

---

# **AgBase Documentation**

*Release 1.0*

**Fiona M. McCarthy**

**Oct 13, 2021**



<b>1 Acknowledgements</b>	<b>3</b>
<b>2 Contact Us</b>	<b>5</b>



AgBase is a curated, open-source, Web-accessible resource for functional analysis of agricultural plant and animal gene products. Our long-term goal is to serve the needs of the agricultural research communities by facilitating post-genome biology for agriculture researchers and for those researchers primarily using agricultural species as biomedical models.

We use controlled vocabularies developed by the [Gene Ontology \(GO\) Consortium](#) to describe molecular function, biological process, and cellular component for genes and gene products in agricultural species. For more information about the AgBase database please visit our [Educational Resources](#) page or refer to our [AgBase publications](#) .

AgBase will also accept annotations from any interested party in the research communities. AgBase develops freely available tools for functional analysis, including tools for using GO. We appreciate any and all questions, comments, and suggestions. Please send us your ideas about how to make AgBase more useful for you.



---

## Acknowledgements

---

AgBase acknowledges the following groups for their help and support: The GO Consortium, especially [DictyBase](#) for providing the database schema and for technical assistance with implementation, [MGI](#) for providing training and continued support with manual curation issues and the [EBI GOA Project](#) for allowing us access to their tools and for their continued help, support and patience.

### 1.1 AgBase has received financial support from

- Mississippi State University
- Office of Research and Economic Development (ORED)
- Division of Agriculture, Forestry and Veterinary Medicine (DAFVM)
- Mississippi Agricultural and Forestry Experiment Station (MAFES)
- College of Veterinary Medicine
- Bagley College of Engineering
- Institute for Genomics, Biocomputing & Biotechnology (IGBB; formerly the Life Sciences and Biotechnology Institute)

### 1.2 Competitive Grants

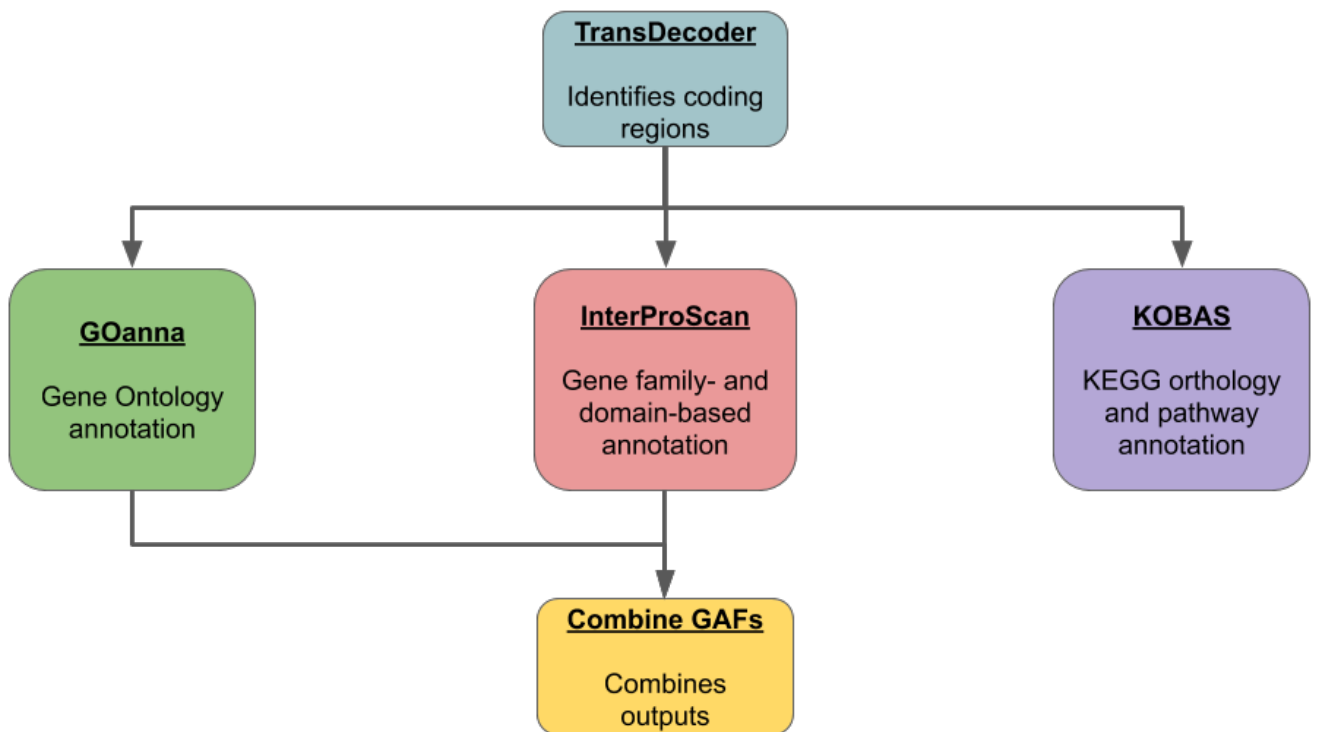
- USDA Agriculture and Food Research Initiative Competitive Grant no. 2011-67015-30332
- National Research Initiative of the USDA Cooperative State Research, Education and Extension Service, grant number 2007-35205-17941
- National Institutes of Health NIGMS project 07111084
- NSF EPSCoR award number EPS 0903787





agbase@email.arizona.edu

## 2.1 Functional Annotation Workflow



This functional annotation workflow employs three annotation tools:

1. **GOanna:** It performs a BLAST search and transfers gene ontology (GO) annotations from BLAST matches to the query gene products.
2. **InterProScan:** InterPro is a database which integrates together predictive information about proteins' function from a number of partner resources, giving an overview of the families that a protein belongs to and the domains and sites it contains. InterProScan can also provide GO and pathway annotations.
3. **KOBAS:** It uses BLAST to annotate the input with KEGG Orthology terms and KEGG pathways

Results and analysis from the application of this functional annotation workflow to the Official gene set v3.0 protein set from *Diaphorina citri* followed by a differential expression analysis was presented at a seminar in the University of Arizona Animal and Comparative Biomedical Sciences in Fall 2020. The [slides](#) and [video](#) are available online.

**Citation:** Please cite the following preprint if you use annotation results from the workflow

Saha, S.; Cooksey, A.M.; Childers, A.K.; Poelchau, M.F.; McCarthy, F.M. Workflows for Rapid Functional Annotation of Diverse Arthropod Genomes. *Insects* 2021, 12, 748. <https://doi.org/10.3390/insects12080748>

---

**Note:** Each of these tools accepts a peptide FASTA file. For those users with nucleotide sequences some documentation has been provided for using **TransDecoder** (although other tools are also acceptable). The [TransDecoder app](#) is available through CyVerse or as a [BioContainer](#) for use on the command line.

---

---

**Note:** As both GOanna and InterProScan provide GO annotations, their outputs are provided in GAF format. The '**Combine GAFs**' tool can then be used to make a single GAF of GO annotations, if desired.

---

## 2.2 Intro

- GOanna performs a BLAST search, allows you to filter based on BLAST match parameters and transfers Gene Ontology (GO) functional annotations from the BLAST matches to your input genes.
- GOanna accepts a protein FASTA file as input.
- BLAST databases are created by AgBase based upon proteins that have GO available and subsetted by phyla. We recommend selecting the database most closely related to the sequence used as input.
- We strongly recommend selecting only GO annotations based on experimental evidence codes. This will ensure the best quality annotations for your data.
- The remaining parameters are standard BLAST parameters. More information on determining the best BLAST parameters for your specific data set can be found in the section below.

Results and analysis from the application of GOanna to the Official gene set v3.0 protein set from *Diaphorina citri* followed by a differential expression analysis was presented at a seminar in the University of Arizona Animal and Comparative Biomedical Sciences in Fall 2020. The [slides](#) and [video](#) are available online.

### 2.2.1 Where to Find GOanna

- [Docker Hub](#)
- [CyVerse Discovery Environment](#)
- [AgBase](#)

## 2.2.2 Getting the GOanna Databases

To run the tool you need some public data. These files are now available as gzipped files to aid downloading. The directories are best downloaded with `iCommands`. Once `iCommands` is `setup` you can use 'iget' to download the data.

- 1) `agbase_database`: species subset to run BLAST against (this command will download the entire directory)

```
iget -rPT /iplant/home/shared/iplantcollaborative/protein_blast_dbs/agbase_database
```

- 2) `go_info`: Uniprot GO annotations (this command will download the entire directory)

```
iget -rPT /iplant/home/shared/iplantcollaborative/protein_blast_dbs/go_info
```

**Note:** Each of these tools accepts a peptide FASTA file. For those users with nucleotide sequences some documentation has been provided for using **TransDecoder** (although other tools are also acceptable). The **TransDecoder app** is available through CyVerse or as a **BioContainer** for use on the command line.

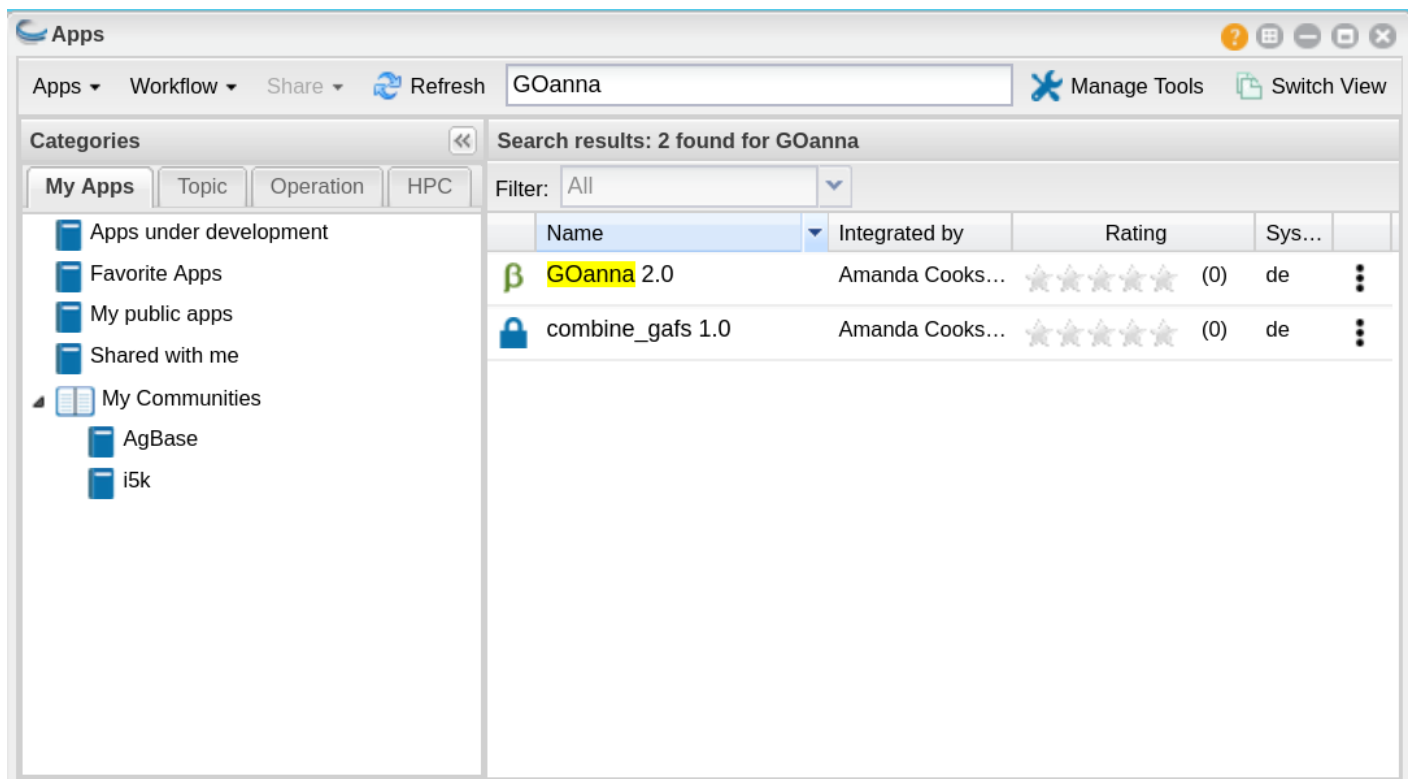
## 2.2.3 Help and Usage Statement

```
Options:
-a BLAST database basename ('arthropod', 'bacteria', 'bird', 'crustacean', 'fish',
↳ 'fungi', 'human', 'insecta',
  'invertebrates', 'mammals', 'nematode', 'plants', 'rodents' 'uniprot_sprot',
↳ 'uniprot_trembl', 'vertebrates'
  or 'viruses')
-c peptide fasta filename
-o output file basename
[-b transfer GO with experimental evidence only ('yes' or 'no'). Default = 'yes'.]
[-d database of query ID. If your entry contains spaces either substitute and
↳ underscore (_) or,
  to preserve the space, use quotes around your entry. Default: 'user_input_db']
[-e Expect value (E) for saving hits. Default is 10.]
[-f Number of aligned sequences to keep. Default: 3]
[-g BLAST percent identity above which match should be kept. Default: keep all
↳ matches.]
[-h help]
[-m BLAST percent positive identity above which match should be kept. Default: keep
↳ all matches.]
[-s bitscore above which match should be kept. Default: keep all matches.]
[-k Maximum number of gap openings allowed for match to be kept. Default: 100]
[-l Maximum number of total gaps allowed for match to be kept. Default: 1000]
[-q Minimum query coverage per subject for match to be kept. Default: keep all
↳ matches]
[-t Number of threads. Default: 8]
[-u 'Assigned by' field of your GAF output file. If your entry contains spaces (eg.
↳ firstname lastname)
  either substitute and underscore (_) or, to preserve the space, use quotes around
↳ your entry (eg. "firstname lastname")
  Default: 'user']
[-x Taxon ID of the query species. Default: 'taxon:0000']
[-p parse_deflines. Parse query and subject bar delimited sequence identifiers]
```

## 2.3 GOanna on CyVerse

### 2.3.1 Accessing GOanna in the Discovery Environment

1. Create an account on CyVerse (free)
2. Open the CyVerse Discovery Environment (DE) and login with your CyVerse credentials.
3. If you are new to the Discovery Environment (DE) the user guide can be found [here](#).
4. Click on the 'Data' button at the left side of the screen to access your files/folders. Upload your data to the DE.
5. To access the [GOanna 2.0](#) app click on the 'Apps' button at the left side of the DE.
6. Search for 'GOanna' in the search bar at the top of the 'apps' window (see below). The contents of the folder will appear in the main pane of the window. The GOanna app is called 'GOanna 2.0'; click on the name to open the app.



#### Find Apps Easily with 'Communities'

The GOanna 2.0 app belongs to the 'i5k' and 'AgBase' communities. You can join either of these communities and they will appear in the left-hand pane of your 'Apps' window (see above).

To join a community click on the person icon in the top-right corner of the Discovery Environment window and select 'Communities'. In the 'Communities' window choose 'all communities' from the drop-down list. A list of communities will appear in the main pane of this window. Select the one you wish to join by clicking on it and then clicking on the 'join' button.

## 2.3.2 Using the GOanna App

### Launching the App

GOanna 2.0

Analysis Name: GOanna\_2.0\_analysis1

Analysis Name:  
GOanna\_2.0\_analysis1

Comments:

Select output folder:  
/iplant/home/amcooksey/analyses Browse

Retain Inputs? Enabling this flag will copy all the input files into the analysis result folder.

\* Input

Parameters

\* Output

Create Quick Launch Launch Analysis

**Analysis Name:GOanna\_2.0\_analysis1:** This menu is used to name the job you will run so that you can find it later. Analysis Name: The default name is “GOanna\_2.0\_analysis1”. We recommend changing the ‘analysis1’ portion of this to reflect the data you are running.

**Comments:** (Optional) You can add additional information in the comments section to distinguish your analyses further.

**Select output folder:** This is where your results will be placed. The default (recommended) is your ‘analyses’ folder.

**Retain Inputs:** Enabling this flag will copy all the input files into the analysis result folder.

**Warning:** Selecting this option will rapidly consume your allocated space. It is not recommended. Your inputs will always remain available in the folder in which you stored them.

### Input

This menu is used to select the BLAST database and your input file.

**BLAST database basename:** BLAST databases are created by AgBase based upon proteins that have GO available and subsetted by phyla. We recommend selecting the database most closely related to the sequence used as input.

**Peptide FASTA file:** Use the Browse button on the right hand side to navigate to your Data folder and select your protein sequence file.

### Parameters

Use this menu to select your BLAST parameters.

**Transfer GO with experimental evidence only:** We strongly recommend selecting the “yes” option from the drop-down menu so that only GO annotations based experimental evidence codes will be transferred . This will ensure the best quality annotations for your data.

The remaining parameters are standard BLAST parameters, and their defaults can be seen by hovering your cursor over the blue i.

### Determining BLAST Parameters to Use

BLAST parameters are contingent on the BLAST database used and the composition of the input file, and so will change for each analysis.

---

**Tip:** Make a subset of 100 randomly selected sequences from your larger dataset and use this as the input for GOanna to test for parameters that give good alignments. The [Split FASTA file](#) app in the CyVerse Discovery Environment can be used to make a subsetted file.

---

1. To test for good parameters use GOanna (either in CyVerse or online at AgBase) by selecting the same database you will use and setting relaxed parameters (e.g., in the CyVerse instance of GOanna, use defaults).
2. Once you have run your subsetted file, use the html file to view alignments, select good alignments and note the parameters for these.

**Parse query and subject bar delimited sequence identifiers:** This option should be selected if you are using a protein fasta file from NCBI. These fasta files have headers with pipes that will not format correctly otherwise.

---

**Note:** This tool uses stand alone BLAST and interprets FASTA defines accordingly. NCBI ‘bar’-delimited defines can be interpreted correctly using the ‘parse-defines’ option. If the ‘parse-defines’ option is not checked then BLAST will interpret the ID to be everything before the first space.

---

### Output

This menu is used to format your GO annotation results into a standard gene association file format.

**Output File basename:** This will be the prefix for your output files. A good name choice is to use fasta file name.

**Database of query ID:** Use the database that sequences were obtained from (Genbank), or a recognizable project name if these sequences are not in a database (e.g., i5k project or Smith Lab).The default is ‘user\_input\_db’.

**‘Assigned by’ field of your GAF output file:** Enter your name. This field is used to track who made the annotations. The default is ‘user’.

**Taxon ID of the query species:** Enter the NCBI taxon number for your species. This can be found by searching for your species name (common or scientific) in the [NCBI taxon database](#). The default is “0000”.

### 2.3.3 Understanding Your Results

If all goes well, you should get 4 output files and a ‘logs’ folder.

**<basename>.asn:** This is standard BLAST output format that allows for conversion to other formats. You probably won't need to look at this output.

**<basename>.html:** This output displays in your web browser so that you can view pairwise alignments to determine BLAST parameters.

**<basename>.tsv:** This is the tab-delimited BLAST output that can be opened and sorted in Excel to determine BLAST parameter values. The file contains the following columns:

- query ID
- query length
- query start
- query end
- subject ID
- subject length
- subject start
- subject end
- e-value
- percent ID
- query coverage
- percent positive ID
- gap openings
- total gaps
- bitscore
- raw score

For more information on the BLAST output parameters see the [NCBI BLAST documentation](#).

