# CS3860 Lab 1 - Database

**MySQL install**

Objective: Understand database management system installation, and gain a high level overview of database creation, defining table schemas, and the pragmatics of how to create tables, import data, and define primary keys and indexes.

**Record your observations and submit the output of your queries in your lab report for each major step (denoted with [major step] below).**

**Deliverables**
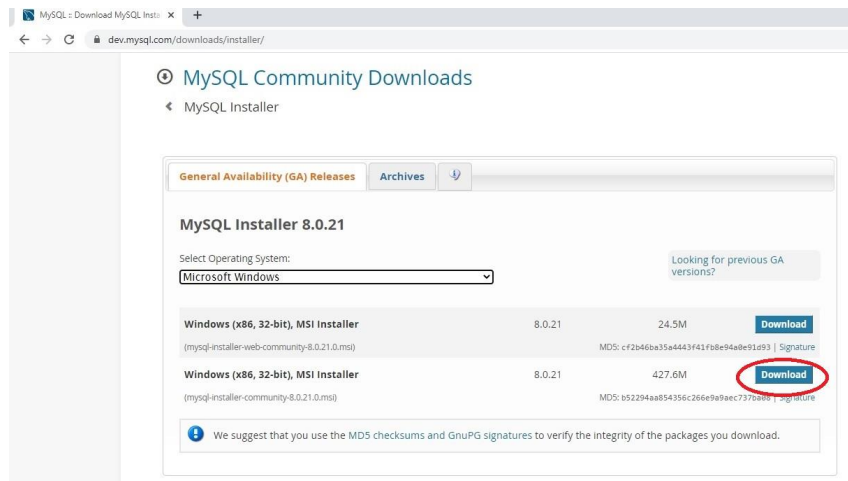Complete a lab report containing the following:
- What you've learned from the lab
- How it could be improved
- Observations and query results for each major step of the lab.
- ERD (Entity Relation Diagram)
- Submit to Canvas in PDF format!

