# CRISPR-GEM v1 Tutorial

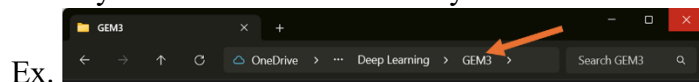**Table of Contents:**

## Overview:

CRISPR-GEM is a machine learning model built to identify promising CRISPR target genes for  specific gene editing applications. There are two individual components of CRISPR-GEM, a graphical user interface (GUI) and a deep neural network (DNN) model in Google Colab. The first is a GUI that is used to compare different cell types and narrow the list of genes to more precisely model CRISPR strategies. The GUI provides a list of input genes and a list of output genes that are candidate genes for the editing strategy or important markers to assess a successful therapy respectively. To accomplish this, users must specify the tissue/cell types to be edited (ex. Diseased Cartilage), the goal for gene editing (ex. Healthy Cartilage), and the type of CRISPR gene editing to be used (ex. CRISPRa) using the GUI.

Once the appropriate cell/tissue types are selected and the GUI is used to obtain the input/output genes for the model, this information can be transferred over to a google drive account for the assessment of each gene candidate. The DNN learns the expression of the output genes provided by the GUI using the expression of the input genes then perturbs input gene expression in a CRISPR-mimetic fashion to assess the downstream effect on output gene expression. The result is a ranked list of CRISPR candidates that score the best for accomplishing a therapeutic goal using CRISPR. These genes should be assessed *in vitro* to verify their accuracy.

It is important to note that this model is highly sensitive to the datasets used to identify gene candidates in the GUI and those fed to the DNN model. When comparing sequencing datasets, it is advised to always use data from the same experiment when possible and at the very least from the same sequencing platform. The DNN is trained only on HiSeq2000 and HiSeq2500 data to reduce any noise resulting from different platforms. If users have their own datasets for these platforms they can be easily integrated into the existing datasets and used for analysis (pg. 10). It is also possible for users to select new datasets for analysis through the Gene Expression Omnibus (GEO) to add to existing datasets (pg. 2). Additionally, users may select target genes manually via differential expression analysis or alternative methods if they have data that is not yet public. Users must be reminded that when obtaining new datasets or using their own, meaningful integration into the model requires that they are sequenced using HiSeq2000 and HiSeq2500. However, adaptations are underway to accommodate alternative datasets.

**Downloading GUI:**

- Download python 3.11.5
  - Python 3.11.5 is required since some modules have been depreciated in newer versions.
  - Make sure that you add it to your systems PATH
  - If a later version of python is downloaded it may be necessary to use a virtual environment to run python 3.11.1.
- The datasets for the GUI & model training are different. They will both have to be personally shared over google drive so email Josh at jog822@lehigh.edu for access.
- The remaining components for the CRISPR-GEM can be downloaded from the CRISPR-GEM Github.
- Download the Github files and dataset to the same location then open a command prompt. (On windows type cmd into the search bar and immediately hit enter).
- Identify the file location of where you downloaded the files.

Ex. 

Right click on the file location and click "Copy Address as Text".
- In the command window type "cd " then paste the address and hit enter.
- Now you can run the script by typing "python GUI.py"

**Using GUI:**

1. Copy the folder location that GUI files and datasets are in
2. In the command window type "cd " then paste the address and hit enter.
3. Now you can run the script by typing "python GUI.py"
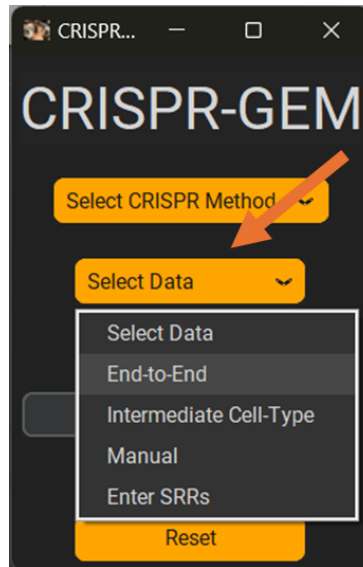4. The GUI will take 1-2 minutes to load and should look like this:



5. Start by selecting the CRISPR strategy based on if you want to perform gene knock-in, knock-out, activation, or repression.
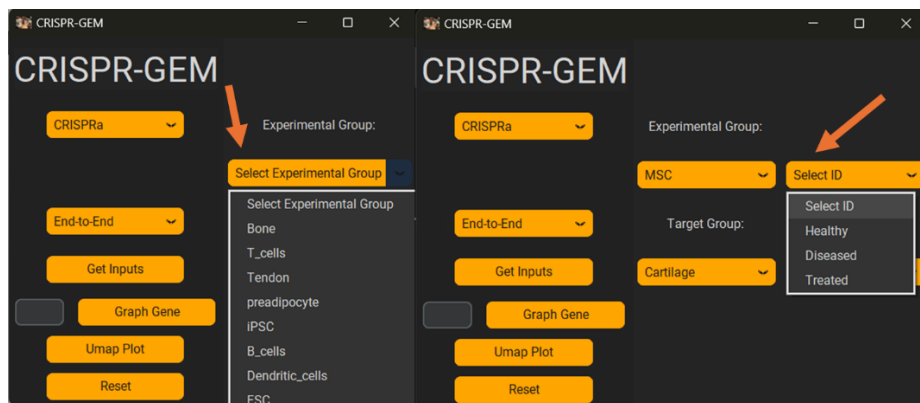


6. Then select the format with which you want to set up the model (Select Data). This has to do with how many cell/tissue types that you want to enter. End-to-End (shown here) uses just two cell tissue types. However, the user can also choose to add a third type using the "Intermediate Cell Type" option. This gives the model a trajectory to work off of. For example, when analyzing chondrogenesis, the Intermediate cell type for MSC chondrogenesis was fetal cartilage which represents the natural
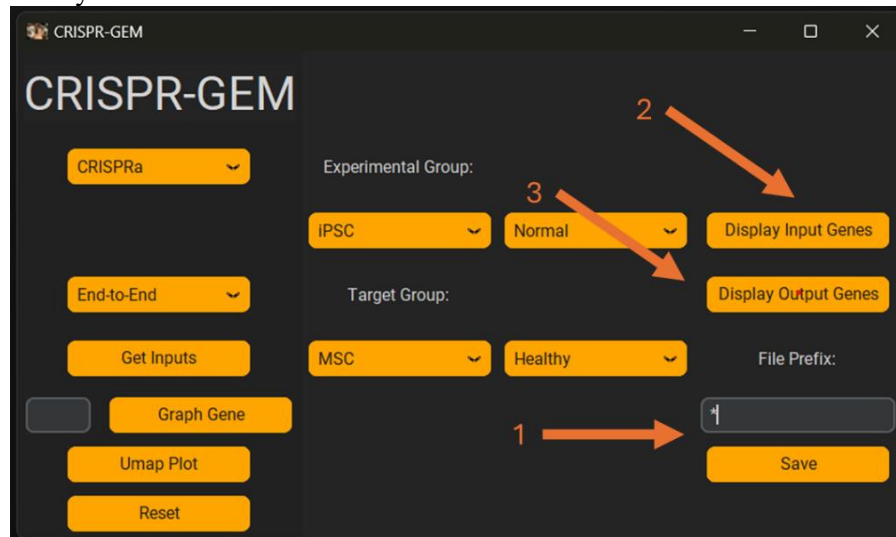
development pathway. Alternatively, the manual method can be used to manually enter genes by giving their Ensemble IDs.



7. The actual cell/tissue types are then selected using two categories. The first is a broad description of the group (ex. MSC, chondrocyte, fibroblast) is labeled 'Select Experimental Group' and pops up first. Then a more specific description such as healthy, diseased, treated, etc. can be selected in the 'Select ID' dropdown that will pop up to the right of the 'Select Experimental Group' tab.