**<basename>\_goanna\_gaf.tsv:** This is the standard tab-separated [GO annotation file format](#) that is used by the GO Consortium and by software tools that accept GO annotation files to do GO enrichment.

If you see more files in your output folder there may have been an error in the analysis or there may have been no GO to transfer. Check the 'condor\_stderr' file in the analysis output 'logs' folder.

If that doesn't clarify the problem contact us at [agbase@email.arizona.edu](mailto:agbase@email.arizona.edu) or [support@cyverse.org](mailto:support@cyverse.org).

## 2.4 GOanna on the Command Line

### 2.4.1 Getting the databases

To run the tool you need some public data. These files are now available as gzipped files to aid downloading. The directories are best downloaded with [iCommands](#). Once [iCommands](#) is [setup](#) you can use 'iget' to download the data.

- 1) `agbase_database`: species subset to run BLAST against (this command will download the entire directory)

```
iget /iplant/home/shared/iplantcollaborative/protein_blast_dbs/agbase_database
```

- 2) `go_info`: Uniprot GO annotations (this command will download the entire directory)

```
iget /iplant/home/shared/iplantcollaborative/protein_blast_dbs/go_info
```

## 2.4.2 Container Technologies

GOanna is provided as a Docker container.

A container is a standard unit of software that packages up code and all its dependencies so the application runs quickly and reliably from one computing environment to another.

There are two major containerization technologies: **Docker** and **Singularity**.

Docker containers can be run with either technology.

## 2.4.3 Running GOanna using Docker

---

### About Docker

- Docker must be installed on the computer you wish to use for your analysis.
- To run Docker you must have ‘root’ permissions (or use sudo).
- Docker will run all containers as ‘root’. This makes Docker incompatible with HPC systems (see Singularity below).
- Docker can be run on your local computer, a server, a cloud virtual machine (such as CyVerse Atmosphere) etc. Docker can be installed quickly on an Atmosphere instance by typing ‘ezd’.
- For more information on installing Docker on other systems see this tutorial: [Installing Docker on your machine](#).

---

### Getting the GOanna container

The GOanna tool is available as a Docker container on Docker Hub: [GOanna container](#)

The container can be pulled with this command:

```
docker pull agbase/goanna:2.0
```

---

### Remember

You must have root permissions or use sudo, like so:

```
sudo docker pull agbase/goanna:2.0
```

---

### Running GOanna with Data

#### Getting the Help and Usage Statement

```
sudo docker run --rm -v $(pwd) :/work-dir agbase/goanna:2.0 -h
```



See *Help and Usage Statement*

**Tip:** There are 3 directories built into this container. These directories should be used to mount data.

- /agbase\_database
- /go\_info
- /work-dir

GOanna has three required parameters:

```
-a BLAST database basename (acceptable options are listed in the help/usage)
-c peptide FASTA file to BLAST
-o output file basename
```

### Example Command

```
sudo docker run \
--rm \
-v /home/amcooksey/i5k/agbase_database:/agbase_database \
-v /home/amcooksey/i5k/go_info:/go_info \
-v $(pwd):/work-dir \
agbase/goanna:2.0 \
-a invertebrates \
-c AROS_10.faa \
-o AROS_10_invert_exponly \
-p \
-g 70 \
-s 900 \
-d RefSeq \
-u "Amanda Cooksey" \
-x 37344 \
-k 9 \
-q 70
```

### Command Explained

**sudo docker run:** tells docker to run

**--rm:** removes the container when the analysis has finished. The image will remain for future use.

**-v /home/amcooksey/i5k/agbase\_database:/agbase\_database:** tells docker to mount the ‘agbase\_database’ directory I downloaded to the host machine to the ‘/agbase\_database’ directory within the container. The syntax for this is: <absolute path on host>:<absolute path in container>

**-v /home/amcooksey/i5k/go\_info:/go\_info:** mounts ‘go\_info’ directory on host machine into ‘go\_info’ directory inside the container

**-v \$(pwd):/work-dir:** mounts my current working directory on the host machine to ‘/work-dir’ in the container

**agbase/goanna:2.0:** the name of the Docker image to use

**Tip:** All the options supplied after the image name are GOanna options

- a invertebrates:** GOanna BLAST database to use—first of three required options.
- c AROS\_10.faa:** input file (peptide FASTA)—second of three required options
- o AROS\_10\_invert\_exponly:** output file basename—last of three required options
- p:** our input file has NCBI deflines. This specifies how to parse them.
- g 70:** tells GOanna to keep only those matches with at least 70% identity
- s 900:** tells GOanna to keep only those matches with a bitscore above 900
- d RefSeq:** database of query ID. This will appear in column 1 of the GAF output file.
- u “Amanda Cooksey”:** name to appear in column 15 of the GAF output file
- x 37344:** NCBI taxon ID of input file species will appear in column 13 of the GAF output file
- k 9:** tells GOanna to keep only those matches with a maximum number of 9 gap openings
- q 70:** tells GOanna to keep only those matches with query coverage of 70 per subject

### Understanding Your Results

If all goes well, you should get 4 output files:

**<basename>.asn:** This is standard BLAST output format that allows for conversion to other formats. You probably won't need to look at this output.

**<basename>.html:** This output displays in your web browser so that you can view pairwise alignments to determine BLAST parameters.

**<basename>.tsv:** This is the tab-delimited BLAST output that can be opened and sorted in Excel to determine BLAST parameter values. The file contains the following columns:

- Query ID
- query length
- query start
- query end
- subject ID
- subject length
- subject start
- subject end
- e-value
- percent ID
- query coverage
- percent positive ID
- gap openings
- total gaps
- bitscore
- raw score

For more information on the BLAST output parameters see the [NCBI BLAST documentation](#).

**<basename>\_goanna\_gaf.tsv:** This is the standard tab-separated [GO annotation file format](#) that is used by the GO Consortium and by software tools that accept GO annotation files to do GO enrichment.

If you see more files in your output folder there may have been an error in the analysis or there may have been no GO to transfer. [Contact us](#).

---

## 2.4.4 Running GOanna using Singularity

---

### About Singularity

- does not require ‘root’ permissions
  - runs all containers as the user that is logged into the host machine
  - HPC systems are likely to have Singularity installed and are unlikely to object if asked to install it (no guarantees).
  - can be run on any machine where it is installed
  - more information about [installing Singularity](#)
  - This tool was tested using Singularity 3.0. Users with Singularity 2.x will need to modify the commands accordingly.
- 

---

### HPC Job Schedulers

Although Singularity can be installed on any computer this documentation assumes it will be run on an HPC system. The tool was tested on a PBSPro system and the job submission scripts below reflect that. Submission scripts will need to be modified for use with other job scheduler systems.

---

### Getting the GOanna Container

The GOanna tool is available as a Docker container on Docker Hub: [GOanna container](#)

The container can be pulled with this command:

```
singularity pull docker://agbase/goanna:20
```

### Running GOanna with Data

### Getting the Help and Usage Statement

#### Example PBS script:

```
#!/bin/bash
#PBS -N goanna
#PBS -W group_list=fionamcc
#PBS -l select=1:ncpus=28:mem=168gb
#PBS -q standard
#PBS -l walltime=6:0:0
#PBS -l cput=168:0:0
```

(continues on next page)

(continued from previous page)

```
module load singularity

cd /rsgrps/shaneburgess/amanda/i5k/GOanna

singularity pull docker://agbase/goanna:2.0

singularity run \
goanna_2.0.sif \
-h
```

See *Help and Usage Statement*

---

**Tip:** There are 3 directories built into this container. These directories should be used to mount data.

- /agbase\_database
- /go\_info
- /work-dir

---

GOanna has three required parameters:

```
-a BLAST database basename (acceptable options are listed in the help/usage)
-c peptide FASTA file to BLAST
-o output file basename
```

### Example PBS Script

```
#!/bin/bash
#PBS -N goanna
#PBS -W group_list=fionamcc
#PBS -l select=1:ncpus=28:mem=168gb
#PBS -q standard
#PBS -l walltime=6:0:0
#PBS -l cput=168:0:0

module load singularity

cd /rsgrps/shaneburgess/amanda/i5k/GOanna

singularity pull docker://agbase/goanna:2.0

singularity run \
-B /rsgrps/shaneburgess/amanda/i5k/agbase_database:/agbase_database \
-B /rsgrps/shaneburgess/amanda/i5k/go_info:/go_info \
-B /rsgrps/shaneburgess/amanda/i5k/goanna:/work-dir \
goanna_2.0.sif \
-a invertebrates \
-c AROS_10.faa \
-o AROS_10_invert_exponly \
-p \
-g 70 \
-s 900 \
```

(continues on next page)

(continued from previous page)

```
-d RefSeq \
-u "Amanda Cooksey" \
-x 37344 \
-t 28 \
-q 70 \
-k 9
```

## Command Explained

**singularity run:** tells Singularity to run

**-B /rsgrps/shaneburgess/amanda/i5k/agbase\_database:/agbase\_database:** tells docker to mount the ‘agbase\_database’ directory I downloaded to the host machine to the ‘/agbase\_database’ directory within the container. The syntax for this is: <absolute path on host>:<absolute path in container>

**-B /rsgrps/shaneburgess/amanda/i5k/go\_info:/go\_info:** mounts ‘go\_info’ directory on host machine into ‘go\_info’ directory inside the container

**-B /rsgrps/shaneburgess/amanda/i5k/goanna:/work-dir:** mounts my current working directory on the host machine to ‘/work-dir’ in the container

**goanna\_2.0.sif:** the name of the Singularity image file to use

---

**Tip:** All the options supplied after the image name are GOanna options

---

**-a invertebrates:** GOanna BLAST database to use—first of three required options.

**-c AROS\_10.faa:** input file (peptide FASTA)—second of three required options

**-o AROS\_10\_invert\_exponly:** output file basename—last of three required options

**-p:** our input file has NCBI defines. This specifies how to parse them.

**-g 70:** tells GOanna to keep only those matches with at least 70% identity

**-s 900:** tells GOanna to keep only those matches with a bitscore above 900

**-d RefSeq:** database of query ID. This will appear in column 1 of the GAF output file.

**-u “Amanda Cooksey”:** name to appear in column 15 of the GAF output file

**-x 37344:** NCBI taxon ID of input file species will appear in column 13 of the GAF output file

**-t 28:** number of threads to use for BLAST. This was run on a node with 28 cores.

**-k 9:** tells GOanna to keep only those matches with a maximum number of 9 gap openings

**-q 70:** tells GOanna to keep only those matches with query coverage of 70 per subject

## Understanding Your Results

If all goes well, you should get 4 output files:

**<basename>.asn:** This is standard BLAST output format that allows for conversion to other formats. You probably won’t need to look at this output.

**<basename>.html:** This output displays in your web browser so that you can view pairwise alignments to determine BLAST parameters.

**<basename>.tsv:** This is the tab-delimited BLAST output that can be opened and sorted in Excel to determine BLAST parameter values. The file contains the following columns:

- Query ID
- query length
- query start
- query end
- subject ID
- subject length
- subject start
- subject end
- e-value
- percent ID
- query coverage
- percent positive ID
- gap openings
- total gaps
- bitscore
- raw score

For more information on the BLAST output parameters see the [NCBI BLAST documentation](#).

**<basename>\_goanna\_gaf.tsv:** This is the standard tab-separated [GO annotation file format](#) that is used by the GO Consortium and by software tools that accept GO annotation files to do GO enrichment.

If you see more files in your output folder there may have been an error in the analysis or there may have been no GO to transfer. [Contact us](#).

## 2.5 GOanna on the ARS Ceres HPC

### 2.5.1 About Ceres/Scinet

- The Scinet VRSC has installed GOanna for ARS use.
- For general information on Scinet/Ceres, how to access it, and how to use it, visit <https://usda-ars-gbru.github.io/scinet-site/>.

### 2.5.2 Accessing the databases

GOanna requires access to some public databases that are already available on Ceres. These need to be in your working directory when you run the program. The best way to set this up is to create symbolic links to the databases from your working directory.

- 1) agbase\_database: species subset to run BLAST against

```
ln -s /reference/data/iplant/2019-09-16/agbase_database
```

2) go\_info: Uniprot GO annotations

```
ln -s /reference/data/iplant/2019-09-16/go_info
```

## 2.5.3 Running GOanna on Ceres

### Running programs on Ceres/Scinet

- You'll need to run GOanna either in interactive mode or batch mode.
- For interactive mode, use the *salloc* command.
- For batch mode, you'll need to write a batch job submission bash script.

### Running GOanna in interactive mode

#### Loading the module

The Scinet VRSC has installed the GOanna module. To load the module in interactive mode, run the command

```
module load agbase
```

#### Getting the Help and Usage Statement

```
goanna -h
```

See *Help and Usage Statement*

GOanna has three required parameters:

```
-a BLAST database basename (acceptable options are listed in the help/usage)
-c peptide FASTA file to BLAST
-o output file basename
```

#### Example Command

```
goanna -c AROS_10.faa -a invertebrates -o goanna_output -p -g 70 -s 900 -d RefSeq -u_
↪Monica -x 37344
```

### Running GOanna in batch mode

#### Running programs on Ceres/Scinet in batch mode

- Before using batch mode, you should review Scinet/Ceres' documentation first, and decide what queue you'll want to use. See <https://usda-ars-gbru.github.io/scinet-site/guide/ceres/>.

Example batch job submission bash script (e.g. *goanna-job.sh*):

```
#!/bin/bash
module load agbase
goanna -c AROS_10.faa -a invertebrates -o goanna_output -p -g 70 -s 900 -d RefSeq -u_
↵Monica -x 37344
```

### Submitting the batch job:

```
sbatch goanna-job.sh
```

## GOanna Commands Explained

- a invertebrates:** GOanna BLAST database to use—first of three required options.
- c AROS\_10.faa:** input file (peptide FASTA)—second of three required options
- o goanna\_output:** output file basename—last of three required options
- p:** our input file has NCBI defines. This specifies how to parse them.
- g 70:** tells GOanna to keep only those matches with at least 70% identity
- s 900:** tells GOanna to keep only those matches with a bitscore above 900
- d RefSeq:** database of query ID. This will appear in column 1 of the GAF output file.
- u “Monica”:** name to appear in column 15 of the GAF output file
- x 37344:** NCBI taxon ID of input file species will appear in column 13 of the GAF output file

## Understanding Your Results

If all goes well, you should get 4 output files:

**<basename>.asn:** This is standard BLAST output format that allows for conversion to other formats. You probably won't need to look at this output.

**<basename>.html:** This output displays in your web browser so that you can view pairwise alignments to determine BLAST parameters.

**<basename>.tsv:** This is the tab-delimited BLAST output that can be opened and sorted in Excel to determine BLAST parameter values. The file contains the following columns:

- Query ID
- query length
- query start<
- query end
- subject ID
- subject length
- subject start
- subject end
- e-value
- percent ID
- query coverage



- percent positive ID
- gap openings
- total gaps
- bitscore
- raw score

For more information on the BLAST output parameters see the [NCBI BLAST documentation](#).

**<basename>\_goanna\_gaf.tsv**: This is the standard tab-separated [GO annotation file format](#) that is used by the GO Consortium and by software tools that accept GO annotation files to do GO enrichment.

If you see more files in your output folder there may have been an error in the analysis or there may have been no GO to transfer. [Contact us](#).

## 2.6 Intro

[InterPro](#) is a database which integrates together predictive information about proteins' function from a number of partner resources, giving an overview of the families that a protein belongs to and the domains and sites it contains.

### Basic functions of this tool

- removes special characters from FASTA sequences
- splits FASTA into groups of 1000 sequences
- runs InterProScan with user-specified options on each of the 1000-sequence files in parallel
- re-combines output files from all groups of 1000
- parses the XML output from InterProScan to generate a gene association file (GAF) (and several other files)

Results and analysis from the application of InterProScan annotation to the Official gene set v3.0 protein set from *Diaphorina citri* followed by a differential expression analysis was presented at a seminar in the University of Arizona Animal and Comparative Biomedical Sciences in Fall 2020. The [slides](#) and [video](#) are available online.

---

**Note:** This tool accepts a peptide FASTA file. For those users with nucleotide sequences some documentation has been provided for using **TransDecoder** (although other tools are also acceptable). The [TransDecoder app](#) is available through CyVerse or as a [BioContainer](#) for use on the command line.

---

---

**Note:** As both GOanna and InterProScan provide GO annotations, their outputs are provided in GAF format. The '**Combine GAFs**' tool can then be used to make a single GAF of GO annotations, if desired.

---

### 2.6.1 Where to Find InterProScan

[Docker Hub](#) (5.36-75, 5.41-78 and 5.45-80)

[InterProScan](#) 5.36-75

## 2.6.2 Getting the InterProScan Data

Partner data are available from the InterProScan FTP site. These data are available as two separate downloads and can be obtained following these instructions:

### 1. Partner Data excluding Panther

```
wget ftp://ftp.ebi.ac.uk/pub/software/unix/iprscan/5/5.41-78.0/alt/interproscan-data-
↪5.41-78.0.tar.gz
wget ftp://ftp.ebi.ac.uk/pub/software/unix/iprscan/5/5.41-78.0/alt/interproscan-data-
↪5.41-78.0.tar.gz.md5
md5sum -c interproscan-data-5.41-78.0.tar.gz.md5
tar -pxvzf interproscan-data-5.41-78.0.tar.gz
```

#### tar options

- p = preserve the file permissions
- x = extract files from an archive
- v = verbosely list the files processed
- z = filter the archive through gzip
- f = use archive file

### 2. Getting the Panther data

```
cd interproscan-5.36-75.0/data
wget ftp://ftp.ebi.ac.uk/pub/software/unix/iprscan/5/data/panther-data-14.1.tar.gz
wget ftp://ftp.ebi.ac.uk/pub/software/unix/iprscan/5/data/panther-data-14.1.tar.gz.md5
md5sum -c panther-data-14.1.tar.gz.md5
tar -pxvzf panther-data-14.1.tar.gz
```

These data will be unarchived into this directory structure. They will need to be mounted to the container at run time.

```
<current_working_directory>/interproscan-5.41-78.0/data
```

## 2.6.3 Help and Usage Statement

```
Options:
-a <ANALYSES>                               Optional, comma separated list of analyses.
↪ If this option                               is not set, ALL analyses will be run.

-b <OUTPUT-FILE-BASE>                         Optional, base output filename (relative or ↪
↪absolute path).                               Note that this option, the output ↪
↪directory (-d) option and                    the output file name (-o) option are ↪
↪mutually exclusive. The                     appropriate file extension for the output ↪
↪format(s) will be                            appended automatically. By default the ↪
↪input file                                   path/name will be used.
```

(continues on next page)

(continued from previous page)