*Notes:*
- ***When installing MySQL, it is important to create and remember your administrator login and password.***
- *Walk Through the Installation Process. Make sure you include options for command line, database administration, and automatic startup. Set and record your administrative password. Include command line tools.*
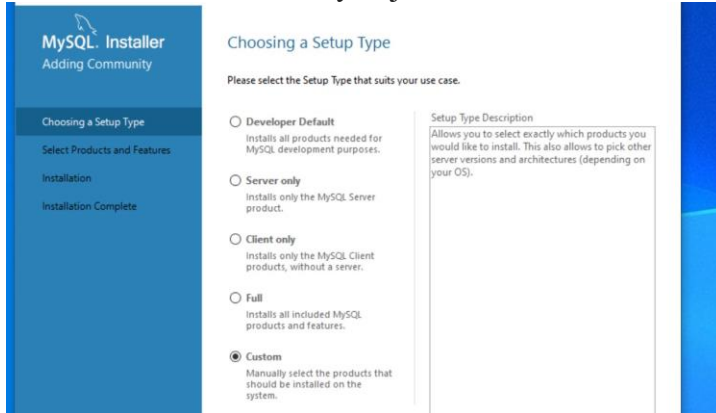
**Download and install MySQL Installer**
MySQL Community server(for all platforms): https://dev.mysql.com/downloads/installer/
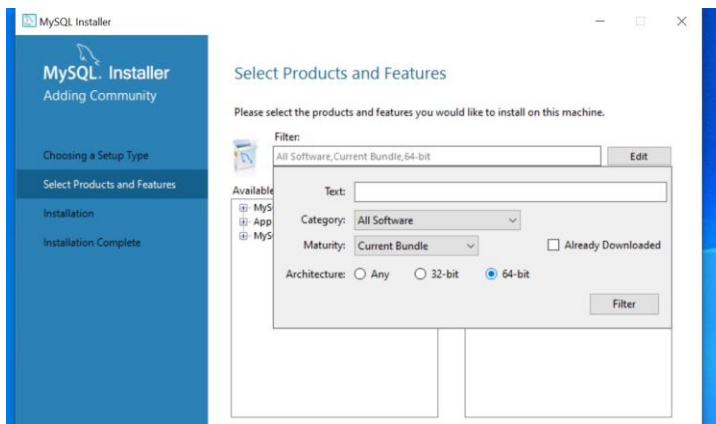
On the next screen, if asked to login/sign up you may simply click "No thanks, just start my download" link.
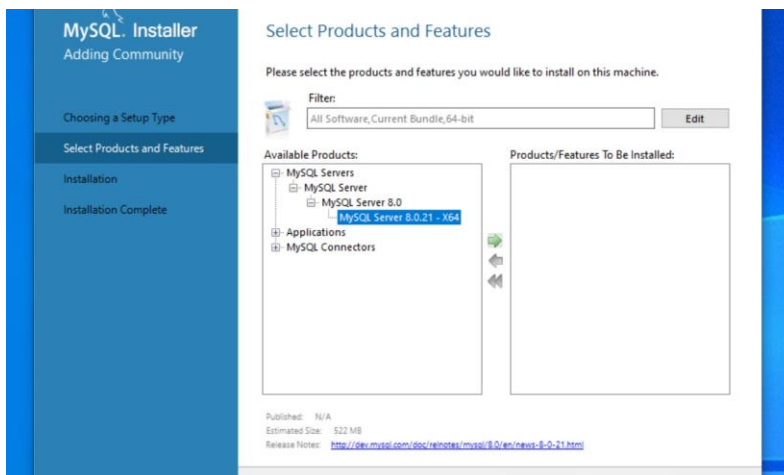
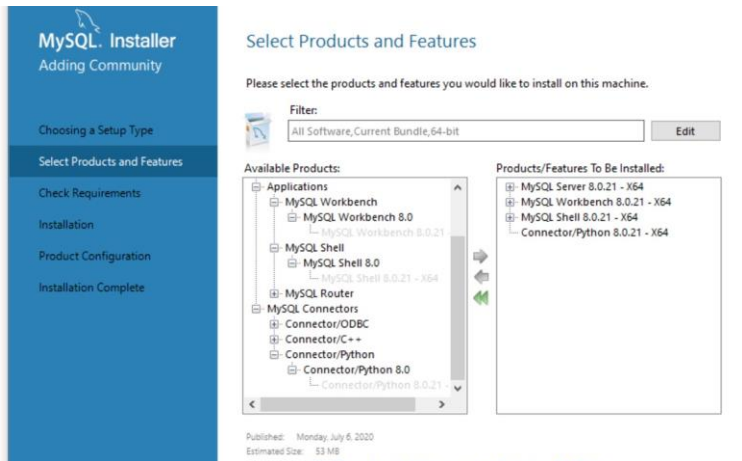Double click the installer you just downloaded, choose the "Custom" setup type.



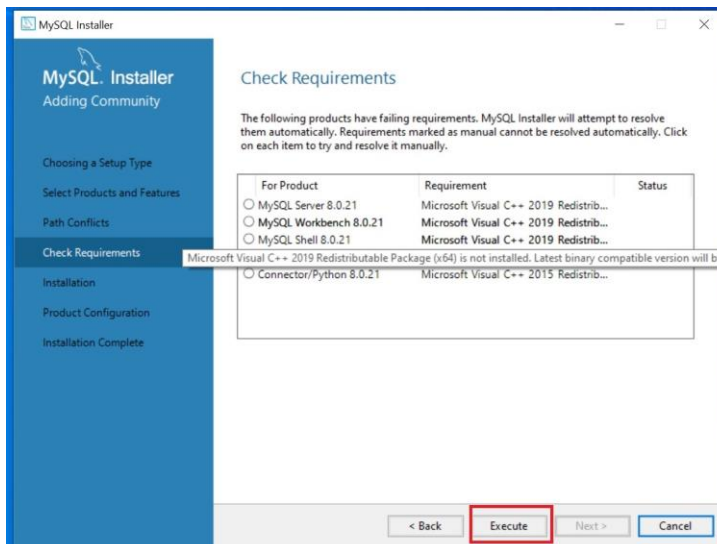You may filter to install products that are 64-bit.