8. Once everything is selected, one can then select for the input/output genes based on those groups using the 'Get Inputs' button.

9. This will take a couple minutes to run but once it is finished a new window will pop up in the corner with an entry box and 'Save' button (1, below). The user should type a short tag that describes the experiment here and save the file. This will create a new folder in the directory. For example, let * represent the short tag, then the output will be files will be in a folder named '*' in the GUI directory with files named '*_input_list.csv', '*_output_list.csv', '*_SRRs_used.csv', '*_input_genes.csv', and '*_output_genes'.
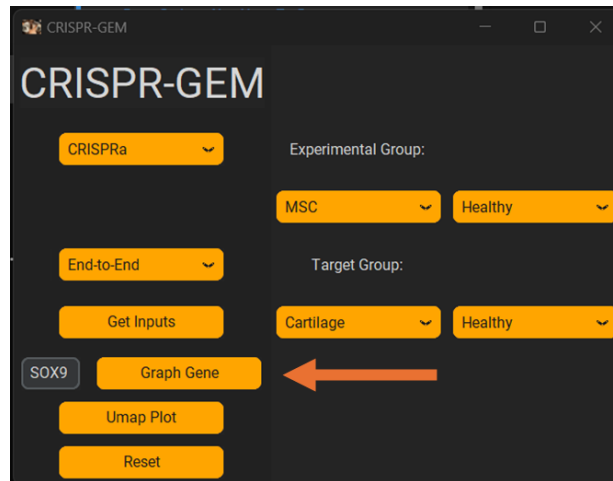
10. The results can be monitored by clicking the 'Display Input Genes' or 'Display Output Genes' buttons (2 & 3 above respectively). These give interactive datasets that you can use to view the results. This can be done by selecting a column (ex. Log2FoldChange) and sorting it by ascending or descending value (1, below). Or a user can input a specific gene (make sure it is all caps) in the entry box (2, below) to find well established genes.for comparison.
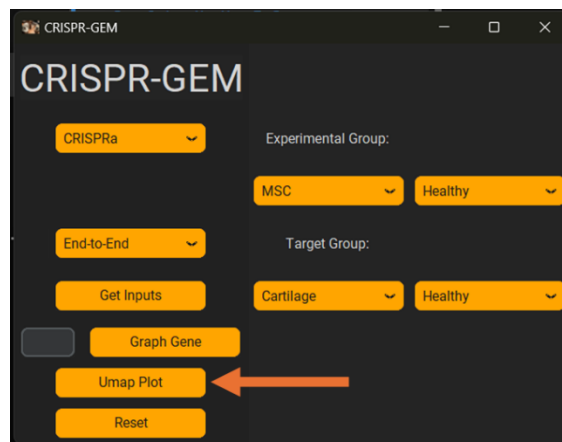