<p>-d &lt;OUTPUT-DIR&gt;  ↳option, the  ↳output file base (-b) option  ↳filename(s) are the  ↳appropriate file  ↳appended automatically .</p> <p>-c  ↳precalculated match lookup  ↳run locally.</p> <p>-C</p> <p>-e  ↳JSON output</p> <p>-f &lt;OUTPUT-FORMATS&gt;  ↳separated list of output  ↳JSON, GFF3, HTML and  ↳ XML and  ↳XML.</p> <p>-g  ↳corresponding Gene Ontology</p> <p>-h</p> <p>-i &lt;INPUT-FILE-PATH&gt;  ↳be loaded on  ↳mode, the</p> <p>-l  ↳InterPro  ↳formats.</p> <p>-m &lt;MINIMUM-SIZE&gt;  ↳to report. Will  ↳sequence type.  ↳specify a too  ↳takes a very long</p>	<p>Optional, output directory. Note that this  output file name (-o) option and the  are mutually exclusive. The output  same as the input filename, with the  extension(s) for the output format(s).</p> <p>Optional. Disables use of the  service. All match calculations will be</p> <p>Optional. Supply the number of cpus to use.</p> <p>Optional, excludes sites from the XML,</p> <p>Optional, case-insensitive, comma  formats. Supported formats are TSV, XML,  SVG. Default for protein sequences are TSV,  GFF3, or for nucleotide sequences GFF3 and</p> <p>Optional, switch on lookup of  annotation (IMPLIES -l lookup option)</p> <p>Optional, display help information</p> <p>Optional, path to fasta file that should  Master startup. Alternatively, in CONVERT  InterProScan 5 XML file to convert.</p> <p>Also include lookup of corresponding  annotation in the TSV and GFF3 output</p> <p>Optional, minimum nucleotide size of ORF  only be considered if n is specified as a  Please be aware of the fact that if you  short value it might be that the analysis  time!</p>
---	---

(continues on next page)

(continued from previous page)

```

-o <EXPLICIT_OUTPUT_FILENAME>           Optional explicit output file name
↳(relative or absolute                    path). Note that this option, the output
↳directory -d option                       and the output file basename -b option are
↳mutually                                   exclusive. If this option is given, you
↳MUST specify a                             single output format using the -f option.
↳The output file                             name will not be modified. Note that
↳specifying an output                       file name using this option OVERWRITES ANY
↳EXISTING FILE.

-p                                         Optional, switch on lookup of
↳corresponding Pathway                     annotation (IMPLIES -l lookup option)

-t <SEQUENCE-TYPE>                       Optional, the type of the input sequences
↳(dna/rna (n))                              or protein (p)). The default sequence
↳type is protein.

-T <TEMP-DIR>                             Optional, specify temporary file directory
↳(relative or                               absolute path). The default location is
↳temp/.

-v                                         Optional, display version number

-r                                         Optional. 'Mode' required ( -r 'cluster')
↳to run in cluster mode. These options     are provided but have not been tested with
↳this wrapper script. For                 more information on running InterProScan
↳in cluster mode:                          https://github.com/ebi-pf-team/
↳interproscan/wiki/ClusterMode

-R                                         Optional. Clusterrunid (cruid) required
↳when using cluster mode.                 -R unique_id

Available analyses:
      TIGRFAM (XX.X) : TIGRFAMS are protein families based on Hidden
↳Markov Models or HMMs
      SFLD (X.X) : SFLDs are protein families based on Hidden
↳Markov Models or HMMs
      ProDom (XXXX.X) : ProDom is a comprehensive set of protein
↳domain families automatically generated from the UniProt Knowledge Database.
      Hamap (XXXXXX.XX) : High-quality Automated and Manual
↳Annotation of Microbial Proteomes
      SMART (X.X) : SMART allows the identification and analysis of
↳domain architectures based on Hidden Markov Models or HMMs
      CDD (X.XX) : Prediction of CDD domains in Proteins
      ProSiteProfiles (XX.XXX) : PROSITE consists of documentation entries
↳describing protein domains, families and functional sites as well as associated
↳patterns and profiles to identify them

```

(continues on next page)

(continued from previous page)

```

ProSitePatterns (XX.XXX) : PROSITE consists of documentation entries
↳describing protein domains, families and functional sites as well as associated
↳patterns and profiles to identify them
SUPERFAMILY (X.XX) : SUPERFAMILY is a database of structural and
↳functional annotation for all proteins and genomes.
PRINTS (XX.X) : A fingerprint is a group of conserved motifs
↳used to characterise a protein family
PANTHER (X.X) : The PANTHER (Protein ANalysis THrough
↳Evolutionary Relationships) Classification System is a unique resource that
↳classifies genes by their functions, using published scientific experimental
↳evidence and evolutionary relationships to predict fu$
Gene3D (X.X.X) : Structural assignment for whole genes and
↳genomes using the CATH domain structure database
PIRSF (X.XX) : The PIRSF concept is being used as a guiding
↳principle to provide comprehensive and non-overlapping clustering of UniProtKB
↳sequences into a hierarchical order to reflect their evolutionary relationships.
Pfam (XX.X) : A large collection of protein families, each
↳represented by multiple sequence alignments and hidden Markov models (HMMs)
Coils (X.X) : Prediction of Coiled Coil Regions in Proteins
MobiDBLite (X.X) : Prediction of disordered domains Regions in
↳Proteins

```

#### OPTIONS FOR XML PARSER OUTPUTS

```

-F <IPRS output directory>          This is the output directory from
↳InterProScan.
-D <database>                       Supply the database responsible for these
↳annotations.
-x <taxon>                          NCBI taxon ID of the ID being annotated
-y <type>                            Transcript or protein
-n <biocurator>                     Name of the biocurator who made these
↳annotations
-M <mapping file>                   Optional. Mapping file.
-B <bad seq file>                   Optional. Bad input sequence file.

```

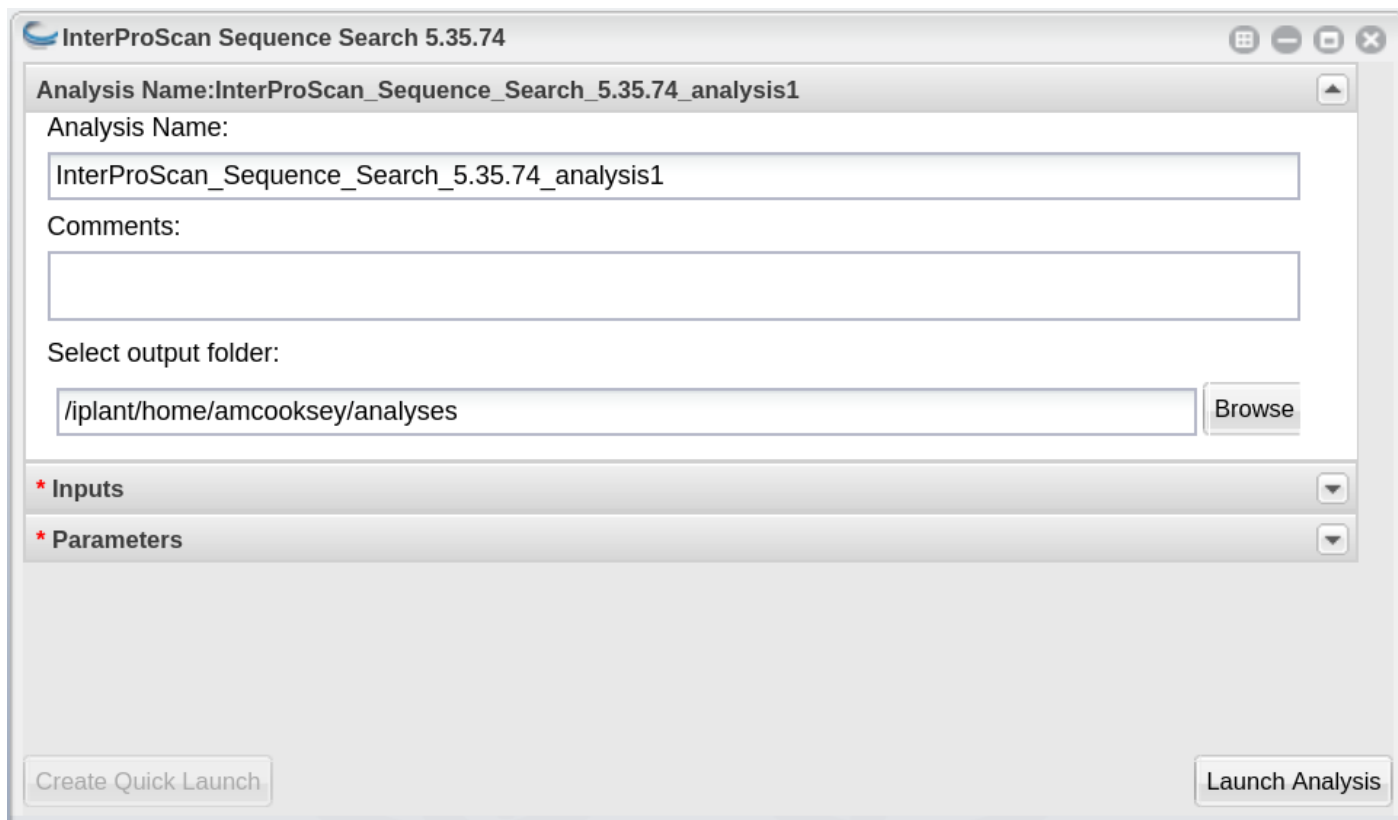
## 2.7 InterProScan on CyVerse

### 2.7.1 Accessing InterProScan in the Discovery Environment

1. Create an account on CyVerse (free)
2. Open the CyVerse Discovery Environment (DE) and login with your CyVerse credentials.
3. If you are new to the Discovery Environment (DE) the user guide can be found [here](#).
4. Click on the 'Data' button at the left side of the screen to access your files/folders. Upload your data to the DE.
5. To access the [InterProScan Sequence Search 5.36-75.0](#) app click on the 'Apps' button at the left side of the DE.
6. Search for 'interproscan' in the search bar at the top of the 'apps' window. The contents of the folder will appear in the main pane of the window. The InterProScan app is called 'InterProScan Sequence Search 5.36-75'; click on the name to open the app.

### 2.7.2 Using the InterProScan App

## Launching the App



InterProScan Sequence Search 5.35.74

Analysis Name: InterProScan\_Sequence\_Search\_5.35.74\_analysis1

Analysis Name:  
InterProScan\_Sequence\_Search\_5.35.74\_analysis1

Comments:

Select output folder:  
/iplant/home/amcooksey/analyses

\* Inputs

\* Parameters

**InterProScan\_Sequence\_Search\_5.36.75\_analysis1:** This menu is used to name the job you will run so that you can find it later. Analysis Name: The default name is “InterProScan\_Sequence\_Search\_5.36.75\_analysis1”. We recommend changing the ‘analysis1’ portion of this to reflect the data you are running.

**Comments:** (Optional) You can add additional information in the comments section to distinguish your analyses further.

**Select output folder:** This is where your results will be placed. The default (recommended) is your ‘analyses’ folder.

**Retain Inputs:** Enabling this flag will copy all the input files into the analysis result folder.

**Warning:** Selecting this option will rapidly consume your allocated space. It is not recommended. Your inputs will always remain available in the folder in which you stored them.

## Input

**Peptide FASTA file:** Use the Browse button on the right hand side to navigate to your Data folder and select your protein sequence file.

## Parameters

**Annotate each peptide with Gene Ontology information:** Be sure this box is checked. This will ensure that you get GO annotations

**Biocurator:** This will be used to fill the ‘assigned by’ field of your GAF output file. If you do not fill it in the default “user” will be used instead.

**Database:** Use the database that sequences were obtained from (Genbank), or a recognizable project name if these sequences are not in a database (e.g., i5k project or Smith Lab).

**Annotate each peptide with biological pathway information:** This is optional. However, if you want pathways annotations be it is checked.

**Taxon:** Enter the NCBI taxon number for your species. This can be found by searching for your species name (common or scientific) in the [NCBI taxon database](#).

**InterProScan output directory name:** This will be the name of the folder for your output files. The default folder name is ‘outdir’.

## 2.7.3 Understanding Your Results

### InterProScan Outputs

This app provides all six of the InterProScan output formats. For more details on the contents of each file please refer to the InterProScan [outputs documentation](#).

**<basename>.gff3**

**<basename>.tsv**

**<basename>.xml**

**<basename>.json**

**<basename>.html.tar.gz**

**<basename>.svg.tar.gz**

This app also runs the ‘InterProScan Results Function’ on the XML output from InterProScan. This tool provides a GAF output file and a variety of summary (count) files described below.

### InterProScan Results Function Outputs

**<basename>\_gaf.txt:** -This table follows the formatting of a gene association file (gaf) and can be used in GO enrichment analyses.

**<basename>\_acc\_go\_counts.txt:** -This table includes input accessions, the number of GO IDs assigned to each accession and GO ID names. GO IDs are split into BP (Biological Process), MF (Molecular Function) and CC (Cellular Component).

**<basename>\_go\_counts.txt:** -This table counts the numbers of sequences assigned to each GO ID so that the user can quickly identify all genes assigned to a particular function.

**<basename>\_acc\_interpro\_counts.txt:** -This table includes input accessions, number of InterPro IDs for each accession, InterPro IDs assigned to each sequence and the InterPro ID name.

**<basename>\_interpro\_counts.txt:** -This table counts the numbers of sequences assigned to each InterPro ID so that the user can quickly identify all genes with a particular motif.

**<basename>\_acc\_pathway\_counts.txt:** -This table includes input accessions, number of pathway IDs for the accession and the pathway names. Multiple values are separated by a semi-colon.

**<basename>\_pathway\_counts.txt:** -This table counts the numbers of sequences assigned to each Pathway ID so that the user can quickly identify all genes assigned to a pathway.

**<basename>.err:** -This file will list any sequences that were not able to be analyzed by InterProScan. Examples of sequences that will cause an error are sequences with a large run of Xs.

If your output doesn't look like you expect please check the 'condor\_stderr' file in the analysis output 'logs' folder. If that doesn't clarify the problem contact us at [agbase@email.arizona.edu](mailto:agbase@email.arizona.edu) or [support@cyverse.org](mailto:support@cyverse.org).

## 2.8 InterProScan on the Command Line

### 2.8.1 Getting the InterProScan Data

Partner data are available from the InterProScan FTP site. These data are available as two separate downloads and can be obtained following these instructions:

#### 1. Partner Data excluding Panther

```
wget ftp://ftp.ebi.ac.uk/pub/software/unix/iprscan/5/5.45-80.0/alt/interproscan-data-5.45-80.0.tar.gz
wget ftp://ftp.ebi.ac.uk/pub/software/unix/iprscan/5/5.45-80.0/alt/interproscan-data-5.45-80.0.tar.gz.md5
md5sum -c interproscan-data-5.45-80.0.tar.gz.md5
tar -pxvzf interproscan-data-5.45-80.0.tar.gz
```

---

#### tar options

- p = preserve the file permissions
- x = extract files from an archive
- v = verbosely list the files processed
- z = filter the archive through gzip
- f = use archive file

---

#### 2. Getting the Panther data

```
cd interproscan-5.45-80.0/data
wget ftp://ftp.ebi.ac.uk/pub/software/unix/iprscan/5/data/panther-data-14.1.tar.gz
wget ftp://ftp.ebi.ac.uk/pub/software/unix/iprscan/5/data/panther-data-14.1.tar.gz.md5
md5sum -c panther-data-14.1.tar.gz.md5
tar -pxvzf panther-data-14.1.tar.gz
```

These data will be unarchived into this directory structure. They will need to be mounted to the container at run time.

```
<current_working_directory>/interproscan-5.45-80.0/data
```

### 2.8.2 Container Technologies

Interproscan is provided as a Docker container.

A container is a standard unit of software that packages up code and all its dependencies so the application runs quickly and reliably from one computing environment to another.

There are two major containerization technologies: **Docker** and **Singularity**.

Docker containers can be run with either technology.



## 2.8.3 Running InterProScan using Docker

---

### About Docker

- Docker must be installed on the computer you wish to use for your analysis.
  - To run Docker you must have ‘root’ permissions (or use sudo).
  - Docker will run all containers as ‘root’. This makes Docker incompatible with HPC systems (see Singularity below).
  - Docker can be run on your local computer, a server, a cloud virtual machine (such as CyVerse Atmosphere) etc. Docker can be installed quickly on an Atmosphere instance by typing ‘ezd’.
  - For more information on installing Docker on other systems see this tutorial: [Installing Docker on your machine](#).
- 

**Important:** We have included this basic documentation for running InterProScan with Docker. However, InterProScan requires quite a lot of compute resources and may need to be run on an HPC system. If you need to use HPC see ‘Singularity’ below.

---

### Getting the InterProScan Container

The InterProScan tool is available as a Docker container on Docker Hub where you can see all the available versions: [InterProScan container](#)

The latest container can be pulled with this command:

```
docker pull agbase/interproscan:5.45-80.0_1
```

---

### Remember

You must have root permissions or use sudo, like so:

```
sudo docker pull agbase/interproscan:5.45-80.0_1
```

---

### Running InterProScan with Data

**Tip:** There is one directory built into this container. This directory should be used to mount your working directory.

- /data
- 

### Getting the Help and Usage Statement

```
sudo docker run --rm -v $(pwd):/work-dir agbase/interproscan:5.45-80.0_1 -h
```

---

See iprsusage

### Example Command

```
sudo docker run \  
-v /rsgrps/shaneburgess/amanda/i5k/interproscan:/data \  
-i /rsgrps/shaneburgess/amanda/i5k/interproscan/pnnl_10000.fasta \  
-v /rsgrps/shaneburgess/amanda/i5k/interproscan/interproscan-5.45-80.0/data:/opt/  
↪interproscan/data \  
agbase/interproscan:5.45-80.0_1 \  
-d outdir_10000 \  
-f tsv,json,xml,html,gff3,svg \  
-g \  
-p \  
-c \  
-n Amanda \  
-x 109069 \  
-D AgBase \  
-l
```

### Command Explained

**sudo docker run:** tells docker to run

**-rm:** removes container when analysis finishes (image will remain for future analyses)

**-v /rsgrps/shaneburgess/amanda/i5k/interproscan:/data:** mount my working directory on the host machine into the /data directory in the container. The syntax for this is <absolute path on host machine>:<absolute path in container>

**-v /rsgrps/shaneburgess/amanda/i5k/interproscan/interproscan-5.45-80.0/data:/opt/interproscan/data:** mounts the InterProScan partner data (downloaded from FTP) on the host machine into the /opt/interproscan/data directory in the container

**agbase/interproscan:5.45-80.0\_1:** the name of the Docker image to use

---

**Tip:** All the options supplied after the image name are Interproscan options

---

**-i /rsgrps/shaneburgess/amanda/i5k/interproscan/pnnl\_10000.fasta:** local path to input FASTA file. You can also use the mounted file path: /data/pnnl\_10000.fasta

**-d outdir\_10000:** output directory name

**-f tsv,json,xml,html,gff3,svg:** desired output file formats

**-g:** tells the tool to perform GO annotation

**-p:** tells tool to perform pathway annotation

**-c:** tells tool to perform local compute and not connect to EBI. This only adds a little to the run time but removes error messages from network time out errors

**-n Amanda:** name of biocurator to include in column 15 of GAF output file

**-x 109069:** taxon ID of query species to be used in column 13 of GAF output file

**-D AgBase:** database of query accession to be used in column 1 of GAF output file

**-l:** tells tools to include lookup of corresponding InterPro annotation in the TSV and GFF3 output formats.

## Understanding Your Results

InterProScan outputs: <https://github.com/ebi-pf-team/interproscan/wiki/OutputFormats>

**Default** - <basename>.gff3 - <basename>.tsv - <basename>.xml

**Optional** - <basename>.json - <basename>.html.tar.gz - <basename>.svg.tar.gz

### Parser Outputs

**<basename>\_gaf.txt:** -This table follows the formatting of a gene association file (gaf) and can be used in GO enrichment analyses.