Select the products to install by clicking on each one as shown in following screenshot, and then clicking the green arrow to move them over to the "to be installed" column
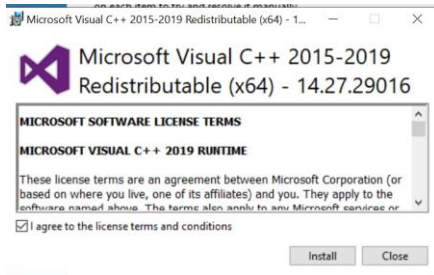
Please select Mysql Sever, Applications → Workbench, Mysql shell and python connector as shown below. **Make sure to match the 4 items shown in the right column before clicking next. Note that if you do not have python installed, it won't let you install the python connector. You can get that here (64bit windows): https://www.python.org/ftp/python/3.7.9/python-3.7.9-amd64.exe



In the next step, if you will find some requirements are not met for few mysql products. But mysql installer will try to download these requirements and install them for you. So, click execute on the below screen.



You will see that mysql tries to install Visual C++ (needed for Workbench), so check "I agree" and click install in the below screen.

Now, since all requirements are met, you are ready to install all Mysql products, so click execute on below screen.



Upon successful installation you will see the following screen.



Click Next, you will see the product configuration screen, click Next again.
On "High Availability" screen, select Standalone Mysql Server (you do not need *Innodb* cluster server or router)

In the "Type and Networking" screen (shown below), select "Development Computer", the default port is 3306, do not change anything on this screen, click Next



Next, choose "Use strong password Encryption for Authentication" option for Authentication Method screen.
Next, enter your password and **remember it.**
Now, set up mysql as windows service as shown below.



Finally, finish your installation by clicking execute on Apply Configuration screen.

**Open MySQL Workbench.**

You may perform the following steps to create a database and selecting a database schema by using either the *MySQL Workbench - SQL Development Tools* or the *command line.* If you use the *Query Browser* in the workbench to create a database, you can **do not need to create the database with the** *command line* (it only needs to be created once). However, you still need to look at the command line instruction section because we will be using the command line to load in a data file to populate our database.

**Setting of a database connection and creating a database schema**

Start *MySQL* Workbench. Under *SQL Development*, select **Connection +**.

Provide a logical name for your **connection**, enter your user name and password (store in vault). Leave *Default Schema* empty. Try *Test Connection*.



*Note: If you have trouble connecting, in the Control Panel, select Windows Firewall. Select the Exceptions tab, add MySQL or Port number: 3306 and name MySQL. If you still have trouble, check to make sure the MySQL service is running by "services.msc" from the windows command line prompt (Start->cmd).*

Click *OK*. Then double click on the connection to open it.

In the future, you will just be able to click on your connections on the main screen.

The SQL Editor will now be displayed showing any existing databases on the schema tab (left about halfway down). The "sys" schema was created by default during installation for me).

7

Select the can, or right-mouse-click on one of the existing databases in the *Object Browser* (left-hand panel) and select *Create New Schema* (the words schema and database are used interchangeably in this program). Name your new schema"*genomics*". Select *Apply*, (And Apply again) then select Finish for *Execute SQL Statements*.

Note: The image below shows the "can" icon:

A new database schema *genomics* should be available. Select the *genomics* schema (double click on it to **bold** it – this sets it to the current database "in use").



The following 2 pages are the command line approach to creating a database, You should review these (even though you just created the genomics database with the GUI using the steps above).

**Connecting to the MySQL Server for the First Time Using *Command line***

Two options:

Open the windows command line (Start->cmd) and navigate using the "cd" command to *C:\\Program Files\MySQL\MySQLServer 8.0\bin*

1. (harder) Use the mysql command to launch mysql

   ```
   %>mysql [options] [database]
   ```

   [database] is the name of your database. To list out the [options]:

   ```
   %>mysql --help
   ```

   For now, we just need to connect to the database server on your system:

   ```
   %>mysql -u root -p --local-infile
   ```

   *root* is the default login for MySQL. local-infile is an option to allow you to load local files. You will be prompted for the password.

   Or 2 (easier). Open the "MySQL 8.0 Command Line Client" directly from the Start menu. You will be prompted for the password.