11. To evaluate datasets before using the identified input/output genes, one can graph the relative expression of a specific gene by typing it into the entry box then clicking the "Graph Gene" button next to it (ex. *SOX9* below). The plots the relative expression in each sample.
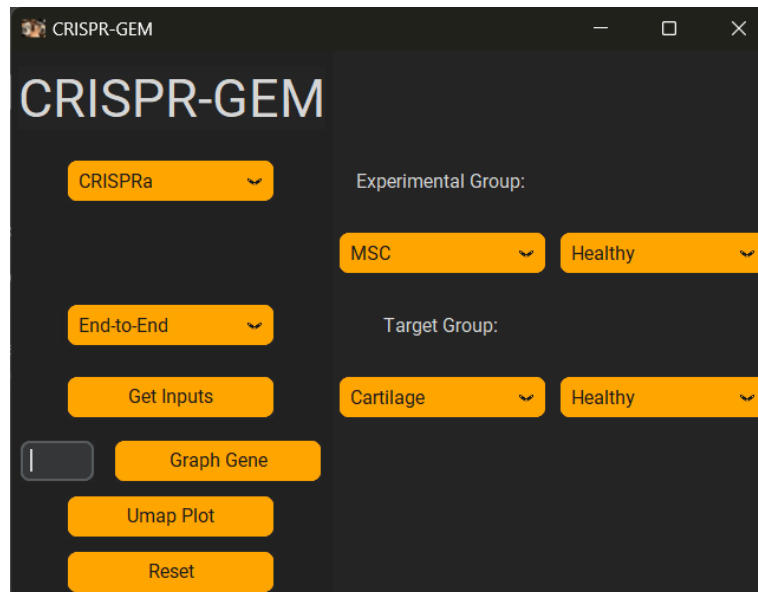
12. Additionally, more global representations of the datasets for each selected cell/tissue type can be represented by clicking the 'Umap Plot' button.
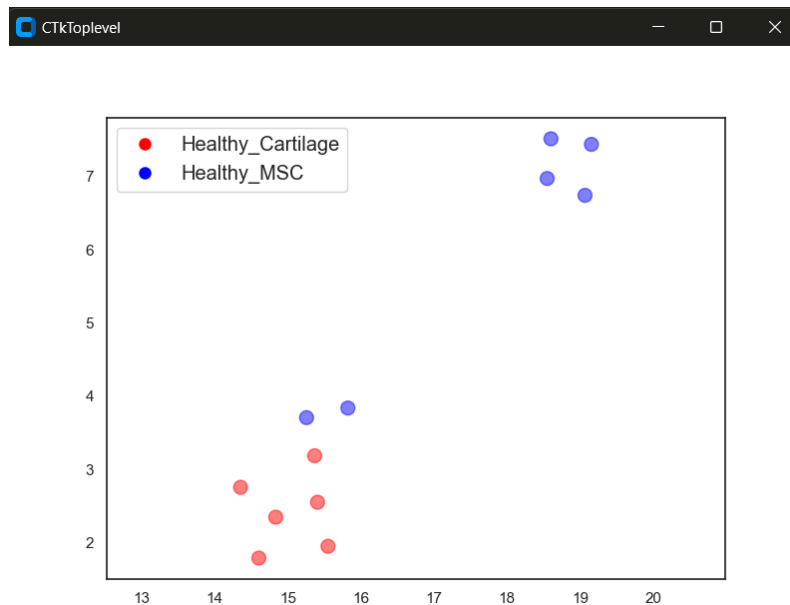
## **Troubleshooting Misaligned Data:**

Ex. Healthy cartilage and healthy MSCs using the end-to-end method



But before pressing the "Get Inputs" button, it is  advised to evaluate the underlying data to ensure that the input and output genes are appropriately assessed. To do so, one can start by using the UMAP function to visualize the data in a reduced dimension to get an idea of how similar the data from each group is. In doing so, the following graph is plotted:
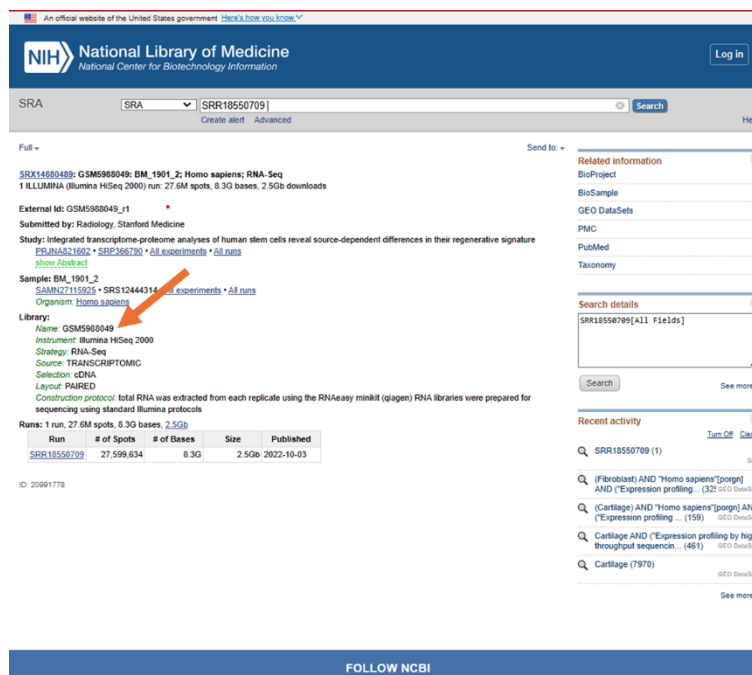