**<basename>\_acc\_go\_counts.txt:** -This table includes input accessions, the number of GO IDs assigned to each accession and GO ID names. GO IDs are split into BP (Biological Process), MF (Molecular Function) and CC (Cellular Component).

**<basename>\_go\_counts.txt:** -This table counts the numbers of sequences assigned to each GO ID so that the user can quickly identify all genes assigned to a particular function.

**<basename>\_acc\_interpro\_counts.txt:** -This table includes input accessions, number of InterPro IDs for each accession, InterPro IDs assigned to each sequence and the InterPro ID name.

**<basename>\_interpro\_counts.txt:** -This table counts the numbers of sequences assigned to each InterPro ID so that the user can quickly identify all genes with a particular motif.

**<basename>\_acc\_pathway\_counts.txt:** -This table includes input accessions, number of pathway IDs for the accession and the pathway names. Multiple values are separated by a semi-colon.

**<basename>\_pathway\_counts.txt:** -This table counts the numbers of sequences assigned to each Pathway ID so that the user can quickly identify all genes assigned to a pathway.

**<basename>.err:** -This file will list any sequences that were not able to be analyzed by InterProScan. Examples of sequences that will cause an error are sequences with a large run of Xs.

If you see more files in your output folder there may have been an error in the analysis or there may have been no GO to transfer. [Contact us](#).

## 2.8.4 Running InterProScan with Singularity (HPC)

---

### About Singularity

- does not require 'root' permissions
- runs all containers as the user that is logged into the host machine
- HPC systems are likely to have Singularity installed and are unlikely to object if asked to install it (no guarantees).
- can be run on any machine where it is installed
- more information about [installing Singularity](#)
- This tool was tested using Singularity 3.0. Users with Singularity 2.x will need to modify the commands accordingly.

### HPC Job Schedulers

Although Singularity can be installed on any computer this documentation assumes it will be run on an HPC system. The tool was tested on a PBSPro system and the job submission scripts below reflect that. Submission scripts will need to be modified for use with other job scheduler systems.

---

### Getting the InterProScan Data

Partner data are available from the InterProScan FTP site. These data are available as two separate downloads and can be obtained following these instructions:

#### 1. Partner Data excluding Panther

```
wget ftp://ftp.ebi.ac.uk/pub/software/unix/iprscan/5/5.45-80.0/alt/interproscan-data-5.45-80.0.tar.gz
wget ftp://ftp.ebi.ac.uk/pub/software/unix/iprscan/5/5.45-80.0/alt/interproscan-data-5.45-80.0.tar.gz.md5
md5sum -c interproscan-data-5.45-80.0.tar.gz.md5
tar -pxvzf interproscan-data-5.45-80.0.tar.gz
```

#### tar options

- p = preserve the file permissions
  - x = extract files from an archive
  - v = verbosely list the files processed
  - z = filter the archive through gzip
  - f = use archive file
- 

#### 2. Getting the Panther data

```
cd interproscan-5.45-80.0/data
wget ftp://ftp.ebi.ac.uk/pub/software/unix/iprscan/5/data/panther-data-14.1.tar.gz
wget ftp://ftp.ebi.ac.uk/pub/software/unix/iprscan/5/data/panther-data-14.1.tar.gz.md5
md5sum -c panther-data-14.1.tar.gz.md5
tar -pxvzf panther-data-14.1.tar.gz
```

These data will be unarchived into this directory structure. They will need to be mounted to the container at run time.

```
<current_working_directory>/interproscan-5.45-80.0/data
```

### Getting the InterProScan Container

The InterProScan tool is available as a Docker container on Docker Hub: [InterProScan container](#)

The container can be pulled with this command:

```
singularity pull docker://agbase/interproscan:5.45-80.0_1
```

## Running InterProScan with Data

### Getting the Help and Usage Statement

#### Example PBS script:

```
#!/bin/bash
#PBS -N 10000j100
#PBS -q standard
#PBS -l select=1:ncpus=28:mem=168gb
#PBS -W group_list=fionamcc
#PBS -l walltime=6:0:0
#PBS -l cput=168:0:0

module load singularity

cd /rsgrps/shaneburgess/amanda/i5k/interproscan

singularity pull docker://agbase/interproscan:5.45-80.0_1

singularity run \
interproscan_5.45-80.0_1.sif \
-h
```

See `iprsusage`

**Tip:** There is one directory built into this container. This directory should be used to mount your working directory.

- `/data`

#### Example PBS Script

```
#!/bin/bash
#PBS -N 10000j100
#PBS -q standard
#PBS -l select=1:ncpus=28:mem=168gb
#PBS -W group_list=fionamcc
#PBS -l walltime=6:0:0
#PBS -l cput=168:0:0

module load singularity

cd /rsgrps/shaneburgess/amanda/i5k/interproscan

singularity pull docker://agbase/interproscan:5.45-80.0_1

singularity run \
-B /rsgrps/shaneburgess/amanda/i5k/interproscan:/data \
-B /rsgrps/shaneburgess/amanda/i5k/interproscan/interproscan-5.45-80.0/data:/opt/
↪interproscan/data \
interproscan_5.45-80.0_1.sif \
-i /rsgrps/shaneburgess/amanda/i5k/interproscan/pnnl_10000.fasta \
-d outdir_10000 \
-f tsv,json,xml,html,gff3,svg \
```

(continues on next page)

(continued from previous page)

```
-g \  
-p \  
-c \  
-n Amanda \  
-x 109069 \  
-D AgBase \  
-l
```

## Command Explained

**singularity run:** tells Singularity to run

**-B /rsgrps/shaneburgess/amanda/i5k/interproscan:/data:** mounts my working directory on the host machine into the /data directory in the container the syntax for this is <absolute path on host machine>:<absolute path in container>

**-B /rsgrps/shaneburgess/amanda/i5k/interproscan/interproscan-5.45-80.0/data:/opt/interproscan/data:** mounts the InterProScan data directory that was downloaded from the FTP site into the InterProScan data directory in the container

**interproscan\_5.45-80.0\_1.sif:** name of the image to use

---

**Tip:** All the options supplied after the image name are options for this tool

---

**-i /rsgrps/shaneburgess/amanda/i5k/interproscan/pnml\_10000.fasta:** input FASTA file

**-d outdir\_10000:** output directory name

**-f tsv,json,xml,html,gff3,svg:** desired output file formats

**-g:** tells the tool to perform GO annotation

**-c:** tells tool to perform local compute and not connect to EBI. This only adds a little to the run time but removes error messages from network time out errors

**-p:** tells tool to perform pathway annoation

**-n Amanda:** name of biocurator to include in column 15 of GAF output file

**-x 109069:** taxon ID of query species to be used in column 13 of GAF output file

**-D AgBase:** database of query accession to be used in column 1 of GAF output file

**-l:** tells tools to include lookup of corresponding InterPro annotation in the TSV and GFF3 output formats.

## Understanding Your Results

**InterProScan outputs:** <https://github.com/ebi-pf-team/interproscan/wiki/OutputFormats>

**Default** - <basename>.gff3 - <basename>.tsv - <basename>.xml

**Optional** - <basename>.json - <basename>.html.tar.gz - <basename>.svg.tar.gz

## Parser Outputs

**<basename>\_gaf.txt:** -This table follows the formatting of a gene association file (gaf) and can be used in GO enrichment analyses.

**<basename>\_acc\_go\_counts.txt:** -This table includes input accessions, the number of GO IDs assigned to each accession and GO ID names. GO IDs are split into BP (Biological Process), MF (Molecular Function) and CC (Cellular Component).

**<basename>\_go\_counts.txt:** -This table counts the numbers of sequences assigned to each GO ID so that the user can quickly identify all genes assigned to a particular function.

**<basename>\_acc\_interpro\_counts.txt:** -This table includes input accessions, number of InterPro IDs for each accession, InterPro IDs assigned to each sequence and the InterPro ID name.

**<basename>\_interpro\_counts.txt:** -This table counts the numbers of sequences assigned to each InterPro ID so that the user can quickly identify all genes with a particular motif.

**<basename>\_acc\_pathway\_counts.txt:** -This table includes input accessions, number of pathway IDs for the accession and the pathway names. Multiple values are separated by a semi-colon.

**<basename>\_pathway\_counts.txt:** -This table counts the numbers of sequences assigned to each Pathway ID so that the user can quickly identify all genes assigned to a pathway.

**<basename>.err:** -This file will list any sequences that were not able to be analyzed by InterProScan. Examples of sequences that will cause an error are sequences with a large run of Xs.

If you see more files in your output folder there may have been an error in the analysis or there may have been no GO to transfer. [Contact us](#).

## 2.9 InterProScan on the ARS Ceres HPC

### 2.9.1 About Ceres/Scinet

- The Scinet VRSC has installed InterProScan for ARS use.
- For general information on Ceres/Scinet, how to access it, and how to use it, visit <https://usda-ars-gbru.github.io/scinet-site/>.

### 2.9.2 Running InterProScan on Ceres/Scinet

---

#### Important:

- InterProScan requires quite a lot of compute resources and should be run in batch mode.
  - Before using batch mode, you should review Ceres/Scinet's documentation first, and decide what queue you'll want to use. See <https://usda-ars-gbru.github.io/scinet-site/guide/ceres/>.
- 

#### Running InterProScan with Data

## Getting the Help and Usage Statement

```
agbase_interproscan -h
```

See iprsusage

## Example batch job submission bash script (e.g. agbase\_interproscan-job.sh):

```
#!/bin/bash
module load agbase
agbase_interproscan -i AROS_10.faa -d outdir -f tsv,json,xml,html,gff3,svg -g -p -n_  
↳Monica -x 109069 -D i5k
```

## Submitting the batch job:

```
sbatch agbase_interproscan-job.sh
```

## Command Explained

- i AROS\_10.faa:** local path to input FASTA file.
- d outdir:** output directory name
- f tsv,json,xml,html,gff3,svg:** desired output file formats
- g:** tells the tool to perform GO annotation
- p:** tells tool to perform pathway annoation
- n Monica:** name of biocurator to include in column 15 of GAF output file
- x 109069:** taxon ID of query species to be used in column 13 of GAF output file
- D i5k:** database of query accession to be used in column 1 of GAF output file

## Understanding Your Results

**InterProScan outputs:** <https://github.com/ebi-pf-team/interproscan/wiki/OutputFormats>

**Default** - <basename>.gff3 - <basename>.tsv - <basename>.xml

**Optional** - <basename>.json - <basename>.html.tar.gz - <basename>.svg.tar.gz

## Parser Outputs

**<basename>\_gaf.txt:** -This table follows the formatting of a gene association file (gaf) and can be used in GO enrichment analyses.

**<basename>\_acc\_go\_counts.txt:** -This table includes input accessions, the number of GO IDs assigned to each accession and GO ID names. GO IDs are split into BP (Biological Process), MF (Molecular Function) and CC (Cellular Component).



**<basename>\_go\_counts.txt:** -This table counts the numbers of sequences assigned to each GO ID so that the user can quickly identify all genes assigned to a particular function.

**<basename>\_acc\_interpro\_counts.txt:** -This table includes input accessions, number of InterPro IDs for each accession, InterPro IDs assigned to each sequence and the InterPro ID name.

**<basename>\_interpro\_counts.txt:** -This table counts the numbers of sequences assigned to each InterPro ID so that the user can quickly identify all genes with a particular motif.

**<basename>\_acc\_pathway\_counts.txt:** -This table includes input accessions, number of pathway IDs for the accession and the pathway names. Multiple values are separated by a semi-colon.

**<basename>\_pathway\_counts.txt:** -This table counts the numbers of sequences assigned to each Pathway ID so that the user can quickly identify all genes assigned to a pathway.

**<basename>.err:** -This file will list any sequences that were not able to be analyzed by InterProScan. Examples of sequences that will cause an error are sequences with a large run of Xs.

If you see more files in your output folder there may have been an error in the analysis or there may have been no GO to transfer. [Contact us](#).

## 2.10 Intro

This tool can be used to combine the gene association file (GAF) outputs from GOanna and InterProScan.

The tool accepts two input files:

1. GOanna GAF output
2. InterProScan GAF output

---

**Note:** InterProScan itself does not produce a GAF file. The [AgBase InterProScan container](#) parses the XML output from InterProScan to produce the GAF file.

---

### 2.10.1 Where to Find Combine GAFs

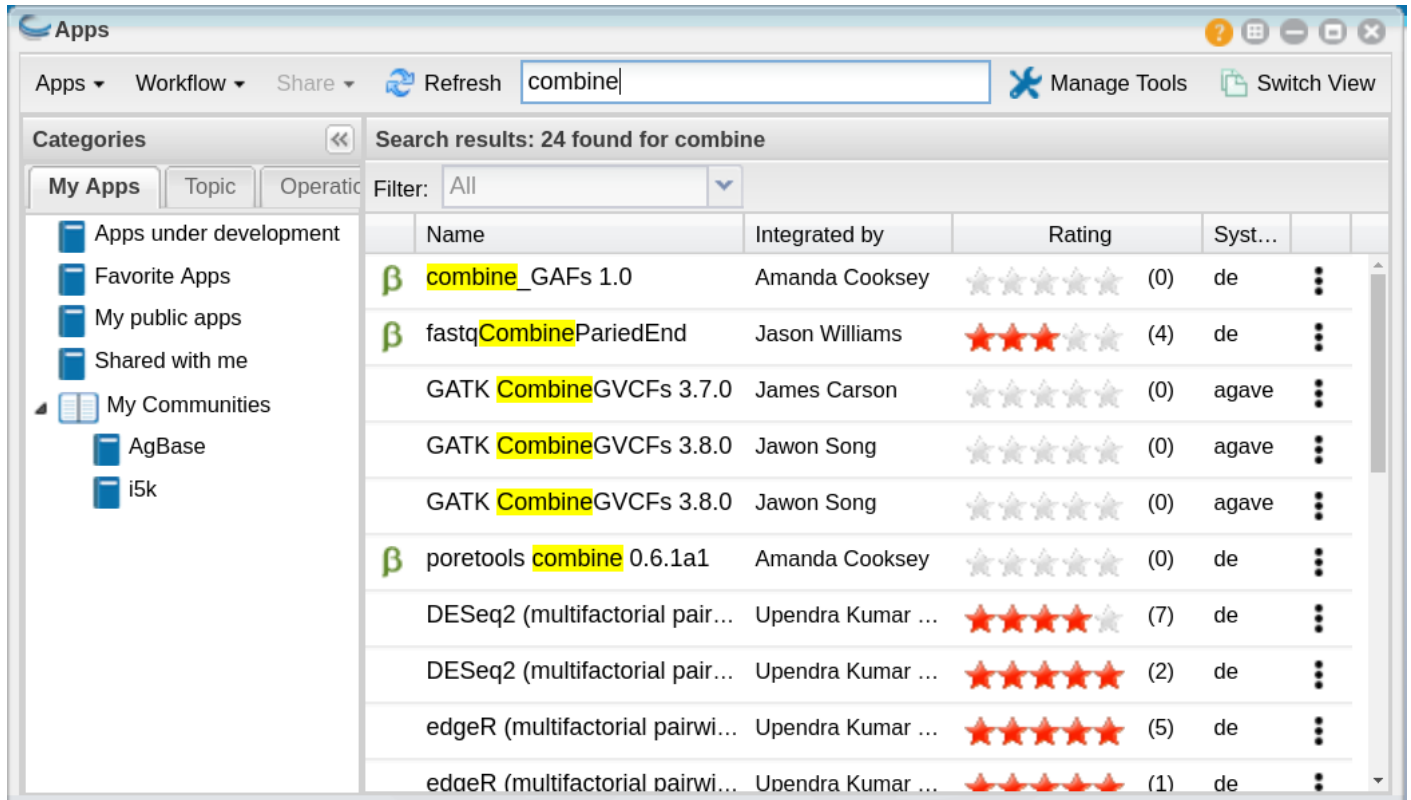
- [Docker Hub](#)
- [Cyverse Discovery Environment](#)

## 2.11 Combine GAFs on CyVerse

### 2.11.1 Accessing GOanna in the Discovery Environment

1. [Create an account on CyVerse](#) (free)
2. Open the CyVerse Discovery Environment (DE) and login with your CyVerse credentials.
3. If you are new to the Discovery Environment (DE) the user guide can be found [here](#).
4. Click on the 'Data' button at the left side of the screen to access your files/folders. Upload your data to the DE.
5. To access the [combine\\_GAFs 1.0](#) app click on the 'Apps' button at the left side of the DE.

- Search for ‘combine’ in the search bar at the top of the ‘apps’ window (see below). The contents of the folder will appear in the main pane of the window. The combine\_GAFs app is called ‘combine\_GAFs 1.0’; click on the name to open the app.



### Find Apps Easily with ‘Communities’

The GOanna 2.0 app belongs to the ‘i5k’ and ‘AgBase’ communities. You can join either of these communities and they will appear in the left-hand pane of your ‘Apps’ window (see above).

To join a community click on the person icon in the top-right corner of the Discovery Environment window and select ‘Communities’. In the ‘Communities’ window choose ‘all communities’ from the drop-down list. A list of communities will appear in the main pane of this window. Select the one you wish to join by clicking on it and then clicking on the ‘join’ button.

### Using the Combine\_GAFs App

## Launching the App

**Analysis Name: Combine\_GAFs\_1.0\_analysis1:** This menu is used to name the job you will run so that you can find it later. Analysis Name: The default name is “Combine\_GAFs\_1.0\_analysis1”. We recommend changing the ‘analysis1’ portion of this to reflect the data you are running.

**Comments:** (Optional) You can add additional information in the comments section to distinguish your analyses further.

**Select output folder:** This is where your results will be placed. The default (recommended) is your ‘analyses’ folder.

**Retain Inputs:** Enabling this flag will copy all the input files into the analysis result folder.

**Warning:** Selecting this option will rapidly consume your allocated space. It is not recommended. Your inputs will always remain available in the folder in which you stored them.

## Input

**GOanna GAF Output File:** This is the GAF file generated by a GOanna analysis.

**InterProScan XML Parser GAF Output File:** This is the GAF output file generated by an InterProScan XML Parser analysis. InterProScan itself does not produce this file, though some InterProScan apps include this analysis. If it is missing from your InterProScan output you can generate it using the InterProScan XML Parser app.

## Output

**Output File Basename:** This will be the prefix for your output file (a .tsv extension will be added).

If your analysis fails please check the 'condor\_stderr' file in the analysis output 'logs' folder. If that doesn't clarify the problem contact us at [agbase@email.arizona.edu](mailto:agbase@email.arizona.edu) or [support@cyverse.org](mailto:support@cyverse.org).

## 2.12 Combine GAFs on the Command Line

### 2.12.1 Container Technologies

GOanna is provided as a Docker container.

A container is a standard unit of software that packages up code and all its dependencies so the application runs quickly and reliably from one computing environment to another.

There are two major containerization technologies: **Docker** and **Singularity**.

Docker containers can be run with either technology.