The following should appear:
```
Welcome to the MySQL monitor.  Commands end with ; or \g.
Blah blah blah
```

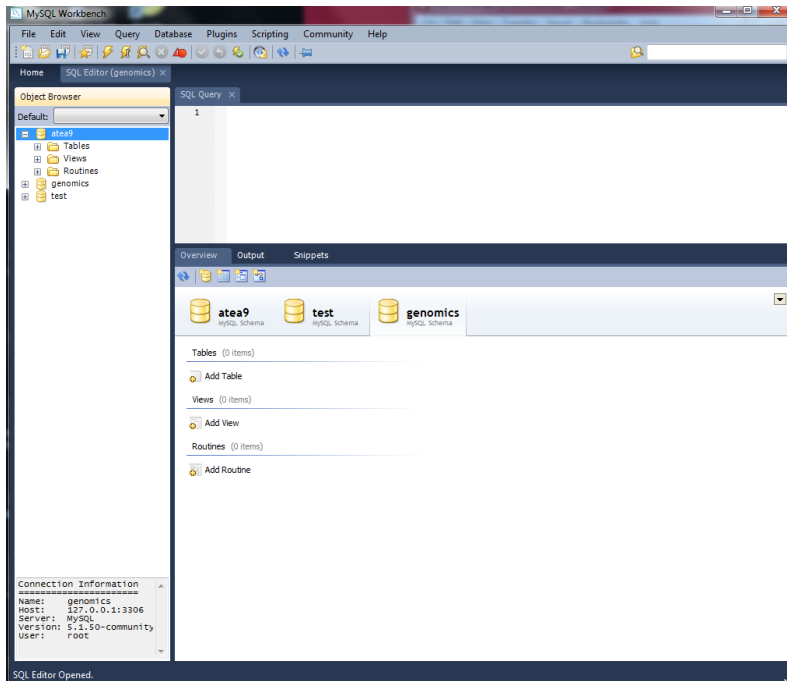**Set the following variable to allow loading local file (this variable must be set to 1 to enable loading data from a local file):**
```
SET GLOBAL local_infile = 1;
```

If you want to, you can change the password to something other than the default or the value you set on install as follows:

```
mysql>SET PASSWORD FOR 'root'@'localhost' =PASSWORD('secret_password');
```

The 'root', which is the username, and 'localhost', which is the hostname, constitute a unique user in MySQL.

'localhost' is a name used to specify the local server. 'root'@'localhost' tells the MySQL server to set the password for a user named 'root' that will connect specifically from the local server (thus 'localhost').

If you changed your password, in order to test the new password, exit the MySQL terminal connection using the following command:

```
mysql>\q
```

This will return you to the system shell. Now log back into the monitor, this time using the following command:

```
%>mysql -u root -p
```

Doing so will result in a prompt for the root user password, as follows:

Enter password:

**Creating and using a database (you likely already did this with the GUI)**

```
mysql> SHOW DATABASES;
information_schema
mysql
test
```

If the `test` database exists, try to access it:

```
mysql> USE test
Database changed
```

Create the genomics database (only do this if you did not create it using the GUI above). In the mysql program:

mysql> create database *genomics*;

```
mysql> SHOW DATABASES;
information_schema
genomics
mysql
test
```

Next, you need to select the database to *use*. Select the genomics database (you can do this not matter which way you created it in the first place):

```
mysql> use genomics;
Database changed
```

**Importing data into a database (use command line for actually importing the data. Either command line or GUI can be used to create the table where the data will go)**

**Note that creating the table can also be done using the command line by following instructions on subsequent pages.**

## To further your understanding, first take a look at the data in excel [record observations about this major step]:

Our goal is to import the *Entrez gene* information database file (a tab-delimited text file) from your *data* directory.

First, we will look at the data in Excel. Open Excel. From within **Excel**, open up the file *gene_info1000*.  In the file dialog, you may need to enable '*All Files'* to see the file. Follow the import steps shown below.

Inspect the different fields. Note that it can be difficult to initially size the fields. *gene_info1000* contains the first 1000 rows of the Entrez Gene – gene information file.

