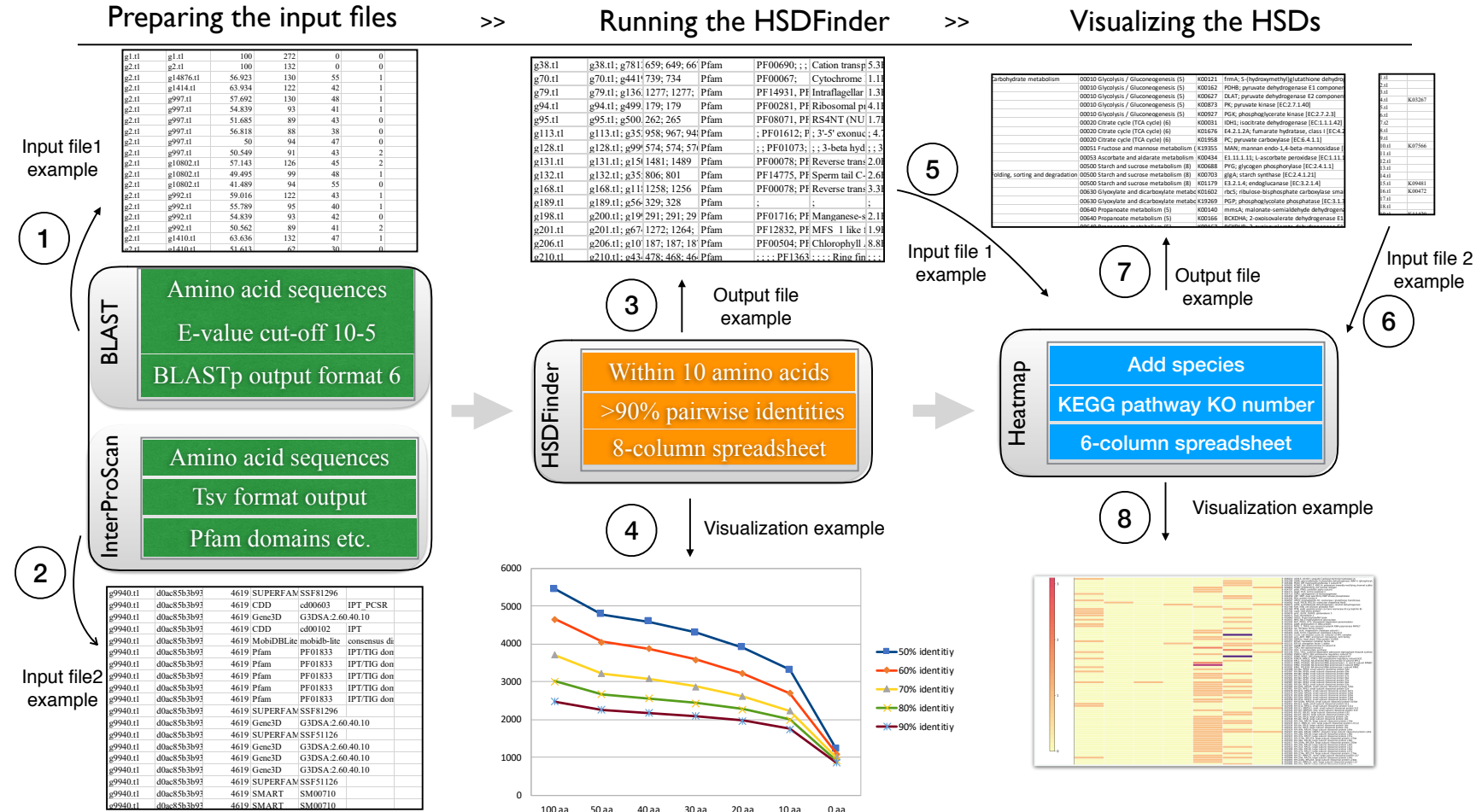


HSDFinder Online server Tutorial: Identification of highly similar duplicates in eukaryotic genomes

1. Upload a protein BLAST search result file of your genome in tab-separated values (tsv) format as the first input file (**File 1**) of HSDFinder.
2. Upload a InterProScan search result file of your genome in tab-separated values (tsv) format as the second input file (**File 2**) of HSDFinder.
3. Yielding the output of HSDFinder with three personalized options.
4. Visualizing the HSDFinder outputs via the Excel tools (optional).
5. Upload the results of HSDFinder from your respective genomes.
6. Upload a gene list with KO annotation from KEGG database.
7. The output files of the online Heatmap Visualization tool.
8. The heatmap of HSDs levels across species.

