

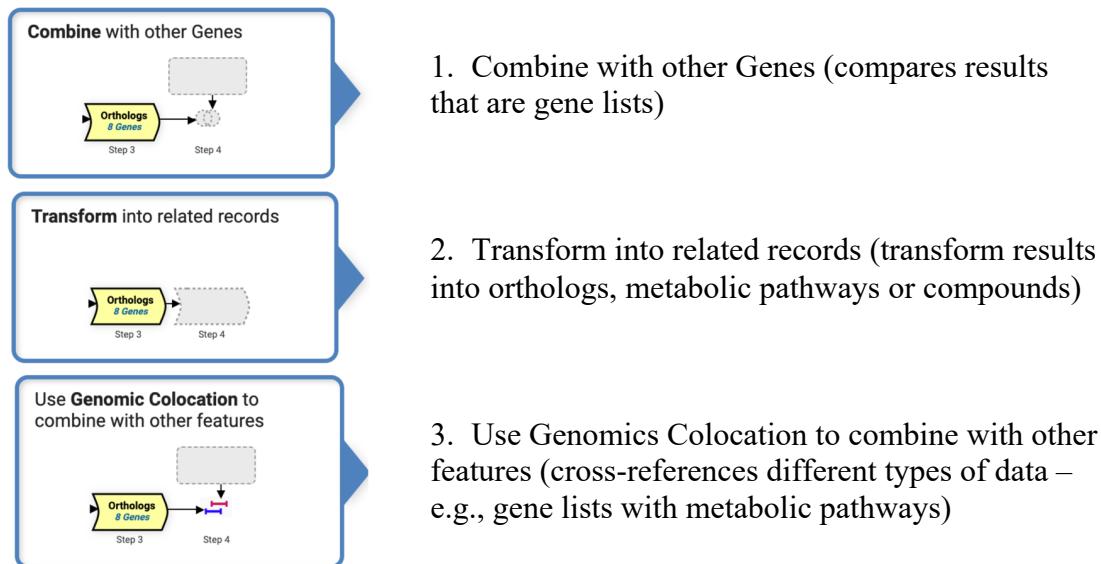
Advanced Search Strategies

Learning objectives:

- Use site search and other types of searches to create a multi-step query across different types of records and genomes.
- Create a multi-step query, save it to your account and share a link in Zoom chat.

The strategy system offers over 100 structured searches that can be combined to produce multi-step strategies. Searches can be deployed from the site search, or ‘Search For...’ menu on the home page and from the ‘Searches’ dropdown menu in the header of every page. Searches listed under Genes will return a list of gene IDs, while searches listed under ‘SNPs’ or ‘Metabolic Pathways’ will return record IDs representing SNPs, or metabolic pathways, respectively, etc.

The searches can be combined via three major approaches:



Strategy steps are connected via the Boolean operators that can intersect, unite, or subtract similar records (e.g., gene lists) and cross-reference different types of data via the genomic colocation option (blue & red parallel lines symbol). Steps can be masked off from the strategy with the help of the “ignore” Boolean operators.

Revise as a boolean operation	
<input checked="" type="radio"/> ● ○ 1 INTERSECT 2	<input type="radio"/> ○ ○ 1 UNION 2
<input type="radio"/> ○ ○ 1 MINUS 2	<input type="radio"/> ○ ○ 2 MINUS 1
Revise as a span operation	
<input type="radio"/> ○ ■ 1 RELATIVE TO 2, using genomic colocation	
Ignore one of the inputs	
<input type="radio"/> ○ ○ IGNORE 2	<input type="radio"/> ○ ○ IGNORE 1
<input type="button" value="Revise"/>	

Malassezia infections are neglected fungal disease with worldwide distribution, including tropics. *M. restricta* is a pathogen that can cause skin disorders and is one of the most common fungal species found on human skin. *Malassezia* cannot produce fatty acids and relies on fatty acid uptake from external sources. Secreted lipases are thought to contribute to *Malassezia* pathogenicity. In this strategy we will identify secreted lipases in *M. restricta* KCTC 27527, cross-reference annotation with InterPro domain annotations and find orthologs of *M. restricta* genes in other species (e.g., *Cryptococcus* and *Candida*).