**Text Import Wizard – Step 2 of 3**

This screen lets you set the delimiters your data contains.  You can see how your text is affected in the preview below.

Delimiters

☑ Tab   ☐ Semicolon   ☐ Comma

☐ Space   ☐ Other: [ | ]

☐ Treat consecutive delimiters as one

Text qualifier: [ " ⬍ ]

Data preview

```
7   5692769 NEWENTRY -                           Record to support submis
9   1246500 repA1    pLeuDn_01 -   -   -   -     putative replication-ass
9   1246501 repA2    pLeuDn_03 -   -   -   -     putative replication-ass
9   1246502 leuA     pLeuDn_04 -   -   -   -     2-isopropylmalate syntha
9   1246503 leuB     pLeuDn_05 -   -   -   -     3-isopropylmalate dehydr
9   1246504 leuC     pLeuDn_06 -   -   -   -     isopropylmalate isomeras
```

[ Cancel ]  [ < Back ]  [ Next > ]  [ Finish ]

**Text Import Wizard – Step 3 of 3**

This screen lets you select each column and set the Data Format.

'General' converts numeric values to numbers, date values to dates, and all remaining values to text.

[ Advanced... ]

Column data format

⦿ General

○ Text

○ Date:  [ MDY ⬍ ]

○ Do not import column (Skip)

Data preview

```
Gener General General  General   GenerGenerGenerGener General
7      5692769 NEWENTRY -                            Record to support submis
9      1246500 repA1    pLeuDn_01 -   -   -   -      putative replication-ass
9      1246501 repA2    pLeuDn_03 -   -   -   -      putative replication-ass
9      1246502 leuA     pLeuDn_04 -   -   -   -      2-isopropylmalate syntha
9      1246503 leuB     pLeuDn_05 -   -   -   -      3-isopropylmalate dehydr
9      1246504 leuC     pLeuDn_06 -   -   -   -      isopropylmalate isomeras
```

[ Cancel ]  [ < Back ]  [ Next > ]  [ Finish ]