Workflow of HSDFinder (8 steps)



Preparing the appropriate input files.

First, before running HSDFinder to acquire the HSDs of your interest genome, there are two spreadsheets (**File 1** and **File 2**) in tab-separated values (tsv) format shall be prepared as input files.

Step 1: Upload a protein BLAST search result file of your genome in tab-separated values (tsv) format as the first input file (**File 1**) of HSDFinder.

- Using BLASTP search of the amino acid sequences of genome against themselves (The BLASTP running parameters: E-value cut-off 10⁻⁵, BLASTP output format 6)
- The BLAST results should be 12-column spreadsheets including the key information from query name to percentage identity etc. For the **File 1 example** please check the red arrows below.
- The BLAST Package can be found via <https://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/LATEST/>
Please download the appropriate tools based on your computer systems (Windows, MacOS or Linux).

Here, this tutorial will go over how to acquire local BLASTP search results via an example FASTA file (The FASTA file can be found via the link <https://github.com/zx0223winner/HSDFinder/blob/master/Example.fasta>). You can ignore this step if you know how to acquire the BLASTP searching results by other ways.

```
$ head Example.fasta
```

```
>g735.tl
```

```
MTMRSPTAAAAALALLLLACTAGIQCARPLLQSEDEVDCSNPPQAPVDLGNASAFVVLGHEVTTGAGVASRILGDLG  
VSPGNSVTGQPDMLEGSVIQKANGASAAGIRDLGIAYNDAAGRTVCPILVAGELGGMTLYPGLYKSSSGLEITDADLT  
SGDGIFIFQMATTFLMTKTMKVTLTNGAQAKNIFWQVGTSATLMDKSVLYGTILADQSITSGTGAVIHGRALARIASVT  
MESAVFSLPAE
```

```
>g741.tl
```

```
MTMRSPTVAAALVLLLLACTAGVQCARPLLQSEDAVDCSNPPQAPVDLRNATTFVVLGHEVTTGAGVASRILSDLG  
VSPGSSVTGQPDMLEGSVIQKANGASAAGIRDLGIAYNDAAGRTVCPILVAGELGGMTLYPGLYKSSSGLEITDANLTLS  
GDGVFIFQMATTFLMTKTMKVTLTNGAQAKNIFWQVGTSATLMDKSVLYGTILADQSITSGTGAVIHGRALARIASVT  
MESAVFSLPAE
```

Build a Database via the example FASTA file. Check command line below:

```
$ makeblastdb -in Example.fasta -dbtype prot Example_database_name
```

Using BLASTP search option to blast the amino acid sequences against themselves:

```
$ blastp -query IBP.FASTA -db Example_database_name -out File1_example_BALSTP_result.txt -evalue 1e-5 -  
outfmt 6
```

This will give a BLAST result file form by 12-column spreadsheets including the key information from query name to percentage identity etc. Check the File 1 example below on the left-hand side.

Run

Text Input

Two spreadsheets in tab-separated values (tsv) format shall be prepared as input files.

The first spreadsheet is from a protein BLAST search of the genome genes against themselves (E-value cut-off 10-5, BLASTp output format 6).

g735.t1	g735.t1	100.000	744	0	0	1	744	1
744	0.0	1375						
g735.t1	g741.t1	96.237	744	28	0	1	744	1
744	0.0	1219						
g735.t1	g8053.t1	90.196	51	3	2	6	55	3
7.50e-13	65.8							
g735.t1	g7171.t1	77.632	608	121	13	144	740	147
750	3.98e-100	355						
g735.t1	g11305.t1	97.500	40	1	0	17	56	14
53	5.80e-14	69.4						
g741.t1	g741.t1	100.000	744	0	0	1	744	1
744	0.0	1375						
g8053.t1	g8053.t1	100.000	747	0	0	1	747	1

See File Examples

File 1

File 2

The second spreadsheet is acquired from InterProScan which is an automatically software providing the protein signatures such as Pfam domain.

g735.t1	c82510c09b797ecced03c40f4da02ffb	247
Pfam	PF11999 Protein of unknown function (DUF3494)	57 241 2.2E-47 T 15-11-2019
IPR021884	Ice-binding protein-like	
g735.t1	c82510c09b797ecced03c40f4da02ffb	247
ProSiteProfiles	PS51257 Prokaryotic membrane lipoprotein lipid attachment site profile. 1	19 5.0 T 15-11-2019
g741.t1	8cf52deba53cb877fbd0af222ed48ce3	247
ProSiteProfiles	PS51257 Prokaryotic membrane lipoprotein lipid attachment site profile. 1	19 5.0 T 15-11-2019
g741.t1	8cf52deba53cb877fbd0af222ed48ce3	247

Users can either paste or upload a BLASTP search result as the input of file 1.

