# Final Project Proposal: GO Simplify
## Beverlie Poblete

**Tool Background**

With the increasing number of advancements in biotechnology, we are able to collect high volumes of genomic data at decreasing costs. This has shifted the need for collecting biological data to a need of analyzing and making sense of the biological data. One way to tackle this problem was through ontology, which has enabled the organization and analysis of large datasets. Ontologies are most particularly useful in creating genomic annotations. The Gene Ontology (GO) database was developed to record annotations of the functional properties of gene products across species, which further contributes to gene function prediction. GO annotations are created by associating a gene or gene product with a GO term. These GO annotations thus capture a "snapshot" about how a gene's molecular function, location in the cell, and what biological processes it is involved in. In the GO database, each ontology consists of a set of GO terms, which are organized in a directed acyclic graph (DAG) where each GO term is a vertex of the graph, and the edges encode the relationship between GO terms. GO annotations stores the currently known functional knowledge of gene products where a positive annotation indicates the gene product carries out the function described by the GO term. A negative annotation in turn indicates the gene product does not perform the function described by the term. The Gene Ontology (GO) knowledgebase (http://geneontology.org/) is the world's largest source of information on the functions of genes. While exploring the knowledgebase, I found I had to click through multiple links to find related genes and gene products annotated to the same GO class of a certain gene product or related annotations. With my final project tool, GO Simplify, I aim to create a user friendly application that will take in a simple inputs of a gene/gene product and organism. The GO simplify tool will then display the associations between GO terms and the gene/gene product for the organism input. There will also be a follow-up option for users to display other genes/gene products annotated to the same GO term.

**Tool Functionality**

This tool requires two basic inputs: a gene or gene product and organism. With the inputs, the tool will query a subset database (that will be a subset of what is offered in the GO knowledgebase) and return and display the associations between your input gene/gene product and GO terms for the select organism. Following this display, a dropdown will be displayed containing the resulting GO terms from the previous search. The end user can select one of the options from the dropdown to see what other genes/gene products were annotated to the selected GO term. With this tool, I aim to simplify the data presented in the GO knowledgebase.

**Tool Description**

The following software technologies will be used in implementing GO Simplify:

1. **SQL relational database/table**
   This database will hold information of gene/product symbol, gene/product name, GO class (direct), and organism. The database will also hold information on genes and gene products annotated to positive regulation of select GO classes. I will manually curate the database to

hold a subset of data for ~50 gene/gene products.  I will be collecting the data to populate my database from the GO knowledgebase.

2. **Python-based Computer Gateway Interface (CGI)**
   The CGI script will ingest the user's input of gene/gene product and organism and query the SQL table for the relevant gene/product symbol, gene/product name, and GO class (direct) for the input gene/gene product and organism.  The CGI will also take in the follow-up user input from the dropdown of relevant GO classes to query the SQL database for related genes/gene products with annotations to the selected GO class.  All information from both queries will be transmitted back to the user via the GUI.

3. **HTML/JavaScript-based graphic user interface (GUI)**
   HTML, JavaScript, and possibly CSS will be used to create a user-friendly web interface for taking in the user's input via HTML form and later drop-down and displaying the results of the query performed in the CGI script.
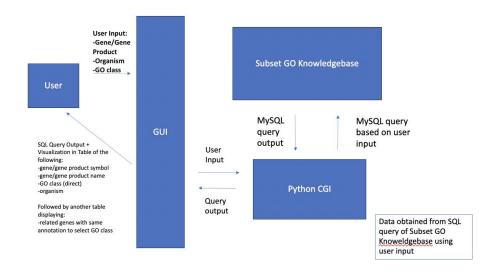


Figure 1. Software components and relative interaction and information exchange

**Tool Design and Development**

The following describe the proposed tool design and development:

1. I am planning on implementing a schema and relational database containing two tables.  The first table will hold gene/gene product symbol, gene/product name, organism, GO class (name), and GO class key.  The second table will contain a row per gene/product containing columns holding data for gene/product name, GO class name, and GO class key (foreign key for first table).  This proposed database will be manually curated through queries in the GO database for a select number of gene/gene products for ~2 or 3 organisms to be able to complete this final project by the deadline.  I will gather the information in an excel sheet prior to creating the database on the bfx3 server.

Small examples of what information the proposed tables will hold are shown in Table 1 and Table 2 below.  Please note I limited the rows to 2-3 rows in the examples below to be concise.

| Gene/gene product symbol | Gene/product name | organism | GO class (direct) | GO class key |
|---|---|---|---|---|
| Tcf7l2 | Transcription factor 7 like 2 | Rattus norvegicus | Positive regulation of insulin secretion | 1 |
| Cckbr | Cholecystokinin A receptor | Rattus norvegicus | Insulin secretion | 2 |
| Cat | catalase | Rattus norvegicus | Response to insulin | 3 |

Table 1. *Proposed Table 1*

| Gene/product | Gene/product name | GO class (direct) | Go class key |
|---|---|---|---|
| Tcf7l2 | Transcription factor 7 like 2 | Positive regulation of insulin secretion | 1 |
| Casr | Calcium-sensing receptor | Positive regulation of insulin secretion | 1 |

Table 2. *Proposed Table 2*

2. One CGI script will ingest the user's input of gene/gene product and organism and query the SQL table for the relevant gene/product symbol, gene/product name, and GO class (direct) for the input gene/gene product and organism.  Another CGI script will also take in the follow-up user input from the dropdown of relevant GO classes to query the SQL database for related genes/gene products with annotations to the selected GO class.  All information from both queries will be transmitted back to the user via the GUI.