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 7 | 5692769 | NEWENTRY | - | - | - | - | - | Record to support submission of GeneRIFs for a | other | - | - | - | - | 200 |
| 2 | 9 | 1246500 | repA1 | pLeuDn_01 | - | - | - | - | putative replication-associated protein | protein-coding | - | - | - | - | 200 |
| 3 | 9 | 1246501 | repA2 | pLeuDn_03 | - | - | - | - | putative replication-associated protein | protein-coding | - | - | - | - | 200 |
| 4 | 9 | 1246502 | leuA | pLeuDn_04 | - | - | - | - | 2-isopropylmalate synthase | protein-coding | - | - | - | - | 200 |
| 5 | 9 | 1246503 | leuB | pLeuDn_05 | - | - | - | - | 3-isopropylmalate dehydrogenase | protein-coding | - | - | - | - | 200 |
| 6 | 9 | 1246504 | leuC | pLeuDn_06 | - | - | - | - | isopropylmalate isomerase large subunit | protein-coding | - | - | - | - | 200 |
| 7 | 9 | 1246505 | leuD | pLeuDn_07 | - | - | - | - | isopropylmalate isomerase small subunit | protein-coding | - | - | - | - | 200 |
| 8 | 9 | 1246509 | ibp | pBPS1_01 | - | - | - | - | ibp protein | protein-coding | - | - | - | - | 200 |
| 9 | 9 | 1246510 | repA1 | pBPS1_02 | - | - | - | - | repA1 protein | protein-coding | - | - | - | - | 200 |
| 10 | 9 | 2827857 | NEWENTRY | - | - | - | - | - | Record to support submission of GeneRIFs for a | other | - | - | - | - | 200 |
| 11 | 9 | 3722426 | pLeuDn_02 | pLeuDn_02 | - | - | - | - | ORF1 | protein-coding | - | - | - | hypothetical | 200 |
| 12 | 19 | 3758873 | NEWENTRY | - | - | - | - | - | Record to support submission of GeneRIFs for a | other | - | - | - | - | 200 |
| 13 | 24 | 5129993 | NEWENTRY | - | - | - | - | - | Record to support submission of GeneRIFs for a | other | - | - | - | - | 200 |
| 14 | 33 | 5961931 | pMF1.19c | pMF1.19c | pMX1.19c | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 15 | 33 | 5961932 | pMF1.12 | pMF1.12 | pMX1.12 | - | - | - | putative protein kinase | protein-coding | - | - | - | - | 200 |
| 16 | 33 | 5961933 | pMF1.22 | pMF1.22 | pMX1.22 | - | - | - | putative cobrynic acid A,C-diamide synthase | protein-coding | - | - | - | - | 200 |
| 17 | 33 | 5961934 | pMF1.17 | pMF1.17 | pMX1.17 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 18 | 33 | 5961935 | pMF1.7 | pMF1.7 | pMX1.7 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 19 | 33 | 5961936 | pMF1.6 | pMF1.6 | pMX1.6 | - | - | - | putative mucin-associated surface protein | protein-coding | - | - | - | - | 200 |
| 20 | 33 | 5961937 | pMF1.23 | pMF1.23 | pMX1.23 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 21 | 33 | 5961938 | pMF1.13 | pMF1.13 | pMX1.13 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 22 | 33 | 5961939 | pMF1.1 | pMF1.1 | pMX1.1 | - | - | - | putative APEG precursor protein | protein-coding | - | - | - | - | 200 |
| 23 | 33 | 5961940 | pMF1.9 | pMF1.9 | pMX1.9 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 24 | 33 | 5961941 | pMF1.5 | pMF1.5 | pMX1.5 | - | - | - | putative RTX toxin-related Ca2+-binding protein | protein-coding | - | - | - | - | 200 |
| 25 | 33 | 5961942 | pMF1.10 | pMF1.10 | pMX1.10 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 26 | 33 | 5961943 | pMF1.21 | pMF1.21 | pMX1.21 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 27 | 33 | 5961944 | pMF1.8 | pMF1.8 | pMX1.8 | - | - | - | glycoside hydrolase family protein 19 | protein-coding | - | - | - | - | 200 |
| 28 | 33 | 5961945 | pMF1.20c | pMF1.20c | pMX1.20c | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 29 | 33 | 5961946 | pMF1.4 | pMF1.4 | pMX1.4 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 30 | 33 | 5961947 | pMF1.14 | pMF1.14 | pMX1.14 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 31 | 33 | 5961948 | pMF1.2 | pMF1.2 | pMX1.2 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 32 | 33 | 5961949 | pMF1.18 | pMF1.18 | pMX1.18 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 33 | 33 | 5961950 | pMF1.16 | pMF1.16 | pMX1.16 | - | - | - | putative resolvase | protein-coding | - | - | - | - | 200 |
| 34 | 33 | 5961951 | pMF1.11 | pMF1.11 | pMX1.11 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 35 | 33 | 5961952 | pMF1.15 | pMF1.15 | pMX1.15 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 36 | 33 | 5961953 | pMF1.3 | pMF1.3 | pMX1.3 | - | - | - | hypothetical protein | protein-coding | - | - | - | - | 200 |
| 37 | 33 | 5999886 | NEWENTRY | - | - | - | - | - | Record to support submission of GeneRIFs for a | other | - | - | - | - | 200 |
| 38 | 34 | 4126706 | NEWENTRY | - | - | - | - | - | Record to support submission of GeneRIFs for a | other | - | - | - | - | 200 |
| 39 | 56 | 5814795 | NEWENTRY | - | - | - | - | - | Record to support submission of GeneRIFs for a | other | - | - | - | - | 200 |

gene_info1000

Normal View    Ready    Sum=7

Back in the MySQL GUI, Open a new SQL query tab.



*Create a database table* to accommodate the data from *gene_info.* Creating a database table with the correct data types and sizes for fields is often a trial and error process.