Or Upload

File1 (tsv.)

The BLAST results should be 12-column spreadsheets including the key information from query name to percentage identity etc.(see more at web [FAQ](#))

Choose File

File1_example.txt

File2 (tsv.)

The output file of InterProScan is tab-separated values (tsv) format in default.

Choose File

File2_example.txt

Step 2: Upload a InterProScan search result file of your genome in tab-separated values (tsv) format as the second input file (**File 2**) of HSDFinder.

- InterProScan which is easy-to-use software providing the protein signatures such as Pfam domain. The latest InterProScan software can be found via the link <https://www.ebi.ac.uk/interpro/download/>
- After finishing downloading, you can test the InterProScan software locally with your interest genome (protein sequence in FASTA file). Check the commands below: More usage instructions can be found via the link <https://interproscan-docs.readthedocs.io/en/latest/HowToRun.html>

```
./interproscan.sh -i proteins_of_your_genome.fasta -f tsv -dp
```
- The output file of InterProScan is tab-separated values (tsv) format in default. Please find the example (File 2) on the right-hand side (see below).

Run

Text Input

Two spreadsheets in tab-separated values (tsv) format shall be prepared as input files.

The first spreadsheet is from a protein BLAST search of the genome genes against themselves (E-value cut-off 10-5, BLASTp output format 6).

g735.t1	g735.t1	100.000	744	0	0	1	744	1
744	0.0	1375						
g735.t1	g741.t1	96.237	744	28	0	1	744	1
744	0.0	1219						
g735.t1	g8053.t1	90.196	51	3	2	6	55	3
7.50e-13	65.8							
g735.t1	g7171.t1	77.632	608	121	13	144	740	147
750	3.98e-100	355						
g735.t1	g11305.t1	97.500	40	1	0	17	56	14
53	5.80e-14	69.4						
g741.t1	g741.t1	100.000	744	0	0	1	744	1
744	0.0	1375						
g8053.t1	g8053.t1	100.000	747	0	0	1	747	1

See File Examples

File 1

File 2

The second spreadsheet is acquired from InterProScan which is an automatically software providing the protein signatures such as Pfam domain.

g735.t1	c82510c09b797ecced03c40f4da02ffb	247
Pfam	PF11999 Protein of unknown function (DUF3494)	57 241 2.2E-47 T 15-11-2019
IPR021884	Ice-binding protein-like	
g735.t1	c82510c09b797ecced03c40f4da02ffb	247
ProSiteProfiles	PS51257 Prokaryotic membrane lipoprotein lipid attachment site profile. 1	19 5.0 T 15-11-2019
g741.t1	8cf52deba53cb877fbd0af222ed48ce3	247
ProSiteProfiles	PS51257 Prokaryotic membrane lipoprotein lipid attachment site profile. 1	19 5.0 T 15-11-2019
g741.t1	8cf52deba53cb877fbd0af222ed48ce3	247

Runing the HSDFinder

Before clicking the submission button, there are three personalized options designing for HSDFinder (Amino acid pairwise identity, amino acid length difference, and protein function database)

Step 3: Yielding the output of HSDFinder with three personalized options.

- The HSDFinder is set default to filter those with near-identical protein lengths (within 10 amino acids) and >90% pairwise identities. By setting such a strict cut-off, user should be able to screen out most of highly identical duplicates, but there are few genuine duplicates might be falling out of the threshold.

Amino acid pairwise identities:

90%

?

Protein function databases:

Pfam

?

Amino acid lengths (aa):

10

?

Click

Submit

- Nevertheless, the users always have an option to try different parameters from 30% to 100% identity or from within 10 aa to 100 aa variances to acquire the duplicates in need.

But the user should be careful with the false-positive rate of HSDs will increase (not genuine duplicates) with larger amino acid variance and lower amino acid pairwise identity.

- The output of this step will be an 8-column spreadsheet integrating with the information of HSD identifier, gene copies number and Pfam domain.