There are two clear clusters for each group, however two of the six MSC datapoints seem to be more similar to the Cartilage group which may interfere with differential expression analysis.  To  overcome this, one can look at the datasets used for analysis using the underlying command window which should look something like this:

.

The datasets used can be identified in the given lists and searched in the SRA for more insight. It is often helpful to find the GEO ID then search the GEO for the dataset as it provides more information.



Once the misaligned dataset is found, the inappropriate datasets can be found in the 'Labels.csv' file and removed by changing the tag by adding a number to the end or by adding a better description. For example, a likely change source of error in the given circumstance is that MSCs were senescent, differentiated, or harvested from an older source so the description 'Healthy' could be changed to 'Senescent' or just 'Healthy1' to remove it from the analysis. Don't forget to change that label back to its original description if it was not the cause of the issue.
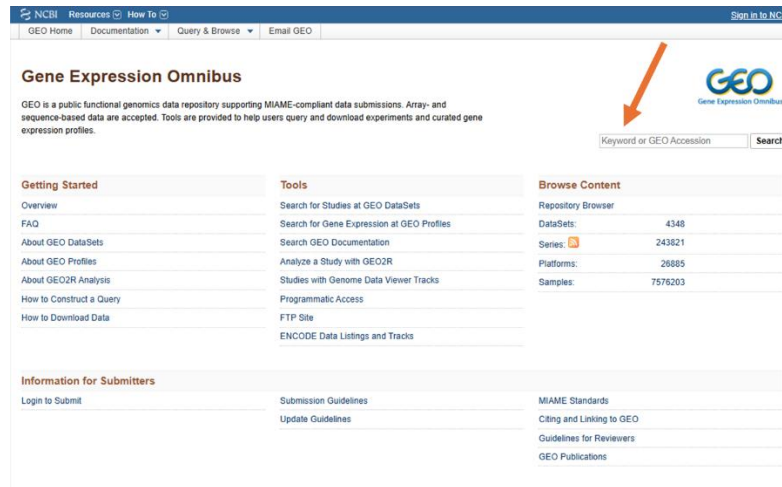
Sometimes it is hard to notice any significant differences just from the descriptions on the SRA or GEO, hence why they were all labeled the same. An alternative option is to eliminate specific datasets one-by-one to monitor which is misaligned. This can be done in two ways. **1)** Use the 'Enter SRRs' option in the "Select Datasets" dropdown menu. This allows you to paste a comma separated list (ex. SRR123, SRR124, …) so it is easy to eliminate a dataset or group of similar datasets one-by-one to identify the inappropriate dataset. **2)** Changing items in the "Labels.csv" file to identify the issue. It is advised to change all datasets containing the same 'GEO' accession (ex. GSE1234…) at once since these are most likely collected at the same time and from the same source. In both cases, once the changes are made, the data can be plotted using the UMAP button again to view the effects. Don't forget to change that label back to its original description if it was not the cause of the issue.