[major step] You can paste the SQL below directly in the SQL Query window and execute it, or you can use the *Add Table* wizard under your *genomics* MySQL schema tab, and enter each field individually. *Tip: Watch out for special characters from Word. Especially quotes.*

```
create table gene_info (
tax_id int,
GeneID int,
Symbol varchar(48),
LocusTag varchar(48),
Synonyms varchar(1000),
dbXrefs varchar(512),
chromosome varchar(48),
map_location varchar(48),
description varchar(4000),
type_of_gene varchar(12),
Symbol_from_nomenclature_authority varchar(48),
Full_name_from_nomenclature_authority varchar(256),
Nomenclature_status varchar(12),
Other_designations varchar(12),
Modification_date varchar(12));
```

To execute your SQL, make sure the genomics database is selected in the Object Browser (left-side panel in Workbench should be **bold**).

Next, select the entire text of the query, and execute the query by selecting the lightning bolt to the left (execute ALL or selected), *or* just click (select) anywhere within your

query text and select the lightning bolt to the right with the big "I" superimposed on it (execute current query).

You can also execute either type of query form the menu under SQL.

You should get the following message:
Create table gene_info(….) 0 rows affected.

You can see the table now in the object browser under your genomics->Tables->gene_info but you may have to right click and refresh the Tables to have it show up.

MySQL reference for create table:
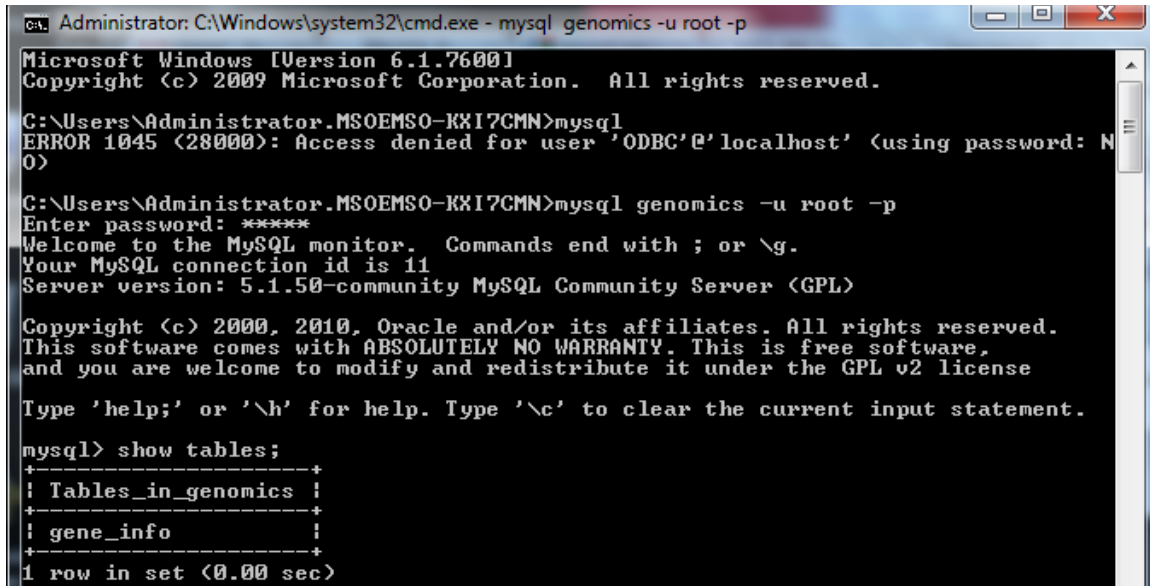http://dev.mysql.com/doc/refman/8.0/en/show-create-table.html

Next perform the *load data* command using the command line SQL query window as follows. I recommend using the command line since workbench can be too picky with file imports.

Command line: Open up a *command line* connection to mySQL as described previously.

And enter your password at the prompt.

Type *show tables* to verify you are connected to the correct database, etc.



Execute:
SELECT @@GLOBAL.`local_infile`;
Make sure it returns a 1.