To build this strategy, use the following approach:

- **Use site search** to identify genes that have “lipase” annotation in *Malassezia restricta* KCTC 27527. This search identifies genes that have “lipase” annotation in several evidence fields.
- **Identify Genes by Signal peptide prediction.** This search returns genes predicted to have signal peptide.
- **Identify Genes with the “LIP” InterPro domain.** This search identifies genes with specific domain signature – secreted lipase (LIP).
- **Transform result into related records in *Cryptococcus* and *Candida* via Orthology.** Lipolytic enzymes have been demonstrated to play an important role in virulence of pathogenic fungi. In this case, we will transform a list of *M. restricta* KCTC 27527 genes into their orthologs in *Cryptococcus neoformans* H99 and *Candida albicans* SC5314.

The individual components of the steps listed above are outlined below:

- **Use site search** to identify genes that have “lipase” annotation in *Malassezia restricta* KCTC 27527
 1. Run site search for genes annotated with “lipase” and filter on Genes.
 2. Use Gene fields to filter your results as shown.
 3. Restrict your search to *M. restricta* KCTC 27527 genes.
 4. Export results as a search strategy.

Hint: if you are not getting the same results or some options are missing, check your organism preferences. If they are enabled, click on the toggle switch to deactivate.

The screenshot shows the BioPax search results for genes matching "lipase".

1. Filter results: A sidebar on the left shows a list of selected filters: EC descriptions and numbers, GO terms, InterPro domains, Notes from annotators, Orthologs, PDB chains, Phenotype, Preferred product description, Product descriptions, and User comments. The "InterPro domains" filter is highlighted with a red circle.

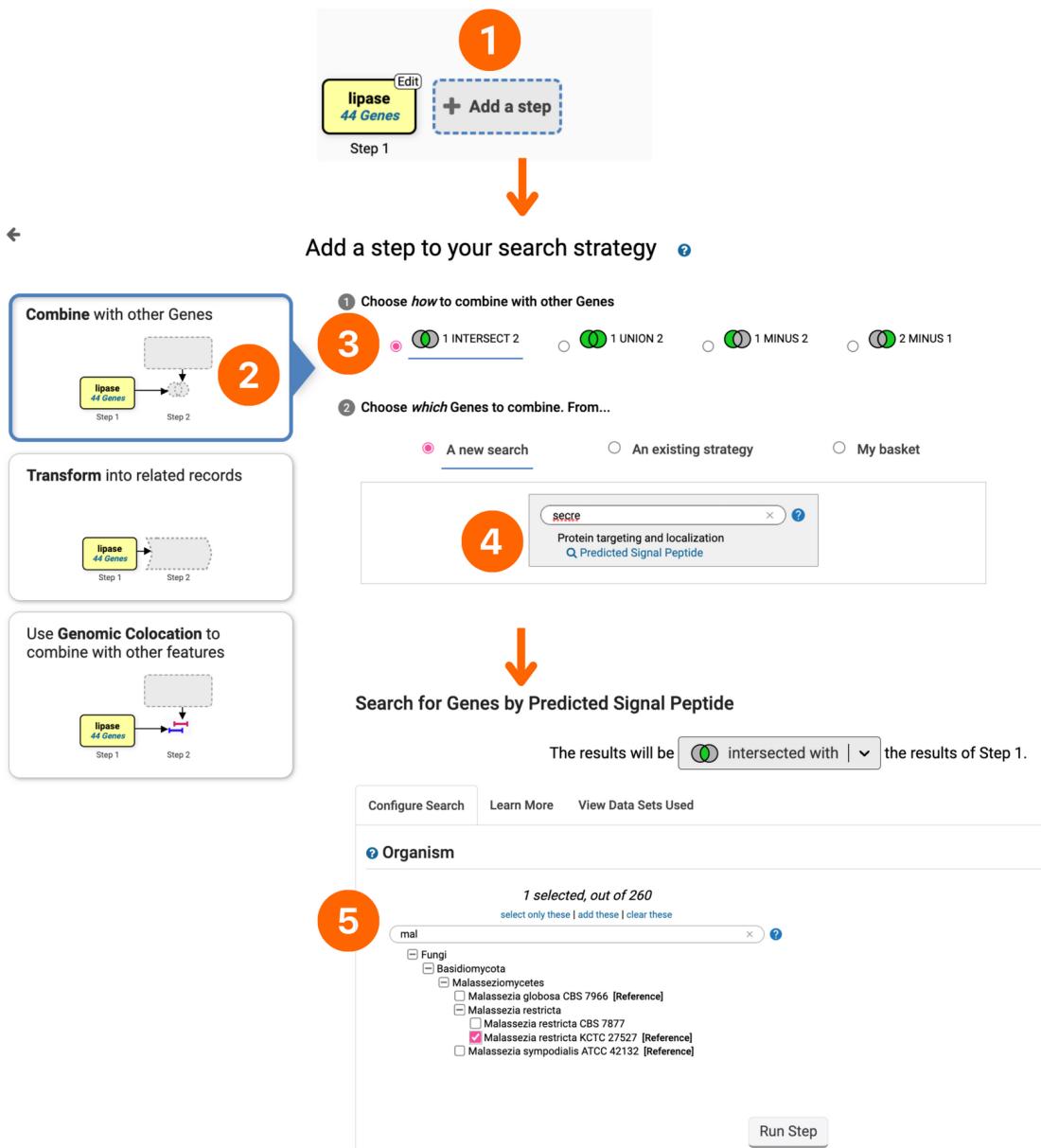
2. Filter Gene fields: A list of gene fields is shown, each with a count and a "Clear filter" button. The "InterPro domains" field is highlighted with a red circle.