Please share any inappropriately labeled datasets with [jog822@lehigh.edu](mailto:jog822@lehigh.edu). In addition to using UMAP. Users can assess datasets by graphing multiple gene, although this strategy is not as encompassing as UMAP.
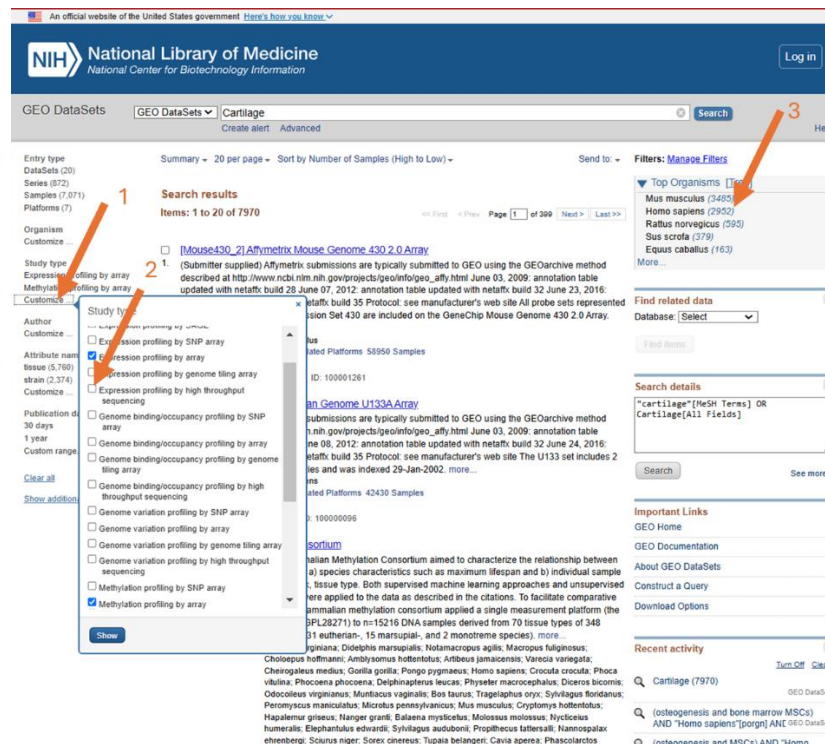
**Adding new datasets into the CRISPR-GEM Database:**

Data may need to be selected in the case that the appropriate datasets are not available on the GUI or if there is not a deep enough data n for that specific cell/tissue type on the GUI.

- All data can be selected from the gene expression omnibus: https://www.ncbi.nlm.nih.gov/geo/.
- Search for appropriate datasets using the search bar.



- In the search results make sure that you filter results for the appropriate data types. Click 'Customize' under Study Type (1), then make sure 'Expression profiling through High Throughput Sequencing' is checked (2). This option will then appear under Study Type. Click 'Expression profiling through High Throughput Sequencing' again to filter results. Then, filter for only human datasets by clicking the number next to 'Homo Sapiens' (3).

- All Data sets used for model training should be taken from HiSeq2000 or HiSeq2500; however, for simply determining the input/output genes for model training other platforms can be used as long as differential expression is always performed across datasets from the same platform. The platform is generally listed on the main page for a dataset (1, below). However, in some cases multiple platforms are listed, and one must click on individual datapoints to identify the platform.



- To download the data click on SRA selector (2, above). The data must be downloadable from SRA for each addition to datasets. All that is needed is the SRA number for each dataset which can be linked to the GSM accession number in the dataset page (3, above).

- These SRA numbers may be either added to an excel file/CSV formatted as follows. Or copied into a comma separated list eg. 'SRA1', 'SRA2', ….