### 2.12.2 Combine GAFs using Docker

---

#### About Docker

- Docker must be installed on the computer you wish to use for your analysis.
  - To run Docker you must have 'root' permissions (or use sudo).
  - Docker will run all containers as 'root'. This makes Docker incompatible with HPC systems (see Singularity below).
  - Docker can be run on your local computer, a server, a cloud virtual machine (such as CyVerse Atmosphere) etc. Docker can be installed quickly on an Atmosphere instance by typing 'ezd'.
  - For more information on installing Docker on other systems see this tutorial: [Installing Docker on your machine](#).
- 

#### Getting the Combine GAFs container

The Combine GAFs tool is available as a Docker container on Docker Hub: [Combine GAFs container](#)

The container can be pulled with this command:

```
docker pull agbase/combine_gafs:1.0
```

---

#### Remember

You must have root permissions or use sudo, like so:

```
sudo docker pull agbase/combine_gafs:1.0
```

---

## Running Combine GAFs with Data

Combine GAFs has three parameters:

```
-i InterProScan XML Parser GAF output
-g GOanna GAF output
-o output file basename
```

### Example Command

```
sudo docker run \
--rm \
-v $(pwd):/work-dir \
agbase/combine_gafs:1.0 \
-i CFLO_1.fa_gaf.txt \
-g clfo1_v_insecta_goanna_gaf.tsv \
-o complete_gaf
```

### Command Explained

**sudo docker run:** tells docker to run

**-rm:** removes the container when the analysis has finished. The image will remain for future use.

**-v \$(pwd):/work-dir:** mounts my current working directory on the host machine to '/work-dir' in the container

**agbase/combine\_gafs:1.0:** the name of the Docker image to use

---

**Tip:** All the options supplied after the image name are Combine\_GAFs options

---

**-i CFLO\_1.fa\_gaf.txt:** InterProScan XML Parser GAF output file.

**-g clfo1\_v\_insecta\_goanna\_gaf.tsv:** GOanna GAF output file.

**-o complete\_gaf:** output file basename—a .tsv extension will be added

## 2.12.3 Combine GAFs using Singularity

### About Singularity

- does not require 'root' permissions
- runs all containers as the user that is logged into the host machine
- HPC systems are likely to have Singularity installed and are unlikely to object if asked to install it (no guarantees).
- can be run on any machine where it is installed
- more information about [installing Singularity](#)
- This tool was tested using Singularity 3.0. Users with Singularity 2.x will need to modify the commands accordingly.

### HPC Job Schedulers

Although Singularity can be installed on any computer this documentation assumes it will be run on an HPC system. The tool was tested on a PBSPro system and the job submission scripts below reflect that. Submission scripts will need to be modified for use with other job scheduler systems.

---

### Getting the Combine GAFs Container

The Combine GAFs tool is available as a Docker container on Docker Hub: [Combine GAFs container](#)

The container can be pulled with this command:

```
singularity pull docker://agbase/combine_gafs:1.0
```

### Running Combine GAFs with Data

Combine GAFs has three parameters:

```
-i InterProScan XML Parser GAF output  
-g GOanna GAF output  
-o output file basename
```

### Example PBS Script

```
#!/bin/bash  
#PBS -N combine_gafs  
#PBS -W group_list=fionamcc  
#PBS -l select=1:ncpus=28:mem=168gb  
#PBS -q standard  
#PBS -l walltime=6:0:0  
#PBS -l cput=168:0:0  
  
module load singularity  
  
cd /rsgrps/shaneburgess/amanda/i5k/combine_gafs  
  
singularity pull docker://agbase/combine_gafs:1.0  
  
singularity run \  
-B /rsgrps/shaneburgess/amanda/i5k/combine_gafs:/work-dir \  
combine_gafs_1.0.sif \  
-i CFLO_1.fa_gaf.txt \  
-g clfol_v_insecta_goanna_gaf.tsv \  
-o complete_gaf
```

### Command Explained

**singularity run:** tells Singularity to run

**-B /rsgrps/shaneburgess/amanda/i5k/combine\_gafs:/work-dir:** mounts my current working directory on the host machine to '/work-dir' in the container

**combine\_gafs\_1.0.sif:** the name of the Singularity image file to use

---

**Tip:** All the options supplied after the image name are GOanna options

---

**-i CFLO\_1.fa\_gaf.txt:** InterProScan XML Parser GAF output file.

**-g clfo1\_v\_insecta\_goanna\_gaf.tsv:** GOanna GAF output file.

**-o complete\_gaf:** output file basename—a .tsv extension will be added

## 2.13 Combine GAFs on the ARS Ceres HPC

### 2.13.1 About Ceres/Scinet

- The Scinet VRSC has installed `combine_gafs` for ARS use.
- For general information on Scinet/Ceres, how to access it, and how to use it, visit <https://usda-ars-gbru.github.io/scinet-site/>.

### 2.13.2 Running GOanna on Ceres

---

#### Running programs on Ceres/Scinet

- You'll need to run `combine_gafs` either in interactive mode or batch mode.
  - For interactive mode, use the `salloc` command.
  - For batch mode, you'll need to write a batch job submission bash script.
- 

#### Running `combine_gafs` in interactive mode

##### Loading the module

The Scinet VRSC has installed the `combine_gafs` program. To load the module in interactive mode, run the command

```
module load agbase
```

#### Running Combine GAFs

Combine GAFs has three parameters:

```
-i InterProScan XML Parser GAF output
-g GOanna GAF output
-o output file basename
```

## Example Command

```
combine_gafs -i CFLO_1.fa_gaf.txt -g clfo1_v_insecta_goanna_gaf.tsv -o complete_gaf
```

## Command Explained

- i CFLO\_1.fa\_gaf.txt:** InterProScan XML Parser GAF output file.
- g clfo1\_v\_insecta\_goanna\_gaf.tsv:** GOanna GAF output file.
- o complete\_gaf:** output file basename—a .tsv extension will be added

## Running combine\_gafs in batch mode

---

### Running programs on Ceres/Scinet in batch mode

- Before using batch mode, you should review Scinet/Ceres' documentation first, and decide what queue you'll want to use. See <https://usda-ars-gbru.github.io/scinet-site/guide/ceres/>.

### Example batch job submission bash script (e.g. combine\_gafs-job.sh):

```
#!/bin/bash
module load agbase
combine_gafs -i CFLO_1.fa_gaf.txt -g clfo1_v_insecta_goanna_gaf.tsv -o complete_gaf
```

### Submitting the batch job:

```
sbatch combine_gafs-job.sh
```

## Command Explained

- i CFLO\_1.fa\_gaf.txt:** InterProScan XML Parser GAF output file.
- g clfo1\_v\_insecta\_goanna\_gaf.tsv:** GOanna GAF output file.
- o complete\_gaf:** output file basename—a .tsv extension will be added

## 2.14 Intro

- KEGG Orthology Based Annotation System (KOBAS) is a standalone Python application.
- **Consists of two modules:**
  1. **annotate**—Assigns appropriate KO terms for queried sequences based on a similarity search.
  2. **identify**—Discovers enriched KO terms among the annotation results by frequency of pathways or statistical significance of pathways.



Results and analysis from the application of KOBAS annotation to the Official gene set v3.0 protein set from *Diaphorina citri* followed by a differential expression analysis was presented at a seminar in the University of Arizona Animal and Comparative Biomedical Sciences in Fall 2020. The [slides](#) and [video](#) are available online.

### 2.14.1 Where to Find KOBAS

KOBAS is provided as a Docker container for use on the command line and as a group of apps in the CyVerse Discovery Environment.

- [Docker Hub](#)
- [KOBAS annotate 3.0.3](#)
- [KOBAS identify 3.0.3](#)
- [KOBAS annotate and identify 3.0.3](#)

**Note:** Each of these tools accepts a peptide FASTA file. For those users with nucleotide sequences some documentation has been provided for using **TransDecoder** (although other tools are also acceptable). The [TransDecoder app](#) is available through CyVerse or as a [BioContainer](#) for use on the command line.

### 2.14.2 Getting the KOBAS Databases

To run the tool you need some public data. The files can be downloaded directly from the [KOBAS homepage](#). These directories are also available as two tar archives in the CyVerse Data Store. The files are best downloaded with [iCommands](#). Once [iCommands](#) is [setup](#) you can use 'iget' to download the data.

- 1) seq\_pep.tar: species-specific BLAST databases used by KOBAS

```
iget /iplant/home/shared/iplantcollaborative/protein_blast_dbs/kobas/seq_pep.tar
tar -xf seq_pep.tar
```

- 2) sqlite3.tar: species-specific annotation databases used by KOBAS

```
iget /iplant/home/shared/iplantcollaborative/protein_blast_dbs/kobas/sqlite3.tar
tar -xf sqlite3.tar
```

**Note:** The above commands should result in two directories (seq\_pep and sqlite3) each containing many files. There is no need to unzip the .gz files.

### 2.14.3 Help and Usage Statement

On the command line the following help statement can be displayed with the option '-h'.

```
Options:
[-h prints this help statement]

[-a runs KOBAS annotate]
KOBAS annotate options:
```

(continues on next page)

(continued from previous page)

```

-i INFILE can be FASTA or one-per-line identifiers. See -t intype for details.
-s SPECIES 3 or 4 letter species abbreviation (can be found here: ftp://ftp.cbi.
↪pku.edu.cn/pub/KOBAS_3.0_DOWNLOAD/species_abbr.txt or here: https://www.kegg.jp/
↪kegg/catalog/org_list.html)
-o OUTPUT file (Default is stdout.)
-t INTYPE (fasta:pro, fasta:nuc, blastout:xml, blastout:tab, id:ncbigi,
↪id:uniprot, id:ensembl, id:ncbigene), default fasta:pro
[-l LIST available species, or list available databases for a specific species]
[-e EVALUE expect threshold for BLAST, default 1e-5]
[-r RANK rank cutoff for valid hits from BLAST result, default is 5]
[-C COVERAGE subject coverage cutoff for BLAST, default 0]
[-z ORTHOLOG whether only use orthologs for cross-species annotation or not,
↪default NO (if only using orthologs, please provide the species abbreviation of
↪your input)]
[-k KOBAS HOME The path to kobas_home, which is the parent directory of sqlite3/
↪and seq_pep/. This is the absolute path in the container.]
[-v BLAST HOME The path to blast_home, which is the parent directory of blastx
↪and blastp. This is the absolute path in the container.]
[-y BLASTDB The path to seq_pep/. This is the absolute path in the container.]
[-q KOBASDB The path to sqlite3/. This is the absolute path in the container.]
[-p BLASTP The path to blastp. This is the absolute path in the container.]
[-x BLASTX The path to blastx. This is the absolute path in the container.]
[-T number of THREADS to use in BLAST search. Default = 8]

[-g runs KOBAS identify]
KOBAS identify options:
-f FGFILE foreground file, the output of annotate
-b BGFILE background file, species abbreviation, see this list for species codes:
↪https://www.kegg.jp/kegg/catalog/org_list.html
-o OUTPUT file (Default is stdout.)
[-d DB databases for selection, 1-letter abbreviation separated by "/": K for
↪KEGG PATHWAY, n for PID, b for BioCarta, R for Reactome, B for BioCyc, p for
↪PANTHER,
  o for OMIM, k for KEGG DISEASE, f for FunDO, g for GAD, N for NHGRI GWAS
↪Catalog and G for Gene Ontology, default K/n/b/R/B/p/o/k/f/g/N/]
[-m METHOD choose statistical test method: b for binomial test, c for chi-square
↪test, h for hypergeometric test / Fisher's exact test, and x for frequency list,
  default hypergeometric test / Fisher's exact test]
[-n FDR choose false discovery rate (FDR) correction method: BH for Benjamini and
↪Hochberg, BY for Benjamini and Yekutieli, QVALUE, and None, default BH]
[-c CUTOFF terms with less than cutoff number of genes are not used for
↪statistical tests, default 5]
[-k KOBAS HOME The path to kobas_home, which is the parent directory of sqlite3/
↪and seq_pep/. This is the absolute path in the container.]
[-v BLAST HOME The path to blast_home, which is the parent directory of blastx
↪and blastp. This is the absolute path in the container.]
[-y BLASTDB The path to seq_pep/. This is the absolute path in the container.]
[-q KOBASDB The path to sqlite3/. This is the absolute path in the container.]
[-p BLASTP The path to blastp. This is the absolute path in the container.]
[-x BLASTX The path to blastx. This is the absolute path in the container.]

[-j runs both KOBAS annotate and identify]

```

## 2.15 KOBAS on CyVerse

### 2.15.1 Accessing KOBAS in the Discovery Environment

1. Create an account on CyVerse (free)
2. Open the CyVerse Discovery Environment (DE) and login with your CyVerse credentials.
3. If you are new to the Discovery Environment (DE) the user guide can be found [here](#).
4. Click on the 'Data' button at the left side of the screen to access your files/folders. Upload your data to the DE.
5. To access the KOBAS apps click on the 'Apps' button at the left side of the DE.
6. Search for 'kobas' in the search bar at the top of the 'apps' window (see below). The contents of the folder will appear in the main pane of the window; click on the name to open the app.

The KOBAS apps are called:

- KOBAS annotate 3.0.3
- KOBAS identify 3.0.3
- KOBAS annotate and identify 3.0.3

The screenshot shows the 'Apps' window in the CyVerse Discovery Environment. The search bar at the top contains 'kobas', and the results show three apps. The left sidebar shows the 'My Communities' section with 'AgBase' and 'i5k' listed.

Name	Integrated by	Rating	Syst...
β KOBAS annotate 3.0.3	Amanda Cooksey	☆☆☆☆☆ (0)	de
β KOBAS annotate and identif...	Amanda Cooksey	☆☆☆☆☆ (0)	de
β KOBAS identify 3.0.3	Amanda Cooksey	☆☆☆☆☆ (0)	de

#### Find Apps Easily with 'Communities'

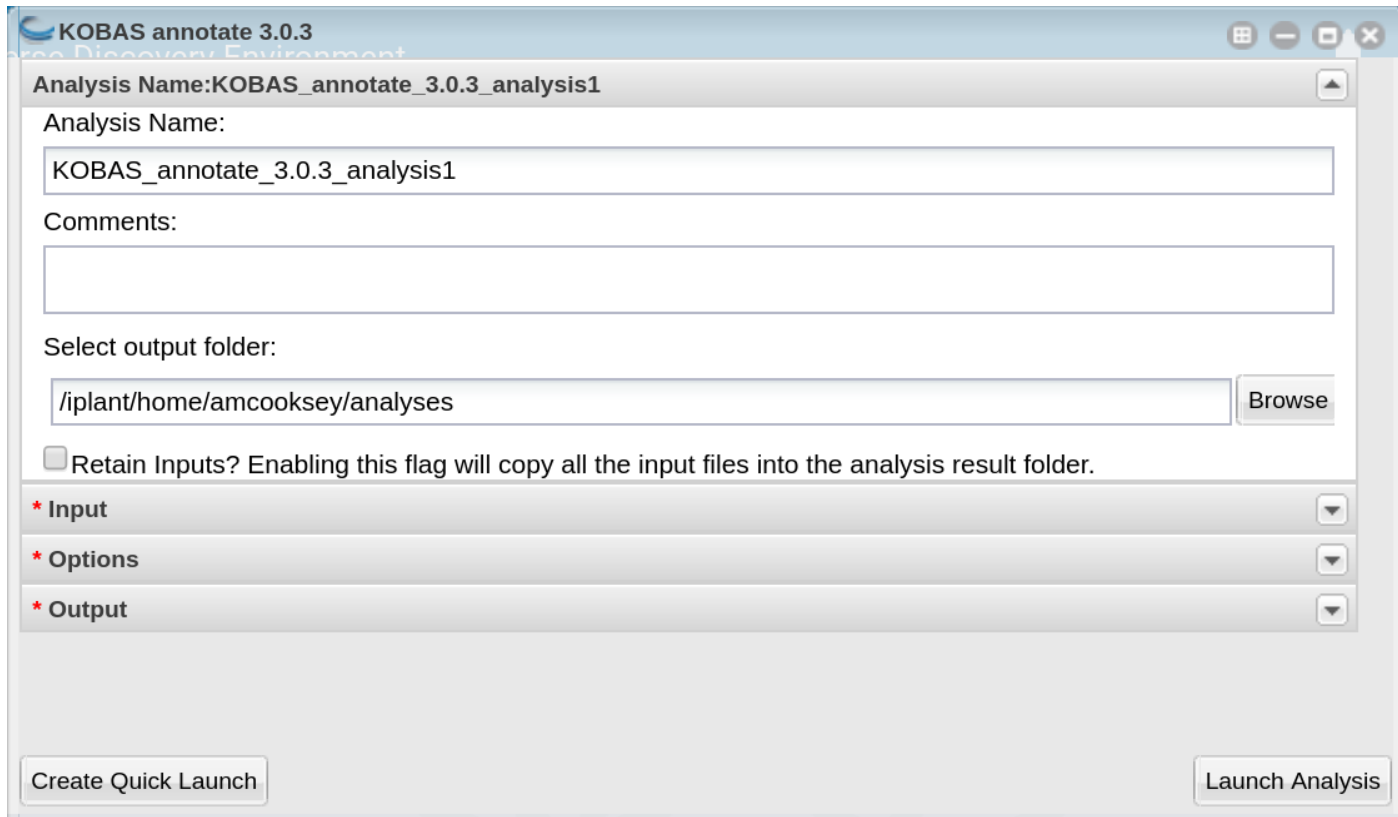
The KOBAS apps belong to the 'i5k' and 'AgBase' communities. You can join either of these communities and they will appear in the left-hand pane of your 'Apps' window (see above). To join a community:

1. Click on the person icon in the top-right corner of the Discovery Environment window
2. Select 'Communities'

3. In the 'Communities' window choose 'all communities' from the drop-down list. A list of communities will appear in the main pane of this window.
  4. Select the one you wish to join by clicking on it and then clicking on the 'join' button.
- 

## 2.15.2 KOBAS annotate 3.0.3

### Launching the app



**Analysis Name: KOBAS\_annotate\_3.0.3\_analysis\_1:** This menu is used to name the job you will run so that you can find it later. Analysis Name: The default name is "KOBAS\_annotate\_3.0.3\_analysis1". We recommend changing the 'analysis1' portion of this to reflect the data you are running.

**Comments:** (Optional) You can add additional information in the comments section to distinguish your analyses further.

**Select output folder:** This is where your results will be placed. The default (recommended) is your 'analyses' folder.

**Retain Inputs:** Enabling this flag will copy all the input files into the analysis result folder.

**Warning:** Selecting this option will rapidly consume your allocated space. It is not recommended. Your inputs will always remain available in the folder in which you stored them.

## Input

**Input File:** Use the ‘browse’ button on the right side of the field to navigate to your input file.

**Input File Type:** Select your input file type from the drop-down list. If your file type isn’t there then the app does not support that file type.

## Options

**Species Code:** Enter the species for the species of the sequences in your input file.

---

**Note:** If you don’t know the code for your species it can be found here: [https://www.kegg.jp/kegg/catalog/org\\_list.html](https://www.kegg.jp/kegg/catalog/org_list.html)

If your species of interest is not available then you should choose the code for the closest-related species available

---

**E value:** This is the evalue to use in the BLAST search. The default is 1e-5.

**Rank:** rank cutoff for valid hits from BLAST result. Default is 5.

**Coverage:** subject coverage cutoff for BLAST. Default is 0.

**Ortholog:** when checked KOBAS will only use orthologs for cross species annotation.

## Output

**Output File Name:** Provide an output file name .

For information on outputs see *Understanding Your Results: Annotate*

## Understanding Your Annotate Results

If all goes well, you should get the following:

- **logs folder:** This folder contains the ‘conder\_stderr’ and ‘condor\_stdout’ files. The files record feedback, progress and, importantly, any errors the app encountered during the analysis. You won’t normally need to look at these but they are very helpful in figuring out what may have happened if your output doesn’t look like you expected.
- **seq\_pep folder:** This folder contains the BLAST database files used in your analysis.
- **sqlite3 folder:** This folder contains the annotation database files used in your analysis
- **<species>.tsv:** This is the tab-delimited output from the BLAST search. It is unlikely that you will need to look at this file.
- **<output\_file\_name\_you\_provided>:** KOBAS-annotate generates a text file with the name you provide. It has two sections.