3. Filter organisms: A tree view of organisms is shown, with "Fungi" expanded. "Basidiomycota" is selected, and "Malasseziomycetes" is expanded. "Malassezia" is selected, and "Malassezia restricta" is expanded. "KCTC 27527" is selected under "Malassezia restricta". The "Malassezia restricta" node is highlighted with a red circle.

4. Results: The main panel displays the results for "Gene - MRET_0019 lipase". It includes the gene ID, type (protein coding gene), organism (Malassezia restricta KCTC 27527), and a list of fields matched. There are 44 such entries in total.

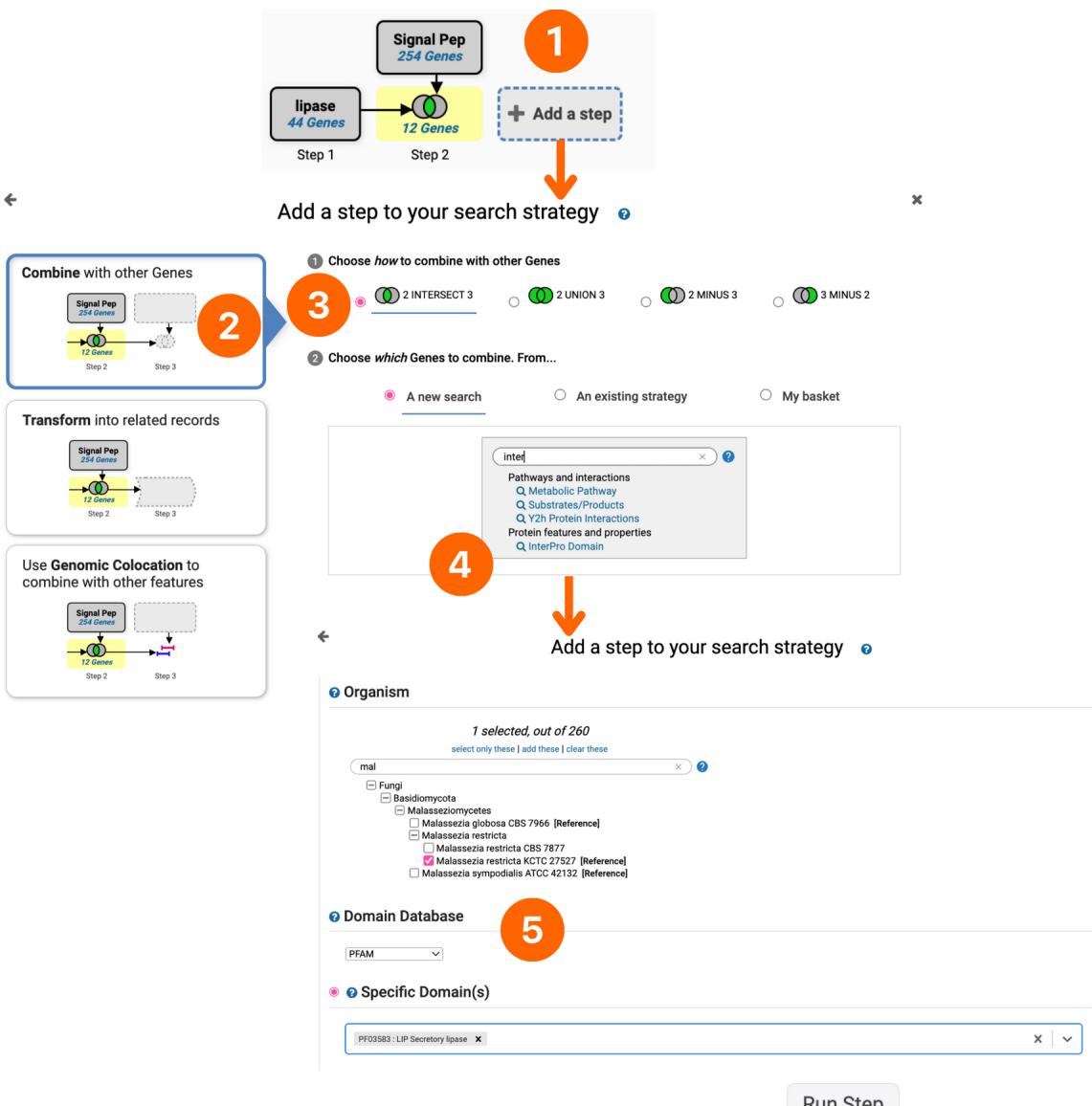
Export as a Search Strategy: A blue button in the top right corner allows users to download or mine their results.

- **Identify Genes also carrying predicted signal peptides.** This step will identify lipases that may be secreted.
 1. Click on the “Add step” button.
 2. Choose the “Combine with other genes” search.
 3. Choose to “intersect” your results with the previous step.
 4. Filter the available searches to identify the “Predicted Signal Peptide” search and click on the link in blue.
 5. Restrict the search to *M. restricta* KCTC 27527 and click on the “Run Step” button.



- **Identify Genes based on InterPro domain.** This search identifies genes annotated with specific domain signature – secreted lipase (LIP).