| SRR ID |
|--------|
| SRR1 |
| SRR2 |

** Any other information can follow such as the GSM #, cell type, disease state, etc.

- Run the code "Kallisto_sra_downloads_Official.ipynb" with the corresponding SRR codes (See line 34 or 36/37). Make sure to follow the prompts.
- To add this data to the dataset for the GUI, add 'Big.csv' to Google Drive in same place as TPM.csv. Then toggle (add and remove '#') back and forth between the pairs of lines below, with the bottom always representing the "Big.csv" and the top line representing "TPM.csv".

```python
hello = pd.read_csv('/content/drive/MyDrive/TIGR_/TPM.csv', index_col=0)
# hello = pd.read_csv('/content/drive/MyDrive/TIGR_/Big.csv', index_col=0)
```

```python
est_counts = tsv_df['tpm'].rename(folder)
# est_counts = tsv_df['est_counts'].rename(folder)
```

```python
TPM.to_csv('/content/drive/MyDrive/TIGR_/TPM.csv')
#TPM.to_csv('/content/drive/MyDrive/TIGR_/Big.csv')
```
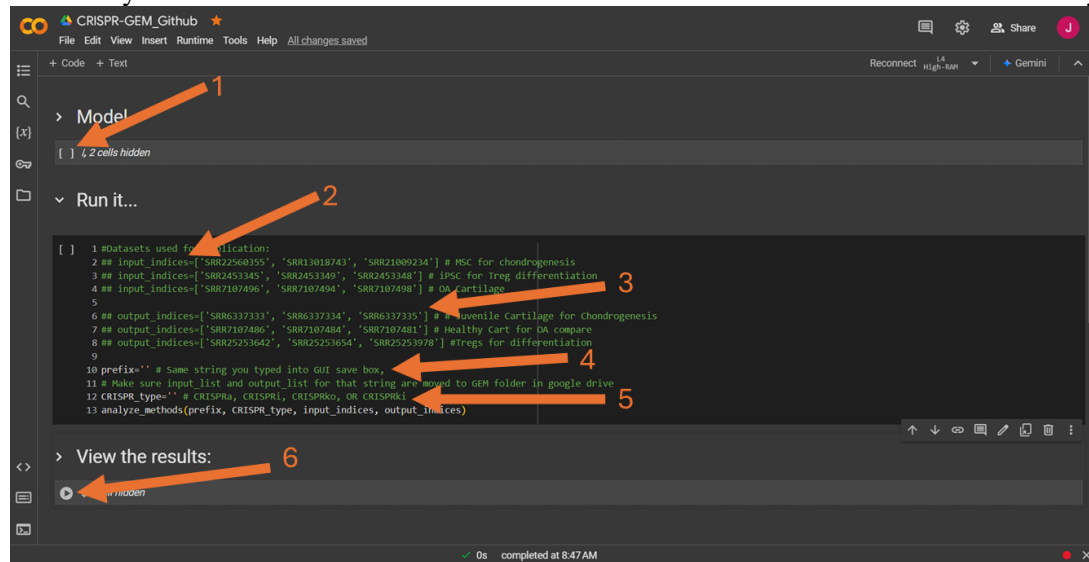
## **Using CRISPR-GEM in Google Colab.**

1. Create a New Folder in Google Drive (Your 'My Drive') Called "GEM".
2. Add the *_input_list and *_output_list into this folder from the GUI.

   * It is also possible to create your own list of genes for each group without the use of the GEO. For example, if one has their own datasets, they can perform differential expression analysis and filter the results using the rules described in table 4 of the manuscript. In this case, ensure that the csv files listing the genes are formatted properly (including capitalization and the index (0,1,2…) in the first column):

|   | inputs |   |
|---|--------|---|
| 0 | A1BG |   |
| 1 | A1CF |   |
| 2 | A2M |   |

3. Open CRISPR-GEM_Github.ipynb in google colab (Available on Github).
4. Run the 'Model' section. Hover your mouse in between the brackets (1) then a play button will appear. Click the play button and let this run. You may need to approve it to access your google drive.

5.  Next the model needs to be set up based on the specific application. Use the SRR IDs output by the GUI and chose three random IDs for both the inputs and outputs. Add them in list format for the respective labels (2 & 3).

6.  The GUI outputs the input and output genes as *_input_list and *_output_list where * is the designated prefix. The exact same prefix should be pasted between the single quotations after 'prefix' (4).

7.  Finally, add the CRISPR type to be assessed (5). This should be the same as the one used in the GUI.

8.  Run the "Run it…" cell and wait for the model to complete running. The time can range from 2-45 minutes depending on the number of input/output genes and the datasets used.

9.  Once the model is finished running, 'View the results' can be run to give a ranked list of the top scoring genes (6).

*Don't forget the datasets you use are very important. Make sure you analyze the datasets using the GUI to make sure you are making meaningful comparisons. If the gene results don't make sense, new datasets may have to be examined.