The first sections looks like this:

```
##dme      Drosophila melanogaster (fruit fly)
##Method: BLAST      Options: evalue <= 1e-05
##Summary: 87 succeed, 0 fail

#Query      Gene ID|Gene name|Hyperlink
```

(continues on next page)

(continued from previous page)

```
lcl|NW_020311285.1_prot_XP_012256083.1_15   dme:Dmel_CG34349|Unc-13-4B|http://www.
↳genome.jp/dbget-bin/www_bget?dme:Dmel_CG34349
lcl|NW_020311286.1_prot_XP_020708336.1_46   dme:Dmel_CG6963|gish|http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG6963
lcl|NW_020311285.1_prot_XP_020707987.1_39   dme:Dmel_CG30403||http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG30403
```

The second section follows a dashed line and looks like this:

```
-----
////
Query:                lcl|NW_020311285.1_prot_XP_012256083.1_15
Gene:                 dme:Dmel_CG34349          Unc-13-4B
Entrez Gene ID:      43002
////
Query:                lcl|NW_020311286.1_prot_XP_020708336.1_46
Gene:                 dme:Dmel_CG6963 gish
Entrez Gene ID:      49701
Pathway:              Hedgehog signaling pathway - fly          KEGG PATHWAY  ↳
↳dme04341
////
Query:                lcl|NW_020311285.1_prot_XP_020707987.1_39
Gene:                 dme:Dmel_CG30403
Entrez Gene ID:      246595
////
Query:                lcl|NW_020311285.1_prot_XP_020707989.1_40
Gene:                 dme:Dmel_CG6148 Past1
Entrez Gene ID:      41569
Pathway:              Endocytosis          KEGG PATHWAY          dme04144
                    Hemostasis          Reactome              R-DME-109582
                    Factors involved in megakaryocyte development and
↳platelet production  Reactome              R-DME-98323
```

If your analysis doesn't complete as you expected please look at your 'condor\_stderr' and 'condor\_stdout' files. If that doesn't clarify the problem contact us at [agbase@email.arizona.edu](mailto:agbase@email.arizona.edu) or [support@cyverse.org](mailto:support@cyverse.org).

### 2.15.3 KOBAS identify 3.0.3

## Launching the App

**Analysis Name: KOBAS\_identify\_3.0.3\_analysis\_1:** This menu is used to name the job you will run so that you can find it later. Analysis Name: The default name is “KOBAS\_identify\_3.0.3\_analysis1”. We recommend changing the ‘analysis1’ portion of this to reflect the data you are running.

**Comments:** (Optional) You can add additional information in the comments section to distinguish your analyses further.

**Select output folder:** This is where your results will be placed. The default (recommended) is your ‘analyses’ folder.

**Retain Inputs:** Enabling this flag will copy all the input files into the analysis result folder.

**Warning:** Selecting this option will rapidly consume your allocated space. It is not recommended. Your inputs will always remain available in the folder in which you stored them.

## Input

**Foreground File:** Use the ‘browse’ button on the right side of the field to navigate to your input file. This should be the output of KOBAS annotate.

**Background:** Enter the species for the species of the sequences in your input file.

**Note:** If you don’t know the code for your species it can be found here: [https://www.kegg.jp/kegg/catalog/org\\_list.html](https://www.kegg.jp/kegg/catalog/org_list.html)

If your species of interest is not available then you should choose the code for the closest-related species available

## Options

**Cutoff:** Annotation terms with less than cutoff number of genes are not used for statistical tests. Default is 5.

**Method:** Choose the statistical method to be used from the drop-down list. Default is hypergeometric/Fisher's Exact.

**FDR:** Method for determining false discovery rate. Default is Benjamini-Hochberg.

## Output

**Output File Name:** Provide an output file name .

## Understanding Your Identify Results

If all goes well, you should get the following:

- **logs folder:** This folder contains the 'condor\_stderr' and 'condor\_stdout' files. The files record feedback, progress and, importantly, any errors the app encountered during the analysis. You won't normally need to look at these but they are very helpful in figuring out what may have happened if your output doesn't look like you expected.
- **sqlite3 folder:** This folder contains the annotation database files used in your analysis
- **<output\_file\_name\_you\_provided>:** KOBAS identify generates a text file with the name you provide.

```
##Databases: PANTHER, KEGG PATHWAY, Reactome, BioCyc
##Statistical test method: hypergeometric test / Fisher's exact test
##FDR correction method: Benjamini and Hochberg

#Term          Database      ID          Input number  Background number  P-Value_
↳Corrected P-Value      Input      Hyperlink
Hedgehog signaling pathway - fly      KEGG PATHWAY      dme04341      12      33      3.
↳20002656734e-18      1.76001461204e-16      lcl|NW_020311286.1_prot_XP_012256678.
↳1_51|lcl|NW_020311286.1_prot_XP_025602973.1_48|lcl|NW_020311286.1_prot_XP_012256683.
↳1_52|lcl|NW_020311286.1_prot_XP_012256679.1_55|lcl|NW_020311286.1_prot_XP_012256674.
↳1_54|lcl|NW_020311286.1_prot_XP_020708336.1_46|lcl|NW_020311285.1_prot_XP_012256108.
↳1_32|lcl|NW_020311286.1_prot_XP_012256682.1_53|lcl|NW_020311286.1_prot_XP_025603025.
↳1_47|lcl|NW_020311286.1_prot_XP_020708334.1_49|lcl|NW_020311285.1_prot_XP_012256109.
↳1_33|lcl|NW_020311286.1_prot_XP_020708333.1_50 http://www.genome.jp/kegg-bin/show_
↳pathway?dme04341/dme:Dmel_CG6963%09red/dme:Dmel_CG6054%09red
Hedgehog signaling pathway      PANTHER      P00025      6      13      3.6166668094e-10      9.
↳94583372585e-09      lcl|NW_020311286.1_prot_XP_025602279.1_78|lcl|NW_020311286.1_
↳prot_XP_025602289.1_76|lcl|NW_020311286.1_prot_XP_025602264.1_79|lcl|NW_020311285.1_
↳prot_XP_012256108.1_32|lcl|NW_020311285.1_prot_XP_012256109.1_33|lcl|NW_020311286.1_
↳prot_XP_012256943.1_77      http://www.pantherdb.org/pathway/pathwayDiagram.jsp?
↳catAccession=P00025
Signaling by NOTCH2      Reactome      R-DME-1980145      3      8      2.00259649553e-05_
↳0.000275357018136      lcl|NW_020311285.1_prot_XP_012256118.1_28|lcl|NW_
↳020311285.1_prot_XP_012256117.1_27|lcl|NW_020311285.1_prot_XP_012256119.1_26
↳http://www.reactome.org/cgi-bin/eventbrowser_st_id?ST_ID=R-DME-1980145
```

If your analysis doesn't complete as you expected please look at your 'condor\_stderr' and 'condor\_stdout' files. If that doesn't clarify the problem contact us at [agbase@email.arizona.edu](mailto:agbase@email.arizona.edu) or [support@cyverse.org](mailto:support@cyverse.org).



## 2.15.4 KOBAS annotate and identify 3.0.3

### Launching the App

This app runs both the annotate and identify analyses together as a convenience for user who wish to run both steps.

**Analysis Name: KOBAS\_annotate\_and\_identify\_3.0.3\_analysis\_1:** This menu is used to name the job you will run so that you can find it later. Analysis Name: The default name is “KOBAS\_annotate\_identify\_3.0.3\_analysis1”. We recommend changing the ‘analysis1’ portion of this to reflect the data you are running.

**Comments:** (Optional) You can add additional information in the comments section to distinguish your analyses further.

**Select output folder:** This is where your results will be placed. The default (recommended) is your ‘analyses’ folder.

**Retain Inputs:** Enabling this flag will copy all the input files into the analysis result folder.

**Warning:** Selecting this option will rapidly consume your allocated space. It is not recommended. Your inputs will always remain available in the folder in which you stored them.

### Input

**Input File:** Use the ‘browse’ button on the right side of the field to navigate to your input file.

**Input File Type:** Select your input file type from the drop-down list. If your file type isn’t there then the app does not support that file type.

### Annotate Options

**Species Code:** Enter the species for the species of the sequences in your input file.

---

**Note:** If you don't know the code for your species it can be found here: [https://www.kegg.jp/kegg/catalog/org\\_list.html](https://www.kegg.jp/kegg/catalog/org_list.html)

If your species of interest is not available then you should choose the code for the closest-related species available

---

**E value:** This is the evaluate to use in the BLAST search. The default is 1e-5.

**Rank:** rank cutoff for valid hits from BLAST result. Default is 5.

**Coverage:** subject coverage cutoff for BLAST. Default is 0.

**Ortholog:** when checked KOBAS will only use orthologs for cross species annotation.

### Identify Options

**Cutoff:** Annotation terms with less than cutoff number of genes are not used for statistical tests. Default is 5.

**Method:** Choose the statistical method to be used from the drop-down list. Default is hypergeometric/Fisher's Exact.

**FDR:** Method for determining false discovery rate. Default is Benjamini-Hochberg.

### Output

**Output File Basename:** This will be the prefix of your output files.

### Understanding Your Annotate and Identify Pipeline Results

If all goes well, you should get the following:

- **logs folder:** This folder contains the 'conder\_stderr' and 'condor\_stdout' files. The files record feedback, progress and, importantly, any errors the app encountered during the analysis. You won't normally need to look at these but they are very helpful in figuring out what may have happened if your output doesn't look like you expected.
- **sqlite3 folder:** This folder contains the annotation database files used in your analysis
- **seq\_pep folder:** This folder contains the BLAST database files used in your analysis.
- **<species>.tsv:** This is the tab-delimited output from the BLAST search. It is unlikely that you will need to look at this file.
- **<basename>\_annotate\_out.txt:** KOBAS annotate generates a text file with the name you provide. It has two sections.

The first sections looks like this:

```
##dme      Drosophila melanogaster (fruit fly)
##Method:  BLAST      Options:  evaluate <= 1e-05
##Summary: 87 succeed, 0 fail

#Query      Gene ID|Gene name|Hyperlink
lc1|NW_020311285.1_prot_XP_012256083.1_15   dme:Dmel_CG34349|Unc-13-4B|http://www.
genome.jp/dbget_bin/www_bget?dme:Dmel_CG34349
```

(continues on next page)

(continued from previous page)

```
lcl|NW_020311286.1_prot_XP_020708336.1_46 dme:Dmel_CG6963|gish|http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG6963
lcl|NW_020311285.1_prot_XP_020707987.1_39 dme:Dmel_CG30403||http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG30403
```

The second section follows a dashed line and looks like this:

```
-----
////
Query:          lcl|NW_020311285.1_prot_XP_012256083.1_15
Gene:           dme:Dmel_CG34349          Unc-13-4B
Entrez Gene ID: 43002
////
Query:          lcl|NW_020311286.1_prot_XP_020708336.1_46
Gene:           dme:Dmel_CG6963 gish
Entrez Gene ID: 49701
Pathway:        Hedgehog signaling pathway - fly          KEGG PATHWAY
↳dme04341
////
Query:          lcl|NW_020311285.1_prot_XP_020707987.1_39
Gene:           dme:Dmel_CG30403
Entrez Gene ID: 246595
////
Query:          lcl|NW_020311285.1_prot_XP_020707989.1_40
Gene:           dme:Dmel_CG6148 Past1
Entrez Gene ID: 41569
Pathway:        Endocytosis          KEGG PATHWAY          dme04144
                Hemostasis          Reactome          R-DME-109582
                Factors involved in megakaryocyte development and
↳platelet production          Reactome          R-DME-98323mel_CG6963
lcl|NW_020311285.1_prot_XP_020707987.1_39 dme:Dmel_CG30403||http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG30403
```

- **<basename>\_identify\_out.txt:** KOBAS identify generates a text file with the name you provide.

```
##Databases: PANTHER, KEGG PATHWAY, Reactome, BioCyc
##Statistical test method: hypergeometric test / Fisher's exact test
##FDR correction method: Benjamini and Hochberg

#Term          Database          ID          Input number          Background number          P-Value
↳Corrected P-Value          Input          Hyperlink
Hedgehog signaling pathway - fly          KEGG PATHWAY          dme04341          12          33          3.
↳20002656734e-18          1.76001461204e-16          lcl|NW_020311286.1_prot_XP_012256678.
↳1_51|lcl|NW_020311286.1_prot_XP_025602973.1_48|lcl|NW_020311286.1_prot_XP_012256683.
↳1_52|lcl|NW_020311286.1_prot_XP_012256679.1_55|lcl|NW_020311286.1_prot_XP_012256674.
↳1_54|lcl|NW_020311286.1_prot_XP_020708336.1_46|lcl|NW_020311285.1_prot_XP_012256108.
↳1_32|lcl|NW_020311286.1_prot_XP_012256682.1_53|lcl|NW_020311286.1_prot_XP_025603025.
↳1_47|lcl|NW_020311286.1_prot_XP_020708334.1_49|lcl|NW_020311285.1_prot_XP_012256109.
↳1_33|lcl|NW_020311286.1_prot_XP_020708333.1_50 http://www.genome.jp/kegg-bin/show_
↳pathway?dme04341/dme:Dmel_CG6963%09red/dme:Dmel_CG6054%09red
Hedgehog signaling pathway          PANTHER          P00025          6          13          3.6166668094e-10          9.
↳94583372585e-09          lcl|NW_020311286.1_prot_XP_025602279.1_78|lcl|NW_020311286.1_
↳prot_XP_025602289.1_76|lcl|NW_020311286.1_prot_XP_025602264.1_79|lcl|NW_020311285.1_
↳prot_XP_012256108.1_32|lcl|NW_020311285.1_prot_XP_012256109.1_33|lcl|NW_020311286.1_
↳prot_XP_012256943.1_77          http://www.pantherdb.org/pathway/pathwayDiagram.jsp?
↳catAccession=P00025
```

(continues on next page)

(continued from previous page)

```
Signaling by NOTCH2 Reactome      R-DME-1980145      3      8      2.00259649553e-05
↪      0.000275357018136      lc1|NW_020311285.1_prot_XP_012256118.1_28|lc1|NW_
↪020311285.1_prot_XP_012256117.1_27|lc1|NW_020311285.1_prot_XP_012256119.1_26
↪http://www.reactome.org/cgi-bin/eventbrowser_st_id?ST_ID=R-DME-1980145
```

If your analysis doesn't complete as you expected please look at your 'condor\_stderr' and 'condor\_stdout' files. If that doesn't clarify the problem contact us at [agbase@email.arizona.edu](mailto:agbase@email.arizona.edu) or [support@cyverse.org](mailto:support@cyverse.org).

## 2.16 KOBAS on the Command Line

### 2.16.1 Getting the databases

To run the tool you need some public data. The files can be downloaded directly from the [KOBAS homepage](#). These directories are also available as two tar archives in the CyVerse Data Store. The files are best downloaded with [iCommands](#). Once iCommands is [setup](#) you can use 'iget' to download the data.

- 1) seq\_pep.tar: species-specific BLAST databases used by KOBAS

```
iget /iplant/home/shared/iplantcollaborative/protein_blast_dbs/kobas/seq_pep.tar
tar -xf seq_pep.tar
```

- 2) sqlite3.tar: species-specific annotation databases used by KOBAS

```
iget /iplant/home/shared/iplantcollaborative/protein_blast_dbs/kobas/sqlite3.tar
tar -xf sqlite3.tar
```

---

**Note:** The above commands should result in two directories (seq\_pep and sqlite3) each containing many files. There is no need to unzip the .gz files.

---

### 2.16.2 Container Technologies

KOBAS is provided as a Docker container.

A container is a standard unit of software that packages up code and all its dependencies so the application runs quickly and reliably from one computing environment to another.

There are two major containerization technologies: **Docker** and **Singularity**.

Docker containers can be run with either technology.

### 2.16.3 Running KOBAS using Docker

---

#### About Docker

- Docker must be installed on the computer you wish to use for your analysis.
- To run Docker you must have 'root' permissions (or use sudo).

- Docker will run all containers as 'root'. This makes Docker incompatible with HPC systems (see Singularity below).
  - Docker can be run on your local computer, a server, a cloud virtual machine (such as CyVerse Atmosphere) etc. Docker can be installed quickly on an Atmosphere instance by typing 'ezd'.
  - For more information on installing Docker on other systems see this tutorial: [Installing Docker on your machine](#).
- 

## Getting the KOBAS container

The KOBAS tool is available as a Docker container on Docker Hub: [KOBAS container](#)

The container can be pulled with this command:

```
docker pull agbase/kobas:3.0.3_0
```

---

### Remember

You must have root permissions or use sudo, like so:

```
sudo docker pull agbase/kobas:3.0.3_0
```

---

## Running KOBAS with Data

### Getting the Help and Usage Statement

```
sudo docker run --rm -v $(pwd):/work-dir agbase/kobas:3.0.3_0 -h
```

---

See *Help and Usage Statement*

---

### Tip:

**There are 3 directories built into this container. These directories should be used to mount data.**

- /work-dir
  - /seq\_pep
  - /sqlite3
- 

KOBAS can perform two tasks - annotate (-a) - identify (enrichment) (-g)

KOBAS can also run both task with a single command (-j).