3. A GUI created using HTML/JavaScript/CSS will be prepared to take user input via HTML form. The first input will be for the gene/gene product string and then organism (most likely radio button to choose amongst 2-3 organisms).  Using this user input, the GUI will display the SQL query output returned by the CGI script as a table on the web page.  After this initial table is displayed, a dropdown will be displayed with a dropdown of the GO classes returned from the first query.  The user can select an option from the dropdown.  The GUI will take this input and possibly call another python CGI script that will perform another query to find related gene/gene products with annotations to the same GO class.  The output from this second query will be displayed to the user in the GUI as a second table of results.
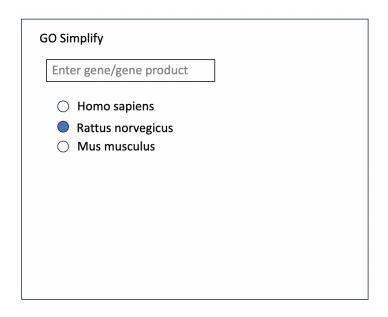
Figure 2. First user input screen

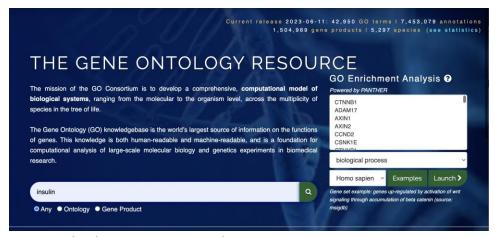

Figure 3. Output/Display after first user input submitted
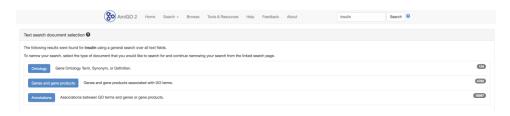
Figure 4. Output/Display after second user input (dropdown)

**Anticipated Challenges/Obstacles:**

I anticipate the data mining section of the final project will be the most difficult and time-consuming. While the GO database contains a large number of entries, I plan on sub-setting my database to 50 gene/gene products for 2-3 organisms. I need to select a limit of how many rows to display in my results to not overload the page with too much data. I also may need to limit the number of genes per GO class key in the proposed second table to keep my database small.
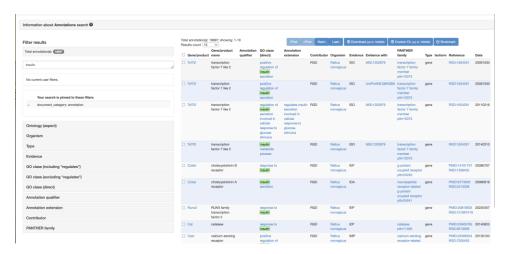
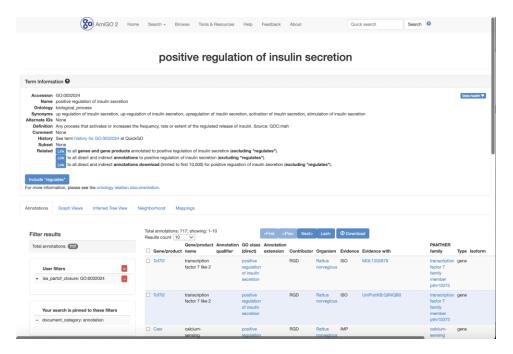## Appendix – How I plan on mining data from GO database:
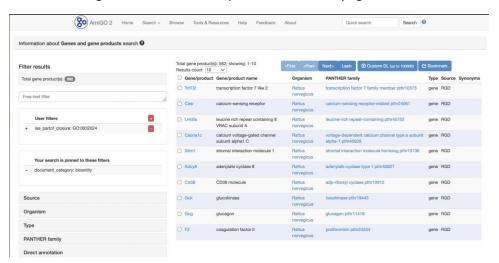


*1. On GO database, enter gene product*



2. Select "Annotations"



3. Use results here to populate table 1 in database

4. Clicking "GO class" link in step 3 will take me to page for selected GO class



5. Clicking on "Link" to all genes and gene products annotated to.... under "Related" in step 4 will take me to table showing genes and gene products with annotation to selected GO class. Use this data to populate proposed table 2 of database.

**Sources:**

Carbon, S., Ireland, A., Mungall, C. J., Shu, S., Marshall, B., & Lewis, S. (2008). Amigo: Online access to ontology and annotation data. *Bioinformatics*, *25*(2), 288–289. https://doi.org/10.1093/bioinformatics/btn615

Hill, D. P., Smith, B., McAndrews-Hill, M. S., & Blake, J. A. (2008). Gene ontology annotations: What they mean and where they come from. *BMC Bioinformatics*, *9*(S5). https://doi.org/10.1186/1471-2105-9-s5-s2

Zhao, Y., Wang, J., Chen, J., Zhang, X., Guo, M., & Yu, G. (2020). A literature review of gene function prediction by modeling gene ontology. *Frontiers in Genetics*, *11*. https://doi.org/10.3389/fgene.2020.00400