to view and sort the .tsv files, and “Raw View” for the .txt files. In the .tsv preview, sort the results by clicking on column headers.

CAVATICA Projects Data Public Apps Public Projects Developer Controlled projects

Files **PhenoDB_Analysis_AD_V_2023_10_18_11-38-28.tsv** [Edit Metadata](#) [Copy](#) [Move](#) [Tags](#) [Download](#) [...](#)

1.9 MB (1,959,573 bytes) - Produced on October 18, 2023 11:38 (Eastern Daylight Time), by PhenoDB Analysis run - 10-18-23 15:32:31 - Hosted on AWS (us-east-1) [M](#)

Metadata Raw View **Preview**

Search [D](#)

Chr	Start	End	Ref	Alt	Func.refGene	Gene.refGene	GeneDetail.refGene
1	970552	970553	GG	-	exonic	PLEKH1	
1	970553	970553	G	-	exonic	PLEKH1	
1	1290276	1290276	C	G	exonic	SCNN1D	
1	1599981	1599981	C	T	exonic	FNDC10	
1	1636044	1636044	G	A	exonic	CDK11B	
1	1916890	1916890	C	T	exonic	CALML6	
1	1922303	1922303	C	T	exonic	CFAP74	

Analysis Inputs Description

- Samples: One or more individuals to be analyzed together. Each Sample has four attributes to be entered:
 - Affected Status: Affected, Unaffected, or Unknown
 - Relationship: Proband, or relationship to proband (Mother, Father, Other Relative)
 - Sex: Male, Female, Unknown or not XX/XY

- VCF: Individual's annotated VCF file from the 'PhenoDB ANNOVAR Annotation' app, which end with "multianno.txt"
- Analysis Type(s): Choose one or more inheritance patterns for the analysis, described below. Recommended default: Autosomal recessive compound heterozygous and Autosomal recessive homozygous
- Exclude minor allele frequency greater than: Choose your cutoff point. Data sources are the ExAC, esp6500siv2, 1000g2014oct, and 1000g2015aug databases. Recommended default: 0.01
- RefGene gene location: Select one or more to include. Recommended default, three locations: "exonic", "exonic;splicing", and "splicing"
- Inheritance Types: Choose as appropriate for your analysis from the below descriptions. Recommended default:

Inheritance Types: Description and Samples Required

AR_CH: Autosomal Recessive Compound Heterozygous

- Requires: At least proband
- Variants in Result: Includes only the heterozygous variants identified in the proband (assumed to be affected); if there is more than one affected family members, the analysis include only the variants that are identified in all affected members; next, the analysis includes only the genes that have more than one variant in the proband but if the same set of variants in a gene is found in one of the parents or in other unaffected family member then this gene (and its variants) is excluded of the analysis (Sobreira et al., 2015)

AR_H: Autosomal Recessive Homozygous

- Requires: At least proband
- Variants in Result: Identifies homozygous variants that are shared by all affected individuals and excludes variants that are homozygous in an unaffected individual (Sobreira et al., 2015)

AD_NM: Autosomal Dominant New Mutation

- Requires: Proband and two parents (trio), and neither parent is affected
- Variants in Result: Includes heterozygous variants that are identified in the proband but not in the parents (Sobreira et al., 2015)

AD_IM: Autosomal Dominant Inherited Mutation

- Requires: Proband plus at least one affected or unaffected relative, not necessarily a parent
- Variants in Result: Retains heterozygous variants that are shared by affected individuals and excludes those found in unaffected individuals (Sobreira et al., 2015)

AD_V: Autosomal Dominant Variant

- Requires: Should be used when only one individual is being analyzed.
- Variants in Result: Retains heterozygous variants with a minor allele frequency (MAF) less than the threshold selected for the ExAC, esp6500siv2, 1000g2014oct, and 1000g2015aug databases