### Annotate Example Command

```
sudo docker run \  
--rm \  
-v /home/amcooksey/i5k/seq_pep:/seq_pep \  
-v /home/amcooksey/i5k/sqlite3:/sqlite3 \  
-v $(pwd):/work-dir \  
agbase/kobas:3.0.3_0 \  

```

---

(continues on next page)

(continued from previous page)

```
-a
-i AROS1000.fa \
-s dme \
-t fasta:pro \
-o AROS1000
```

## Command Explained

**sudo docker run:** tells docker to run

**-rm:** removes the container when the analysis has finished. The image will remain for future use.

**-v /home/amcooksey/i5k/seq\_pep:/seq\_pep:** tells docker to mount the 'seq\_pep' directory I downloaded to the host machine to the '/seq\_pep' directory within the container. The syntax for this is: <absolute path on host>:<absolute path in container>

**-v /home/amcooksey/i5k/sqlite3:/sqlite3:** mounts 'sqlite3' directory on host machine into 'sqlite3' directory inside the container

**-v \$(pwd):/work-dir:** mounts my current working directory on the host machine to '/work-dir' in the container

**agbase/kobas:3.0.3\_0:** the name of the Docker image to use

---

**Tip:** All the options supplied after the image name are KOBAS options

---

**-a:** Tells KOBAS to run the 'annotate' process.

**-i AROS1000.fa:** input file (peptide FASTA)

**-s dme:** Enter the species for the species of the sequences in your input file.

---

**Note:** If you don't know the code for your species it can be found here: [https://www.kegg.jp/kegg/catalog/org\\_list.html](https://www.kegg.jp/kegg/catalog/org_list.html)

If your species of interest is not available then you should choose the code for the closest-related species available

---

**-t:** input file type; in this case, protein FASTA.

**-o AROS1000:** name of output file

For information on output files see *Understanding Your Results: Annotate*

## Identify Example Command

```
sudo docker run \
--rm \
-v /home/amcooksey/i5k/seq_pep:/seq_pep \
-v /home/amcooksey/i5k/sqlite3:/sqlite3 \
-v $(pwd) :/work-dir \
agbase/kobas:3.0.3_0 \
-g \
-f AROS1000 \
-b dme \
-o ident_out
```

## Command Explained

**sudo docker run:** tells docker to run

**-rm:** removes the container when the analysis has finished. The image will remain for future use.

**-v /home/amcooksey/i5k/seq\_pep:/seq\_pep:** tells docker to mount the 'seq\_pep' directory I downloaded to the host machine to the '/seq\_pep' directory within the container. The syntax for this is: <absolute path on host>:<absolute path in container>

**-v /home/amcooksey/i5k/sqlite3:/sqlite3:** mounts 'sqlite3' directory on host machine into 'go\_info' directory inside the container

**-v \$(pwd):/work-dir:** mounts my current working directory on the host machine to '/work-dir' in the container

**agbase/kobas:3.0.3\_0:** the name of the Docker image to use

---

**Tip:** All the options supplied after the image name are KOBAS options

---

**-g:** Tells KOBAS to run the 'identify' process.

**-f AROS1000:** output file from KOBAS annotate

**-b dme:** background; enter the species code for the species of the sequences in your input file.

---

**Note:** If you don't know the code for your species it can be found here: [https://www.kegg.jp/kegg/catalog/org\\_list.html](https://www.kegg.jp/kegg/catalog/org_list.html)

If your species of interest is not available then you should choose the code for the closest-related species available

---

**-o ident\_out:** name of output file

For information on outputs see *Understanding Your Results: Identify*

## Annotate and Identify Pipeline Example Command

```
sudo docker run \  
--rm \  
-v /home/amcooksey/i5k/seq_pep:/seq_pep \  
-v /home/amcooksey/i5k/sqlite3:/sqlite3 \  
-v $(pwd):/work-dir \  
agbase/kobas:3.0.3_0 \  
-j \  
-i AROS1000.fa \  
-s dme \  
-t fasta:pro \  
-o AROS1000
```

## Command Explained

**sudo docker run:** tells docker to run

**-rm:** removes the container when the analysis has finished. The image will remain for future use.

**-v /home/amcooksey/i5k/seq\_pep:/seq\_pep:** tells docker to mount the 'seq\_pep' directory I downloaded to the host machine to the '/seq\_pep' directory within the container. The syntax for this is: <absolute path on host>:<absolute path in container>

**-v /home/amcooksey/i5k/sqlite3:/sqlite3:** mounts 'sqlite3' directory on host machine into 'go\_info' directory inside the container

**-v \$(pwd):/work-dir:** mounts my current working directory on the host machine to '/work-dir' in the container

**agbase/kobas:3.0.3\_0:** the name of the Docker image to use

---

**Tip:** All the options supplied after the image name are KOBAS options

---

**-j:** Tells KOBAS to run the 'annotate' process.

**-i AROS1000.fa:** input file (peptide FASTA)

**-s dme:** Enter the species for the species of the sequences in your input file.

---

**Note:** If you don't know the code for your species it can be found here: [https://www.kegg.jp/kegg/catalog/org\\_list.html](https://www.kegg.jp/kegg/catalog/org_list.html)

If your species of interest is not available then you should choose the code for the closest-related species available

---

**-t:** input file type; in this case, protein FASTA.

**-o AROS1000:** basename of output files

---

**Note:** This pipeline will automatically use the output of 'annotate' as the -f foreground input for 'identify'. This will also use your species option as the -b background input for 'identify'.

---

For more information on outputs see *Understanding Your Results: Annotate and Identify*

### 2.16.4 Running KOBAS using Singularity

---

#### About Singularity

- does not require 'root' permissions
  - runs all containers as the user that is logged into the host machine
  - HPC systems are likely to have Singularity installed and are unlikely to object if asked to install it (no guarantees).
  - can be run on any machine where it is installed
  - more information about [installing Singularity](#)
  - This tool was tested using Singularity 3.0. Users with Singularity 2.x will need to modify the commands accordingly.
- 

#### HPC Job Schedulers

---



Although Singularity can be installed on any computer this documentation assumes it will be run on an HPC system. The tool was tested on a PBSPro system and the job submission scripts below reflect that. Submission scripts will need to be modified for use with other job scheduler systems.

## Getting the KOBAS container

The KOBAS tool is available as a Docker container on Docker Hub: [KOBAS container](#)

The container can be pulled with this command:

```
singularity pull docker://agbase/kobas:3.0.3_0
```

## Running KOBAS with Data

### Getting the Help and Usage Statement

#### Example PBS script:

```
#!/bin/bash
#PBS -N kobas
#PBS -W group_list=fionamcc
#PBS -l select=1:ncpus=28:mem=168gb
#PBS -q standard
#PBS -l walltime=6:0:0
#PBS -l cput=168:0:0

module load singularity

cd /rsgrps/shaneburgess/amanda/i5k/kobas

singularity pull docker://agbase/kobas:3.0.3_0

singularity run \
kobas_3.0.3_0.sif \
-h
```

See *Help and Usage Statement*

**Tip:** There are 3 directories built into this container. These directories should be used to mount data.

- /seq\_pep
- /sqlite3
- /work-dir

## Example PBS Script for Annotate Process

```
#!/bin/bash
#PBS -N kobas
#PBS -W group_list=fionamcc
```

(continues on next page)

(continued from previous page)

```

#PBS -l select=1:ncpus=28:mem=168gb
#PBS -q standard
#PBS -l walltime=6:0:0
#PBS -l cput=168:0:0

module load singularity

cd /rsgrps/shaneburgess/amanda/i5k/kobas

singularity pull docker://agbase/kobas:3.0.3_0

singularity run \
-B /rsgrps/shaneburgess/amanda/i5k/seq_pep:/seq_pep \
-B /rsgrps/shaneburgess/amanda/i5k/sqlite3:/sqlite3 \
-B /rsgrps/shaneburgess/amanda/i5k/kobas:/work-dir \
kobas_3.0.3_0.sif \
-a
-i AROS1000.fa \
-s dme \
-t fasta:pro \
-o AROS1000 \

```

## Command Explained

**singularity run:** tells Singularity to run

**-B /rsgrps/shaneburgess/amanda/i5k/seq\_pep:/seq\_pep:** tells docker to mount the 'seq\_pep' directory I downloaded to the host machine to the '/seq\_pep' directory within the container. The syntax for this is: <absolute path on host>:<absolute path in container>

**-B /rsgrps/shaneburgess/amanda/i5k/sqlite3:/sqlite3:** mounts 'sqlite3' directory on host machine into 'go\_info' directory inside the container

**-B /rsgrps/shaneburgess/amanda/i5k/kobas:/work-dir:** mounts my current working directory on the host machine to '/work-dir' in the container

**kobas\_3.0.3\_0.sif:** the name of the Singularity image to use

---

**Tip:** All the options supplied after the image name are KOBAS options

---

**-a:** Tells KOBAS to run the 'annotate' process.

**-i AROS1000.fa:** input file (peptide FASTA)

**-s dme:** Enter the species for the species of the sequences in your input file.

---

**Note:** If you don't know the code for your species it can be found here: [https://www.kegg.jp/kegg/catalog/org\\_list.html](https://www.kegg.jp/kegg/catalog/org_list.html)

If your species of interest is not available then you should choose the code for the closest-related species available

---

**-t:** input file type; in this case, protein FASTA.

**-o AROS1000:** name of output file

For information on output files see *Understanding Your Results: Annotate*

### Example PBS Script for Identify Process

```
#!/bin/bash
#PBS -N kobas
#PBS -W group_list=fionamcc
#PBS -l select=1:ncpus=28:mem=168gb
#PBS -q standard
#PBS -l walltime=6:0:0
#PBS -l cput=168:0:0

module load singularity

cd /rsgrps/shaneburgess/amanda/i5k/kobas

singularity pull docker://agbase/kobas:3.0.3_0

singularity run \
-B /rsgrps/shaneburgess/amanda/i5k/seq_pep:/seq_pep \
-B /rsgrps/shaneburgess/amanda/i5k/sqlite3:/sqlite3 \
-B /rsgrps/shaneburgess/amanda/i5k/kobas:/work-dir \
kobas_3.0.3_0.sif \
-g
-f AROS1000 \
-b dme \
-o ident_out
```

### Command Explained

**singularity run:** tells Singularity to run

**-B /rsgrps/shaneburgess/amanda/i5k/seq\_pep:/seq\_pep:** tells docker to mount the 'seq\_pep' directory I downloaded to the host machine to the '/seq\_pep' directory within the container. The syntax for this is: <absolute path on host>:<absolute path in container>

**-B /rsgrps/shaneburgess/amanda/i5k/sqlite3:/sqlite3:** mounts 'sqlite3' directory on host machine into 'go\_info' directory inside the container

**-B /rsgrps/shaneburgess/amanda/i5k/kobas:/work-dir:** mounts my current working directory on the host machine to '/work-dir' in the container

**kobas\_3.0.3\_0.sif:** the name of the Singularity image to use

---

**Tip:** All the options supplied after the image name are KOBAS options

---

**-g:** Tells KOBAS to run the 'identify' process.

**-f AROS1000:** output file from 'annotate'

**-b dme:** background; enter the species for the species of the sequences in your input file.

---

**Note:** If you don't know the code for your species it can be found here: [https://www.kegg.jp/kegg/catalog/org\\_list.html](https://www.kegg.jp/kegg/catalog/org_list.html)

If your species of interest is not available then you should choose the code for the closest-related species available

---

**-o ident\_out:** name of output file

For information on output see *Understanding Your Results: Identify*

### Example PBS Script for Annotate and Identify Pipeline

```
#!/bin/bash
#PBS -N kobas
#PBS -W group_list=fionamcc
#PBS -l select=1:ncpus=28:mem=168gb
#PBS -q standard
#PBS -l walltime=6:0:0
#PBS -l cput=168:0:0

module load singularity

cd /rsgrps/shaneburgess/amanda/i5k/kobas

singularity pull docker://agbase/kobas:3.0.3_0

singularity run \
-B /rsgrps/shaneburgess/amanda/i5k/seq_pep:/seq_pep \
-B /rsgrps/shaneburgess/amanda/i5k/sqlite3:/sqlite3 \
-B /rsgrps/shaneburgess/amanda/i5k/kobas:/work-dir \
kobas_3.0.3_0.sif \
-j
-i AROS1000.fa \
-s dme \
-t fasta:pro \
-o AROS1000 \
```

### Command Explained

**singularity run:** tells Singularity to run

**-B /rsgrps/shaneburgess/amanda/i5k/seq\_pep:/seq\_pep:** tells docker to mount the 'seq\_pep' directory I downloaded to the host machine to the 'seq\_pep' directory within the container. The syntax for this is: <absolute path on host>:<absolute path in container>

**-B /rsgrps/shaneburgess/amanda/i5k/sqlite3:/sqlite3:** mounts 'sqlite3' directory on host machine into 'go\_info' directory inside the container

**-B /rsgrps/shaneburgess/amanda/i5k/kobas:/work-dir:** mounts my current working directory on the host machine to '/work-dir' in the container

**kobas\_3.0.3\_0.sif:** the name of the Singularity image to use

---

**Tip:** All the options supplied after the image name are KOBAS options

---

**-j:** Tells KOBAS to run the 'annotate' process.

**-i AROS1000.fa:** input file (peptide FASTA)

**-s dme:** Enter the species for the species of the sequences in your input file.

**Note:** If you don't know the code for your species it can be found here: [https://www.kegg.jp/kegg/catalog/org\\_list.html](https://www.kegg.jp/kegg/catalog/org_list.html)

If your species of interest is not available then you should choose the code for the closest-related species available

**-t:** input file type; in this case, protein FASTA.

**-o AROS1000:** name of output file

**Note:** This pipeline will automatically use the output of 'annotate' as the -f foreground input for 'identify'. This will also use your species option as the -b background input for 'identify'.

For information on outputs see *Understanding Your Results: Annotate and Identify*

## 2.16.5 Understanding Your Results

### Annotate

If all goes well, you should get the following:

- **seq\_pep folder:** This folder contains the BLAST database files used in your analysis.
- **sqlite3 folder:** This folder contains the annotation database files used in your analysis
- **<species>.tsv:** This is the tab-delimited output from the BLAST search. It is unlikely that you will need to look at this file.
- **<output\_file\_name\_you\_provided>:** KOBAS-annotate generates a text file with the name you provide. It has two sections.

The first sections looks like this:

```
##dme      Drosophila melanogaster (fruit fly)
##Method: BLAST      Options: evaluate <= 1e-05
##Summary:  87 succeed, 0 fail

#Query      Gene ID|Gene name|Hyperlink
lcl|NW_020311285.1_prot_XP_012256083.1_15    dme:Dmel_CG34349|Unc-13-4B|http://www.
↳genome.jp/dbget-bin/www_bget?dme:Dmel_CG34349
lcl|NW_020311286.1_prot_XP_020708336.1_46    dme:Dmel_CG6963|gish|http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG6963
lcl|NW_020311285.1_prot_XP_020707987.1_39    dme:Dmel_CG30403||http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG30403
```

The second section follows a dashed line and looks like this:

```
-----
////
Query:          lcl|NW_020311285.1_prot_XP_012256083.1_15
Gene:           dme:Dmel_CG34349      Unc-13-4B
Entrez Gene ID: 43002
////
Query:          lcl|NW_020311286.1_prot_XP_020708336.1_46
```

(continues on next page)

(continued from previous page)

```

Gene:                dme:Dmel_CG6963 gish
Entrez Gene ID:      49701
Pathway:             Hedgehog signaling pathway - fly          KEGG PATHWAY
↳dme04341
////
Query:               lcl|NW_020311285.1_prot_XP_020707987.1_39
Gene:                dme:Dmel_CG30403
Entrez Gene ID:      246595
////
Query:               lcl|NW_020311285.1_prot_XP_020707989.1_40
Gene:                dme:Dmel_CG6148 Past1
Entrez Gene ID:      41569
Pathway:             Endocytosis          KEGG PATHWAY      dme04144
                    Hemostasis          Reactome           R-DME-109582
                    Factors involved in megakaryocyte development and
↳platelet production Reactome           R-DME-98323
    
```

## Identify

If all goes well, you should get the following:

- **sqlite3 folder:** This folder contains the annotation database files used in your analysis
- **<output\_file\_name\_you\_provided>:** KOBAS identify generates a text file with the name you provide.

```

##Databases: PANTHER, KEGG PATHWAY, Reactome, BioCyc
##Statistical test method: hypergeometric test / Fisher's exact test
##FDR correction method: Benjamini and Hochberg

#Term          Database          ID          Input number    Background number    P-Value
↳Corrected P-Value      Input      Hyperlink
Hedgehog signaling pathway - fly      KEGG PATHWAY      dme04341          12          33          3.
↳20002656734e-18          1.76001461204e-16          lcl|NW_020311286.1_prot_XP_012256678.
↳1_51|lcl|NW_020311286.1_prot_XP_025602973.1_48|lcl|NW_020311286.1_prot_XP_012256683.
↳1_52|lcl|NW_020311286.1_prot_XP_012256679.1_55|lcl|NW_020311286.1_prot_XP_012256674.
↳1_54|lcl|NW_020311286.1_prot_XP_020708336.1_46|lcl|NW_020311285.1_prot_XP_012256108.
↳1_32|lcl|NW_020311286.1_prot_XP_012256682.1_53|lcl|NW_020311286.1_prot_XP_025603025.
↳1_47|lcl|NW_020311286.1_prot_XP_020708334.1_49|lcl|NW_020311285.1_prot_XP_012256109.
↳1_33|lcl|NW_020311286.1_prot_XP_020708333.1_50 http://www.genome.jp/kegg-bin/show_
↳pathway?dme04341/dme:Dmel_CG6963%09red/dme:Dmel_CG6054%09red
Hedgehog signaling pathway      PANTHER      P00025      6          13          3.6166668094e-10          9.
↳94583372585e-09          lcl|NW_020311286.1_prot_XP_025602279.1_78|lcl|NW_020311286.1_
↳prot_XP_025602289.1_76|lcl|NW_020311286.1_prot_XP_025602264.1_79|lcl|NW_020311285.1_
↳prot_XP_012256108.1_32|lcl|NW_020311285.1_prot_XP_012256109.1_33|lcl|NW_020311286.1_
↳prot_XP_012256943.1_77          http://www.pantherdb.org/pathway/pathwayDiagram.jsp?
↳catAccession=P00025
Signaling by NOTCH2      Reactome          R-DME-1980145      3          8          2.00259649553e-05
↳0.000275357018136          lcl|NW_020311285.1_prot_XP_012256118.1_28|lcl|NW_
↳020311285.1_prot_XP_012256117.1_27|lcl|NW_020311285.1_prot_XP_012256119.1_26
↳http://www.reactome.org/cgi-bin/eventbrowser_st_id?ST_ID=R-DME-1980145
    
```

## Annotate and Identify Pipeline

If all goes well, you should get the following:

- **logs folder:** This folder contains the ‘conder\_stderr’ and ‘condor\_stdout’ files. The files record feedback, progress and, importantly, any errors the app encountered during the analysis. You won’t normally need to look at these but they are very helpful in figuring out what may have happened if your output doesn’t look like you expected.
- **sqlite3 folder:** This folder contains the annotation database files used in your analysis
- **seq\_pep folder:** This folder contains the BLAST database files used in your analysis.
- **<species>.tsv:** This is the tab-delimited output from the BLAST search. It is unlikely that you will need to look at this file.
- **<basename>\_annotate\_out.txt:** KOBAS annotate generates a text file with the name you provide. It has two sections.