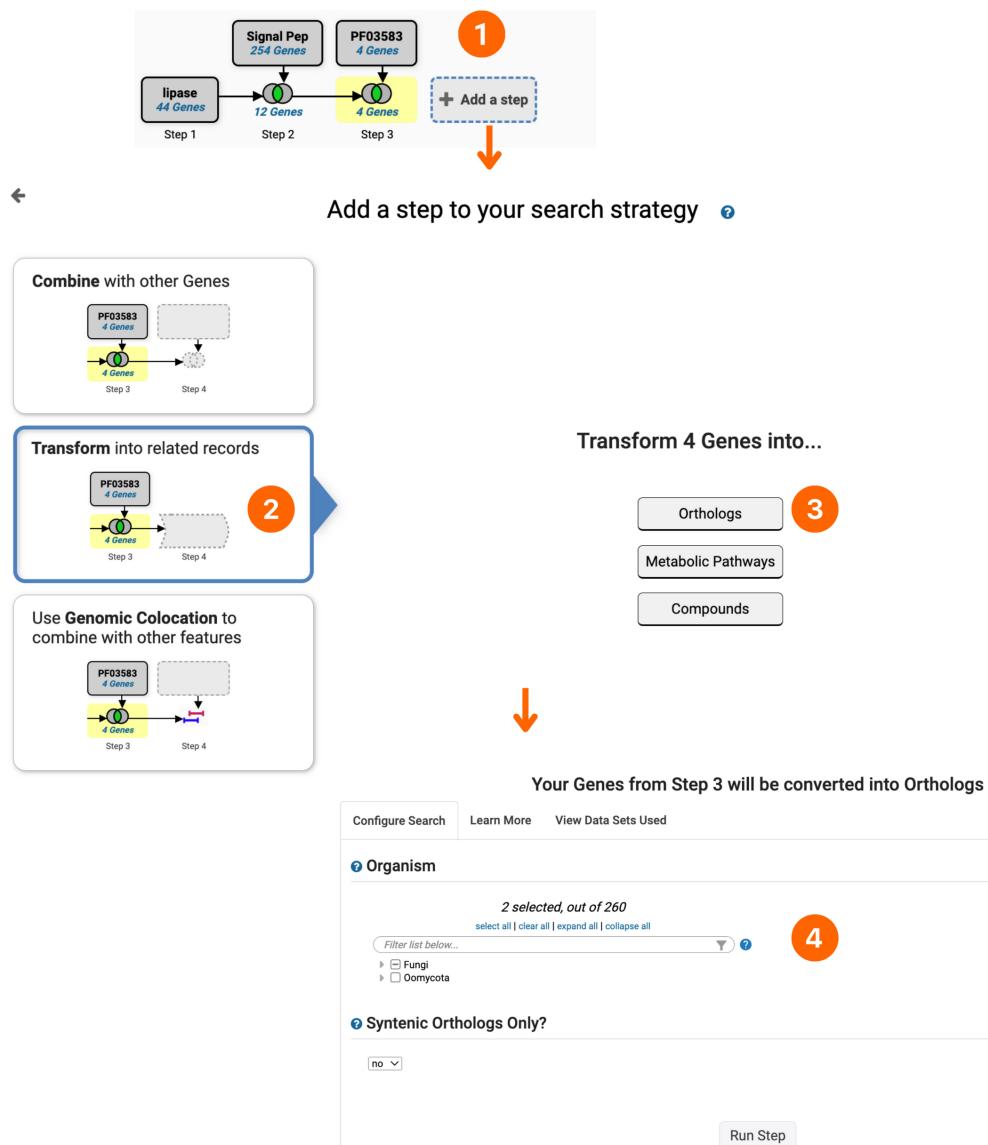
1. Click on the “Add step” button.
2. Choose the “Combine with other genes” search.
3. Choose to “intersect” your results with the previous step.
4. Filter the available searches by entering “inter” to quickly find the “InterPro domain” search and click on the link in blue.
5. Restrict the search to *M. restricta* KCTC 27527, select “Secretory lipase” domain (PF03583 : LIP Secretory lipase), and click on the “Run Step” button.



- Transform into related records in *Cryptococcus* and *Candida* via the Orthologs search.

The transformation step performed here will convert a list of genes in one organism (*M. restricta*) to their orthologs in another organisms. This search is particularly useful if you are working with a poorly annotated genome and want to take advantage of annotations from another, better annotated, genome. In this exercise, we will practice finding orthologs in *Cryptococcus neoformans* H99 and *Candida albicans* SC5314. Orthologs are predicted by the OrthoMCL algorithm, which clusters proteins into ortholog groups based on BLAST similarity across at 150 genomes that span the tree of life.

