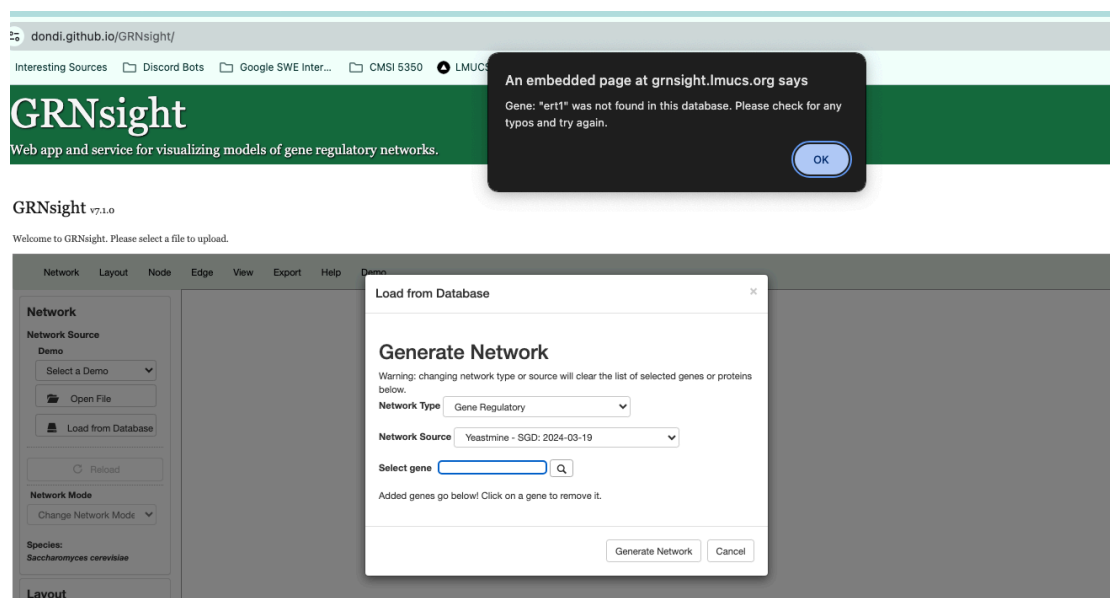


# GRNsight Database

## Background

GRNsight includes four primary databases: the expression database, grnsettings database, network database, and protein-protein database. This document will focus on the network database and the protein-protein database. Both of these databases contain a gene table, a source table, and a network table (referred to as physical interactions in the protein-protein database). Additionally, the protein-protein database features a protein table. All scripts for downloading and uploading data to these databases are located in the database folder.

## Problems



ERT1 is not in the network database. After investigating the queries, I discovered that our original approach involved first querying the list of regulators, followed by querying the targets for each regulator. However, ERT1 was not identified as a regulator in this process. See [here](#) for the list of regulators.

In contrast, when I queried the list of regulations from AllianceMine, I found approximately 15 connections where ERT1 is recognized as a regulator. You can view the relevant data in [Box](#).

**Comparison of Databases:** I conducted a comparison between the server network database and the AllianceMine network:

	Server Database	AllianceMine
--	-----------------	--------------

Missing network	2961	17
Missing regulator genes	270	1
Missing target genes	44	6

For further details on the missing networks or genes, please refer to the provided [link](#).

**Migration Note:** Also, [YeastMine is no longer supported](#), and we need to transition to AllianceMine. The service link in the sample code is currently non-functional.

## Proposed Solutions

### ERT1 issue

Solution 1: Get the network by querying the list of regulations in one request.

As we can see from the AllianceMine network, we can see more regulatory networks from our current query. However, we don't know whether this is because of the query or because new data has been updated since Yeastmine is down.

Pros	Cons
<ul style="list-style-type: none"> <li>✓ Query all regulatory networks in one query.</li> <li>✓ Capture more networks than querying all the regulators and find the targets for each query.</li> <li>✓ Don't need to perform multiple manual queries if you can't request automatically by the service.</li> </ul>	<ul style="list-style-type: none"> <li>✗ Possibly missing some regulators</li> </ul>

Solution 2: Combine the network database from the list of regulations query and the list of regulators and their corresponding target queries

To address the issue in solution 1, we can first query the list of regulations query. Then we query the list of regulators. If we don't see any gene in the regulators list in regulations, then we would query the target for each missing regulator.

Pros	Cons
<ul style="list-style-type: none"> <li>✓ Capture all the regulatory networks</li> </ul>	<ul style="list-style-type: none"> <li>✗ Require a lot of queries if there are a lot of missing regulators (even worse if we can't automate the task since AllianceMine server is not working).</li> </ul>

## AllianceMine

Solution 1: Manually download data from AllianceMine, then process the file

Download the required data file using the AllianceMine query builder. Write a Python script to process the downloaded file, extracting relevant details such as genes and regulators (would be the same for protein-protein-interactions). Depending on the approach taken to resolve the ERT1 issue, we would either develop and process the data from the sheet or reuse the existing code.

Pros	Cons
<ul style="list-style-type: none"><li>✓ Simple, don't need to worry about web scraping</li><li>✓ Know the exact structure of the downloaded file</li><li>✓ Changes to the AllianceMine web interface or the service API won't affect the process (as long as the file download format remains the same)</li></ul>	<ul style="list-style-type: none"><li>✗ Required a manual download process, which might have errors and time-consuming.</li></ul>

Solution 2: Automating Data Retrieval with Web Scraping

I'm not familiar with this process, and I require more investigation about this approach as how we can download the whole table.

Pros	Cons
<ul style="list-style-type: none"><li>✓ Prevent errors from manual work.</li></ul>	<ul style="list-style-type: none"><li>✗ Can break when there are changes to the website structure</li><li>✗ Violate the terms of service of the website (need to double check AllianceMine policies)</li><li>✗ Requires frequent monitoring and maintenance to ensure the scraper works as expected.</li></ul>

## Appendix

### 1. GRN Network AllianceMine Query Builder

The screenshot shows the 'Query Editor' tab of the GRN Network AllianceMine Query Builder. The interface displays a tree view of available fields for a query. The fields are organized into several categories:

- Gene**
  - Systematic Name
  - Standard Name
- Organism**
  - Species
  - Taxon Id
- Regulatory Regions (RegulatoryRegion)**
  - TF Binding Site** (highlighted)
  - Annotation Type
  - Datasource
  - Experiment Condition
  - Regulation Direction
  - Strain Background
- Publications (Publication)**
  - PubMed ID
- Reg Evidence (RegulationEvidence)**
  - Ontology Term**
    - Identifier
    - Name
- Regulator (Gene)**
  - Systematic Name
  - Standard Name

### 2. Regulators Query

The screenshot shows the 'Query Editor' tab of the GRN Network AllianceMine Query Builder. The interface displays a tree view of available fields for a query. The fields are organized into several categories:

- Gene**
  - Name
  - Systematic Name
  - Sgd Alias
  - Standard Name
- Regulation Summary**
  - Summary Paragraph
- Publications (Publication)**
  - Citation
  - PubMed ID