The first sections looks like this:

```
##dme      Drosophila melanogaster (fruit fly)
##Method: BLAST      Options: evaluate <= 1e-05
##Summary: 87 succeed, 0 fail

#Query      Gene ID|Gene name|Hyperlink
lcl|NW_020311285.1_prot_XP_012256083.1_15   dme:Dmel_CG34349|Unc-13-4B|http://www.
↳genome.jp/dbget-bin/www_bget?dme:Dmel_CG34349
lcl|NW_020311286.1_prot_XP_020708336.1_46   dme:Dmel_CG6963|gish|http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG6963
lcl|NW_020311285.1_prot_XP_020707987.1_39   dme:Dmel_CG30403||http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG30403
```

The second section follows a dashed line and looks like this:

```
-----
////
Query:      lcl|NW_020311285.1_prot_XP_012256083.1_15
Gene:      dme:Dmel_CG34349      Unc-13-4B
Entrez Gene ID:      43002
////
Query:      lcl|NW_020311286.1_prot_XP_020708336.1_46
Gene:      dme:Dmel_CG6963 gish
Entrez Gene ID:      49701
Pathway:      Hedgehog signaling pathway - fly      KEGG PATHWAY      ↳
↳dme04341
////
Query:      lcl|NW_020311285.1_prot_XP_020707987.1_39
Gene:      dme:Dmel_CG30403
Entrez Gene ID:      246595
////
Query:      lcl|NW_020311285.1_prot_XP_020707989.1_40
Gene:      dme:Dmel_CG6148 Past1
Entrez Gene ID:      41569
Pathway:      Endocytosis      KEGG PATHWAY      dme04144
      Hemostasis      Reactome      R-DME-109582
      Factors involved in megakaryocyte development and↳
↳platelet production      Reactome      R-DME-98323mel_CG6963
lcl|NW_020311285.1_prot_XP_020707987.1_39   dme:Dmel_CG30403||http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG30403
```

- **<basename>\_identify\_out.txt:** KOBAS identify generates a text file with the name you provide.

```
##Databases: PANTHER, KEGG PATHWAY, Reactome, BioCyc
##Statistical test method: hypergeometric test / Fisher's exact test
##FDR correction method: Benjamini and Hochberg

#Term          Database      ID          Input number  Background number  P-Value_
↳Corrected P-Value      Input      Hyperlink
Hedgehog signaling pathway - fly      KEGG PATHWAY      dme04341          12          33          3.
↳20002656734e-18          1.76001461204e-16          lcl|NW_020311286.1_prot_XP_012256678.
↳1_51|lcl|NW_020311286.1_prot_XP_025602973.1_48|lcl|NW_020311286.1_prot_XP_012256683.
↳1_52|lcl|NW_020311286.1_prot_XP_012256679.1_55|lcl|NW_020311286.1_prot_XP_012256674.
↳1_54|lcl|NW_020311286.1_prot_XP_020708336.1_46|lcl|NW_020311285.1_prot_XP_012256108.
↳1_32|lcl|NW_020311286.1_prot_XP_012256682.1_53|lcl|NW_020311286.1_prot_XP_025603025.
↳1_47|lcl|NW_020311286.1_prot_XP_020708334.1_49|lcl|NW_020311285.1_prot_XP_012256109.
↳1_33|lcl|NW_020311286.1_prot_XP_020708333.1_50 http://www.genome.jp/kegg-bin/show_
↳pathway?dme04341/dme:Dmel_CG6963%09red/dme:Dmel_CG6054%09red
Hedgehog signaling pathway      PANTHER      P00025      6          13          3.6166668094e-10          9.
↳94583372585e-09          lcl|NW_020311286.1_prot_XP_025602279.1_78|lcl|NW_020311286.1_
↳prot_XP_025602289.1_76|lcl|NW_020311286.1_prot_XP_025602264.1_79|lcl|NW_020311285.1_
↳prot_XP_012256108.1_32|lcl|NW_020311285.1_prot_XP_012256109.1_33|lcl|NW_020311286.1_
↳prot_XP_012256943.1_77          http://www.pantherdb.org/pathway/pathwayDiagram.jsp?
↳catAccession=P00025
Signaling by NOTCH2      Reactome          R-DME-1980145      3          8          2.00259649553e-05_
↳0.000275357018136          lcl|NW_020311285.1_prot_XP_012256118.1_28|lcl|NW_
↳020311285.1_prot_XP_012256117.1_27|lcl|NW_020311285.1_prot_XP_012256119.1_26
↳http://www.reactome.org/cgi-bin/eventbrowser_st_id?ST_ID=R-DME-1980145
```

Contact us.

## 2.17 KOBAS on the ARS Ceres HPC

### 2.17.1 About Ceres/Scinet

- The Scinet VRSC has installed KOBAS for ARS use.
- For general information on Scinet/Ceres, how to access it, and how to use it, visit <https://usda-ars-gbru.github.io/scinet-site/>.

### 2.17.2 Getting the databases

To run the tool you need some public data in your working directory. The files can be downloaded directly from the [KOBAS homepage](#). These directories are also available on Ceres/Scinet under `/reference/data/iplant/2019-09-16/kobas`. This directory will need to be **copied** to your working directory, since KOBAS requires write-access to the subdirectories within this folder.

- 1) `seq_pep/`: species-specific BLAST databases used by KOBAS

```
cp -r /reference/data/iplant/2019-09-16/kobas/seq_pep/ .
chmod -R 777 seq_pep
```

- 2) `sqlite3`: species-specific annotation databases used by KOBAS

```
cp -r /reference/data/iplant/2019-09-16/kobas/sqlite3/ .
chmod -R 777 sqlite3
```



**Note:** The above commands should result in two directories (seq\_pep and sqlite3) each containing many files. There is no need to unzip the .gz files.

---

## 2.17.3 Running KOBAS on Ceres

---

### Running programs on Ceres/Scinet

- You'll need to run KOBAS either in interactive mode or batch mode.
  - For interactive mode, use the *salloc* command.
  - For batch mode, you'll need to write a batch job submission bash script.
- 

### Running KOBAS in interactive mode

#### Loading the module

The Scinet VRSC has installed the KOBAS program. To load the module in interactive mode, run the command

```
module load agbase
```

---

#### Getting the Help and Usage Statement

```
kobas -h
```

---

See *Help and Usage Statement*

KOBAS can perform two tasks: - annotate (-a) - identify (enrichment) (-g) KOBAS can also run both task with a single command (-j).

#### Annotate Example Command - interactive mode

```
kobas -a -i AROS_10.faa -s dme -t fasta:pro -o kobas_output -k /project/nal_genomics/  
↪monica.poelchau
```

---

#### Annotate Example Command - batch mode

---

### Running programs on Ceres/Scinet in batch mode

- Before using batch mode, you should review Ceres/Scinet's documentation first, and decide what queue you'll want to use. See <https://usda-ars-gbru.github.io/scinet-site/guide/ceres/>.
- 

**Example batch job submission bash script (e.g. kobas-job.sh):**

```
#!/bin/bash
module load agbase
kobas -a -i AROS_10.faa -s dme -t fasta:pro -o kobas_output -k /project/nal_genomics/
↪monica.poelchau
```

### Submitting the batch job:

```
SBATCH kobas-job.sh
```

### Command Explained

**-a:** Tells KOBAS to run the ‘annotate’ process.

**-i AROS\_10.faa:** input file (peptide FASTA)

**-s dme:** Enter the species for the species of the sequences in your input file.

---

**Note:** If you don’t know the code for your species it can be found here: [https://www.kegg.jp/kegg/catalog/org\\_list.html](https://www.kegg.jp/kegg/catalog/org_list.html)

If your species of interest is not available then you should choose the code for the closest-related species available

---

**-t:** input file type; in this case, protein FASTA.

**-o kobas\_output:** name of output file

**-k /project/nal\_genomics/monica.poelchau:** KOBAS HOME. The path to kobas\_home, which is the parent directory of sqlite3/ and seq\_pep/.

For information on output files see *Understanding Your Results: Annotate*

### Identify Example Command - interactive mode

```
kobas -g -f kobas_output -b dme -k /project/nal_genomics/monica.poelchau -o ident_out
```

### Annotate Example Command - batch mode

---

#### Running programs on Ceres/Scinet in batch mode

- Before using batch mode, you should review Ceres/Scinet’s documentation first, and decide what queue you’ll want to use. See <https://usda-ars-gbru.github.io/scinet-site/guide/ceres/>.
- 

#### Example batch job submission bash script (e.g. kobas-job.sh):

```
#!/bin/bash
module load agbase
kobas -g -f kobas_output -b dme -k /project/nal_genomics/monica.poelchau -o ident_out
```

### Submitting the batch job:

```
sbatch kobas-job.sh
```

### Command Explained

**-g:** Tells KOBAS to run the 'identify' process.

**-f kobas\_output:** output file from KOBAS annotate

**-b dme:** background; enter the species code for the species of the sequences in your input file.

**Note:** If you don't know the code for your species it can be found here: [https://www.kegg.jp/kegg/catalog/org\\_list.html](https://www.kegg.jp/kegg/catalog/org_list.html)

If your species of interest is not available then you should choose the code for the closest-related species available

**-o ident\_out:** name of output file

For information on outputs see *Understanding Your Results: Identify*

### Annotate and Identify Pipeline Example Command - interactive mode

```
kobas -j -i AROS_10.faa -s dme -t fasta:pro -k /project/nal_genomics/monica.poelchau -
↳o kobas_output
```

### Annotate Example Command - batch mode

#### Running programs on Ceres/Scinet in batch mode

- Before using batch mode, you should review Ceres/Scinet's documentation first, and decide what queue you'll want to use. See <https://usda-ars-gbru.github.io/scinet-site/guide/ceres/>.

#### Example batch job submission bash script (e.g. kobas-job.sh):

```
#!/bin/bash
module load agbase
kobas -j -i AROS_10.faa -s dme -t fasta:pro -k /project/nal_genomics/monica.poelchau -
↳o kobas_output
```

#### Submitting the batch job:

```
sbatch kobas-job.sh
```

### Command Explained

**-j:** Tells KOBAS to run the 'annotate' process.

**-i AROS\_10.faa:** input file (peptide FASTA)

**-s dme:** Enter the species for the species of the sequences in your input file.

**Note:** If you don't know the code for your species it can be found here: [https://www.kegg.jp/kegg/catalog/org\\_list.html](https://www.kegg.jp/kegg/catalog/org_list.html)

If your species of interest is not available then you should choose the code for the closest-related species available

---

**-t:** input file type; in this case, protein FASTA.

**-o kobas\_output:** basename of output files

---

**Note:** This pipeline will automatically use the output of 'annotate' as the -f foreground input for 'identify'. This will also use your species option as the -b background input for 'identify'.

---

For more information on outputs see *Understanding Your Results: Annotate and Identify*

## 2.17.4 Understanding Your Results

### Annotate

If all goes well, you should get the following:

- **seq\_pep folder:** This folder contains the BLAST database files used in your analysis.
- **sqlite3 folder:** This folder contains the annotation database files used in your analysis
- **<species>.tsv:** This is the tab-delimited output from the BLAST search. It is unlikely that you will need to look at this file.
- **<output\_file\_name\_you\_provided>:** KOBAS-annotate generates a text file with the name you provide. It has two sections.

The first sections looks like this:

```
##dme      Drosophila melanogaster (fruit fly)
##Method: BLAST      Options: evalue <= 1e-05
##Summary:  87 succeed, 0 fail

#Query      Gene ID|Gene name|Hyperlink
lcl|NW_020311285.1_prot_XP_012256083.1_15   dme:Dmel_CG34349|Unc-13-4B|http://www.
↳genome.jp/dbget-bin/www_bget?dme:Dmel_CG34349
lcl|NW_020311286.1_prot_XP_020708336.1_46   dme:Dmel_CG6963|gish|http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG6963
lcl|NW_020311285.1_prot_XP_020707987.1_39   dme:Dmel_CG30403||http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG30403
```

The second section follows a dashed line and looks like this:

```
-----
////
Query:          lcl|NW_020311285.1_prot_XP_012256083.1_15
Gene:           dme:Dmel_CG34349      Unc-13-4B
Entrez Gene ID: 43002
////
Query:          lcl|NW_020311286.1_prot_XP_020708336.1_46
Gene:           dme:Dmel_CG6963  gish
```

(continues on next page)

(continued from previous page)

```

Entrez Gene ID:          49701
Pathway:                 Hedgehog signaling pathway - fly          KEGG PATHWAY
↳dme04341
////
Query:                  lcl|NW_020311285.1_prot_XP_020707987.1_39
Gene:                   dme:Dmel_CG30403
Entrez Gene ID:         246595
////
Query:                  lcl|NW_020311285.1_prot_XP_020707989.1_40
Gene:                   dme:Dmel_CG6148 Past1
Entrez Gene ID:         41569
Pathway:                Endocytosis           KEGG PATHWAY      dme04144
                       Hemostasis           Reactome         R-DME-109582
                       Factors involved in megakaryocyte development and
↳platelet production   Reactome         R-DME-98323

```

## Identify

If all goes well, you should get the following:

- **sqlite3 folder:** This folder contains the annotation database files used in your analysis
- **<output\_file\_name\_you\_provided>:** KOBAS identify generates a text file with the name you provide.

```

##Databases: PANTHER, KEGG PATHWAY, Reactome, BioCyc
##Statistical test method: hypergeometric test / Fisher's exact test
##FDR correction method: Benjamini and Hochberg

#Term          Database          ID          Input number    Background number    P-Value
↳Corrected P-Value    Input    Hyperlink
Hedgehog signaling pathway - fly    KEGG PATHWAY    dme04341        12        33        3.
↳20002656734e-18        1.76001461204e-16        lcl|NW_020311286.1_prot_XP_012256678.
↳1_51|lcl|NW_020311286.1_prot_XP_025602973.1_48|lcl|NW_020311286.1_prot_XP_012256683.
↳1_52|lcl|NW_020311286.1_prot_XP_012256679.1_55|lcl|NW_020311286.1_prot_XP_012256674.
↳1_54|lcl|NW_020311286.1_prot_XP_020708336.1_46|lcl|NW_020311285.1_prot_XP_012256108.
↳1_32|lcl|NW_020311286.1_prot_XP_012256682.1_53|lcl|NW_020311286.1_prot_XP_025603025.
↳1_47|lcl|NW_020311286.1_prot_XP_020708334.1_49|lcl|NW_020311285.1_prot_XP_012256109.
↳1_33|lcl|NW_020311286.1_prot_XP_020708333.1_50 http://www.genome.jp/kegg-bin/show_
↳pathway?dme04341/dme:Dmel_CG6963%09red/dme:Dmel_CG6054%09red
Hedgehog signaling pathway    PANTHER    P00025    6        13        3.6166668094e-10    9.
↳94583372585e-09        lcl|NW_020311286.1_prot_XP_025602279.1_78|lcl|NW_020311286.1_
↳prot_XP_025602289.1_76|lcl|NW_020311286.1_prot_XP_025602264.1_79|lcl|NW_020311285.1_
↳prot_XP_012256108.1_32|lcl|NW_020311285.1_prot_XP_012256109.1_33|lcl|NW_020311286.1_
↳prot_XP_012256943.1_77        http://www.pantherdb.org/pathway/pathwayDiagram.jsp?
↳catAccession=P00025
Signaling by NOTCH2    Reactome          R-DME-1980145    3        8        2.00259649553e-05
↳0.000275357018136        lcl|NW_020311285.1_prot_XP_012256118.1_28|lcl|NW_
↳020311285.1_prot_XP_012256117.1_27|lcl|NW_020311285.1_prot_XP_012256119.1_26
↳http://www.reactome.org/cgi-bin/eventbrowser_st_id?ST_ID=R-DME-1980145

```

## Annotate and Identify Pipeline

If all goes well, you should get the following:

- **logs folder:** This folder contains the 'conder\_stderr' and 'conder\_stdout' files. The files record feedback, progress and, importantly, any errors the app encountered during the analysis. You won't normally need to look at these but they are very helpful in figuring out what may have happened if your output doesn't look like you expected.
- **sqlite3 folder:** This folder contains the annotation database files used in your analysis
- **seq\_pep folder:** This folder contains the BLAST database files used in your analysis.
- **<species>.tsv:** This is the tab-delimited output from the BLAST search. It is unlikely that you will need to look at this file.
- **<basename>\_annotate\_out.txt:** KOBAS annotate generates a text file with the name you provide. It has two sections.

The first sections looks like this:

```
##dme      Drosophila melanogaster (fruit fly)
##Method: BLAST      Options: evaluate <= 1e-05
##Summary: 87 succeed, 0 fail

#Query      Gene ID|Gene name|Hyperlink
lcl|NW_020311285.1_prot_XP_012256083.1_15  dme:Dmel_CG34349|Unc-13-4B|http://www.
↳genome.jp/dbget-bin/www_bget?dme:Dmel_CG34349
lcl|NW_020311286.1_prot_XP_020708336.1_46  dme:Dmel_CG6963|gish|http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG6963
lcl|NW_020311285.1_prot_XP_020707987.1_39  dme:Dmel_CG30403||http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG30403
```

The second section follows a dashed line and looks like this:

```
-----
////
Query:      lcl|NW_020311285.1_prot_XP_012256083.1_15
Gene:      dme:Dmel_CG34349      Unc-13-4B
Entrez Gene ID:      43002
////
Query:      lcl|NW_020311286.1_prot_XP_020708336.1_46
Gene:      dme:Dmel_CG6963 gish
Entrez Gene ID:      49701
Pathway:      Hedgehog signaling pathway - fly      KEGG PATHWAY      ↳
↳dme04341
////
Query:      lcl|NW_020311285.1_prot_XP_020707987.1_39
Gene:      dme:Dmel_CG30403
Entrez Gene ID:      246595
////
Query:      lcl|NW_020311285.1_prot_XP_020707989.1_40
Gene:      dme:Dmel_CG6148 Past1
Entrez Gene ID:      41569
Pathway:      Endocytosis      KEGG PATHWAY      dme04144
      Hemostasis      Reactome      R-DME-109582
      Factors involved in megakaryocyte development and↳
↳platelet production      Reactome      R-DME-98323mel_CG6963
lcl|NW_020311285.1_prot_XP_020707987.1_39  dme:Dmel_CG30403||http://www.genome.jp/
↳dbget-bin/www_bget?dme:Dmel_CG30403
```

- **<basename>\_identify\_out.txt:** KOBAS identify generates a text file with the name you provide.

```

##Databases: PANTHER, KEGG PATHWAY, Reactome, BioCyc
##Statistical test method: hypergeometric test / Fisher's exact test
##FDR correction method: Benjamini and Hochberg

#Term          Database      ID          Input number  Background number  P-Value_
↪Corrected P-Value      Input      Hyperlink
Hedgehog signaling pathway - fly      KEGG PATHWAY      dme04341      12      33      3.
↪20002656734e-18      1.76001461204e-16      lcl|NW_020311286.1_prot_XP_012256678.
↪1_51|lcl|NW_020311286.1_prot_XP_025602973.1_48|lcl|NW_020311286.1_prot_XP_012256683.
↪1_52|lcl|NW_020311286.1_prot_XP_012256679.1_55|lcl|NW_020311286.1_prot_XP_012256674.
↪1_54|lcl|NW_020311286.1_prot_XP_020708336.1_46|lcl|NW_020311285.1_prot_XP_012256108.
↪1_32|lcl|NW_020311286.1_prot_XP_012256682.1_53|lcl|NW_020311286.1_prot_XP_025603025.
↪1_47|lcl|NW_020311286.1_prot_XP_020708334.1_49|lcl|NW_020311285.1_prot_XP_012256109.
↪1_33|lcl|NW_020311286.1_prot_XP_020708333.1_50 http://www.genome.jp/kegg-bin/show_
↪pathway?dme04341/dme:Dmel_CG6963%09red/dme:Dmel_CG6054%09red
Hedgehog signaling pathway      PANTHER      P00025      6      13      3.6166668094e-10      9.
↪94583372585e-09      lcl|NW_020311286.1_prot_XP_025602279.1_78|lcl|NW_020311286.1_
↪prot_XP_025602289.1_76|lcl|NW_020311286.1_prot_XP_025602264.1_79|lcl|NW_020311285.1_
↪prot_XP_012256108.1_32|lcl|NW_020311285.1_prot_XP_012256109.1_33|lcl|NW_020311286.1_
↪prot_XP_012256943.1_77      http://www.pantherdb.org/pathway/pathwayDiagram.jsp?
↪catAccession=P00025
Signaling by NOTCH2      Reactome      R-DME-1980145      3      8      2.00259649553e-05_
↪0.000275357018136      lcl|NW_020311285.1_prot_XP_012256118.1_28|lcl|NW_
↪020311285.1_prot_XP_012256117.1_27|lcl|NW_020311285.1_prot_XP_012256119.1_26
↪http://www.reactome.org/cgi-bin/eventbrowser_st_id?ST_ID=R-DME-1980145

```

Contact us.