Submit

Output:

HSDs results example (tsv. format)

```
g735.t1 g735.t1; g741.t1; g8053.t1 744; 744; 747 Pfam PF11999; PF11999; PF11999 Protein of unknown
function (DUF3494); Protein of unknown function (DUF3494); Protein of unknown function (DUF3494) 2.2E-47; 7.8E-
47; 2.5E-47 IPR021884; IPR021884; IPR021884 Ice-binding protein-like ; Ice-binding protein-like ; Ice-binding
protein-like
```

Output:

HSDFinder generates one output files: 8-column spreadsheet integrating with the information of HSD identifier, gene copies number and Pfam domain.

Example of the 8-column spreadsheet:

```
g735.t1 g735.t1; g741.t1; g8053.t1 744; 744; 747 Pfam PF11999; PF11999; PF11999
Protein of unknown function (DUF3494); Protein of unknown function (DUF3494); Protein of unknown
function (DUF3494) 2.2E-47; 7.8E-47; 2.5E-47 IPR021884; IPR021884; IPR021884 Ice-binding
protein-like; Ice-binding protein-like ; Ice-binding protein-like
```

Column explanation:

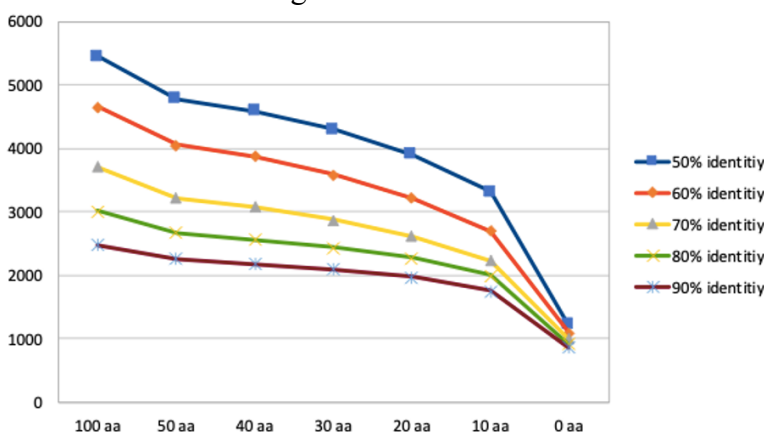
1. Highly Similar Duplicates (HSDs) identifiers: The first gene model of the duplicate gene copies is used as the HSD identifiers in default.
2. Duplicate gene copies (within 10 amino acids, $\geq 90\%$ pairwise identities): g735.t1; g741.t1; g8053.t1
3. Amino acid length of duplicate gene copies (aa): 744; 744; 747
4. Pfam identifier: PF11999; PF11999; PF11999
5. Pfam Description: Protein of unknown function (DUF3494); Protein of unknown function (DUF3494); Protein of unknown function (DUF3494)
6. InterPro Entry Identifier: IPR021884; IPR021884; IPR021884
7. InterPro Entry Description: Ice-binding protein-like ; Ice-binding protein-like ; Ice-binding protein-like

Step 4: Visualizing the HSDFinder outputs via the Excel tools (Optional)

Additionally, the user can conveniently set different values to create a trendline graph of the gene copies numbers under different criteria. Check the example we used below. The genome datasets are from a psychrophilic green alga *Chlamydomonas* sp. UWO241 (NCBI BioProject: PRJNA547753).

	50% identity	60% identity	70% identity	80% identity	90% identity
100 aa	5453	4654	3719	3012	2476
50 aa	4783	4061	3228	2672	2257
40 aa	4594	3879	3088	2565	2181
30 aa	4310	3596	2875	2438	2091
20 aa	3912	3224	2614	2270	1973
10 aa	3320	2699	2234	2012	1339
0 aa	1223	1097	999	926	859

The table of total gene copies number at different thresholds of amino acid pairwise identity and deduced amino acid length.



The line graph of duplicates set to different thresholds of amino acid pairwise identity and deduced amino acid length. The X-axis indicates the deduced amino acid length (aa) of each duplicate, the Y-axis tells the number of gene copies.

Visualizing the HSDs across species

If you have more than one species needing to compare the levels of HSDs in each genome, the following option of Online Heatmap tool is a great choice. User can use the online heatmap plotting option to visualize the HSDs results in different KEGG pathway category.

Step 5: Upload the results of HSDFinder from your respective genomes.