1. Click on the “Add step” button.
2. Choose the “Transform into related records” search.
3. Click on the “Orthologs” button to deploy the search.
4. Restrict the orthologs search to *Cryptococcus neoformans* H99 and *Candida albicans* SC5314 and click on the “Run Step” button.



Examine your results. Do they make sense?

ACE1 *

Step 1: lipase (44 Genes) → Step 2: Signal Pep (254 Genes) → Step 3: PF03583 (4 Genes) → Step 4: Orthologs (10 Genes) → Add a step

10 Genes (1 ortholog groups) [Revise this search](#)

Gene Results | Genome View | Analyze Results

Rows per page: 1000

Gene ID	Transcript ID	Organism	Product Description	Input Ortholog(s)	Ortholog Group	Paralog count
C1_09420W_A	C1_09420W_A-T	<i>Candida albicans</i> SC5314	Triacylglycerol lipase [Source:UniProtKB/TrEMBL;Acc:Q5APG1]	MRET_0930,MRET_1179,MRET_4098,MRET_4099	066_130000	9
C1_09580C_A	C1_09580C_A-T	<i>Candida albicans</i> SC5314	Triacylglycerol lipase [Source:UniProtKB/TrEMBL;Acc:AOA1D8PER6]	MRET_0930,MRET_1179,MRET_4098,MRET_4099	066_130000	9
C1_09590C_A	C1_09590C_A-T	<i>Candida albicans</i> SC5314	Lipase 10 [Source:UniProtKB/Swiss-Prot;Acc:Q9P4E5]	MRET_0930,MRET_1179,MRET_4098,MRET_4099	066_130000	9
C1_09600C_A	C1_09600C_A-T	<i>Candida albicans</i> SC5314	Lipase 6 [Source:UniProtKB/Swiss-Prot;Acc:Q9P4E8]	MRET_0930,MRET_1179,MRET_4098,MRET_4099	066_130000	9
C1_09900W_A	C1_09900W_A-T	<i>Candida albicans</i> SC5314	Triacylglycerol lipase [Source:UniProtKB/TrEMBL;Acc:AOA1D8PEQ3]	MRET_0930,MRET_1179,MRET_4098,MRET_4099	066_130000	9
C6_04490W_A	C6_04490W_A-T	<i>Candida albicans</i> SC5314	Lipase 4 [Source:UniProtKB/Swiss-Prot;Acc:Q9PBW1]	MRET_0930,MRET_1179,MRET_4098,MRET_4099	066_130000	9
C7_02830C_A	C7_02830C_A-T	<i>Candida albicans</i> SC5314	Lipase 5 [Source:UniProtKB/Swiss-Prot;Acc:Q9PBW0]	MRET_0930,MRET_1179,MRET_4098,MRET_4099	066_130000	9
C7_02880C_A	C7_02880C_A-T	<i>Candida albicans</i> SC5314	Lipase 9 [Source:UniProtKB/Swiss-Prot;Acc:Q9P4E6]	MRET_0930,MRET_1179,MRET_4098,MRET_4099	066_130000	9

How can you lower the stringency of the search by removing the third step from the search without deleting it? (Hint: you will need to use a Boolean operator).

ACE1 *

Step 1: lipase (44 Genes) → Step 2: Signal Pep (254 Genes) → Step 3: PF03583 (4 Genes) → Step 4: Orthologs (10 Genes) → Add a step

Details for step Combine Gene results

4 Genes

Revise as a boolean operation

1 2 INTERSECT 3 2 UNION 3 2 MINUS 3 3 MINUS 2

Revise as a span operation

2 RELATIVE TO 3, using genomic colocation

Ignore one of the inputs

2 IGNORE 3 3 IGNORE 2

3 Revise

- Save the strategy and share the link in the Zoom chat.



Strategy URL:

<https://fungidb.org/fungidb/app/workspace/strategies/import/de6db119526563ad>

References:

Park et al. J. Microbiol. Biotechnol. 2021; 31(5): 637-644 doi:10.4014/jmb.2012.12048