There will be two files needed to plot the heatmap for each species. First input file is the output of your interest species after running the HSDFinder (the step 3), file examples are given to guide the appropriate input file (see the File 1 example below in the image).

Step 6: Upload a gene list with KO annotation from KEGG database.

- The second file is retrieved from the KEGG database documented the correlation of KO accession with each gene model identifier.
- Using the Ghost KOALA analysis tool of KEGG to acquire the KO annotation file of your genome. <https://www.kegg.jp/ghostkoala/> File examples are given to guide the appropriate input file (see the File 2 example below).
- Fill in the organism name with respective to your species, this is the identifier to compare levels of HSDs in different species.

Visualization

To comparative analyze the HSDs across different species, we developed an online heat map plotting option to visualize the HSDs results in different KEGG pathway category.

Creat Heatmap

HSD File ? File1 Example
 no file selected

Gene list with KO annotation ? File2 Example
 no file selected

Organism name ? e.g. UW0241

Click and drag the file to import

Click and drag the file to import

To add more species, you can simply touch the **+add species** button and fill in the respective files (Repeat step 5 and 6). To best present your visualization result, at least two species shall be submitted. You can still visualize the result via a single species though, if you simply want to try.

Creat Heatmap

HSD File ? File1 Example
 Heatmap_example1.txt

Gene list with KO annotation ? File2 Example
 Heatmap_example2.txt

Organism name ? green algae 1

HSD File
 no file selected

Gene list with KO annotation
 no file selected

Organism name
green algae 2

Step 7: The output files of the online Heatmap Visualization tool.

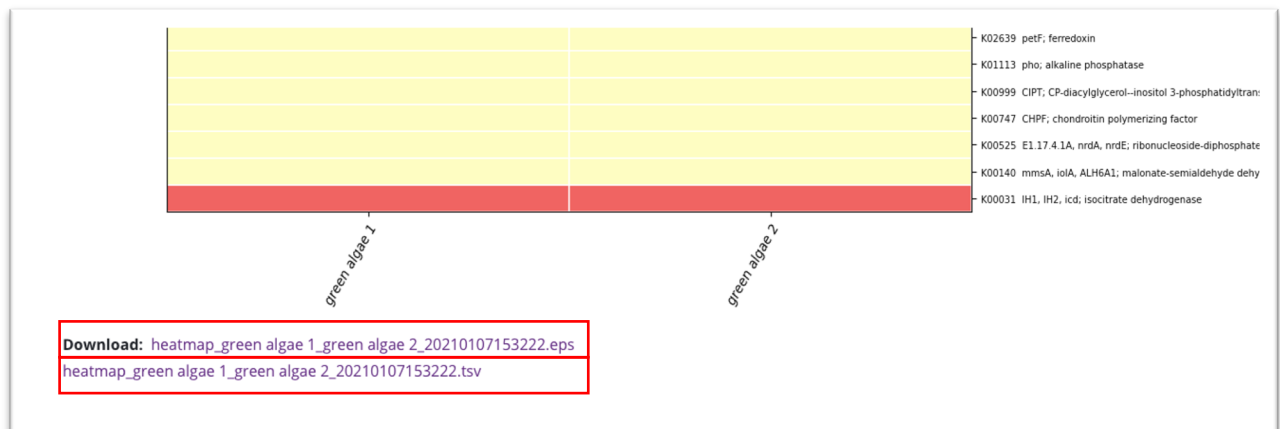
User can feel free to set the Figure size via the options of row and col. After clicking the Create Heatmap button. There will be a heatmap below as well as two documents generated. One is the high resolution heatmap in eps. format (image). Another one is the HSDs content details in tsv. format (text), which user can be conveniently used for future analysis. (This spreadsheet file includes information about HSDs under different KEGG function category such as carbohydrate metabolism, energy metabolism, translation etc.)

Figure Size: row col

Click

Creat Heatmap

Once the input files have been submitted, the HSDs numbers for each species will be displayed in a heatmap under different KEGG function category. On the left side, the color bar indicates a broad category of HSDs who have pathway function matches, such as carbohydrate metabolism, energy metabolism, translation etc. The color for the matrix indicates the number of HSDs across species.



Step 8: The heatmap of HSDs levels across species.

Once the input files have been submitted, the HSDs numbers for each species will be displayed in a heatmap under different KEGG function category. On the left side, the color bar indicates a broad category of HSDs who have pathway function matches, such as carbohydrate metabolism, energy metabolism, translation etc. The color for the matrix indicates the number of HSDs across species.