If not, enter the following:
SET GLOBAL local_infile=1

[major step] Enter the following *load data local infile* command – you have to change it to the location of the file on your system (no carriage returns, unix-style slashes!):

```
Load Data local Infile
'D:/Dropbox/MSOE/www/cs386/labs/lab1/entrez/gene_info1000'
into table genomics.gene_info fields terminated by '\t' lines
terminated by '\n';
```

If you get "ERROR 2068 (HY000): LOAD DATA LOCAL INFILE file request rejected due to restrictions on access." Then type exit on the command prompt.
Reopen mysql command prompt with following (so not via the link on the start menu):
`mysql -u root -p --local-infile=1`  and try the load data command again after running "use genomics;"

*Note: Always check the data file for column delimiters and end-of-line character(s).*
- *In Unix and all Unix-like systems (OS-X), \n is the code for end-of-line, \r means nothing special.*
- *In Windows (and some old OS's), the code for end of line is 2 characters, \r\n, in this order (CR, LF)*

```
mysql> Load Data local Infile 'D:/Jay/MSOE/www/cs386/labs/lab1/entrez/gene_info1
000' into table genomics. gene_info fields terminated by '\t' lines terminated b
y '\n';
Query OK, 1000 rows affected, 1103 warnings (0.10 sec)
Records: 1000  Deleted: 0  Skipped: 0  Warnings: 1103
```

*Note: the large number of warnings! In MySQL workbench, you hover over the message column with your cursor in the output window towards the bottom.*

*Note: gene_info1000 is the first 1000 rows of the Entrez gene database table. We will import the whole table after we fix a few problems first.*

*Note: I deleted the header block from the top of the text file.*

MySQL reference for load data: http://dev.mysql.com/doc/refman/8.0/en/load-data.html

[major step] In your MySQL workbench query window (or *cmd* window) execute:

```
show warnings;
```

You should see several error messages as follows:
```
Data truncated for column 'type_of_gene' at row 2
Data truncated for column 'type_of_gene' at row 3
…
Data truncated for column 'Other_designations' at row 11
…
```

This indicates that we did not size the *type_of_gene* and *Other_designations* database fields correctly. This may not be important if you do not care about the data in these fields, but let's assume you want to import the entire file correctly.

From the command line, or your MySQL Workbench SQL window in the GUI:

First, drop (delete) the current gene_info  database table:

```
drop table gene_info;
```

[major step] Create a new table with larger fields for *type_of_gene* and *Other_designations.* We see from inspecting the file, that neither of these entries is very commonly supplied. A typical entry for *type_of_gene* is *'protein coding',* and *Other_designations* may have an entry like '*putative F plasmid gene 32-like protein'*.

```
create table gene_info (
tax_id int,
GeneID int,
Symbol varchar(48),
LocusTag varchar(48),
Synonyms varchar(1000),
dbXrefs varchar(512),
chromosome varchar(48),
map_location varchar(48),
description varchar(4000),
type_of_gene varchar(24),
Symbol_from_nomenclature_authority varchar(48),
Full_name_from_nomenclature_authority varchar(256),
Nomenclature_status varchar(12),
Other_designations varchar(4000),
Modification_date varchar(12));
```

[major step] Reload the entire full length Entrez gene data set (not just the 1000 record version, and be a little patient this can take a couple minutes) using load data local

```
Load Data local infile
'D:/Dropbox/MSOE/www/cs386/labs/lab1/entrez/gene_info' into table
gene_info fields terminated by '\t' lines terminated by '\n' IGNORE 1
LINES;
```

```
Note: I added the "IGNORE 1 LINES" option to skip the first line on
when loading the whole file. For the abbreviated file gene_info1000, I
had removed the column header line. In the original file, the header is
still there and we need to deal with it.
```

[major step] Check for warnings and verify the length of the database table. You may still see some warnings about rows not having all columns, which you can ignore.

```
select count(*) from gene_info;
4354758
```

**Creating our working table [major steps – entire section, but feel free to abbreviate results shown in report]**

Issue the following select statements, limiting the result set to the first 100 records, and inspect the table.

```
select * from gene_info
limit 100;
```

MySQL reference for select: http://dev.mysql.com/doc/refman/8.0/en/select.html

We can see that several fields including *dbXrefs*, *chromosome*, *map_location*, and others are typically blank, and may not be needed for our experiments. So let's create a new table with just the information we need.

```
create table gene_info2 (
tax_id int,
GeneID int,
Symbol varchar(48),
LocusTag varchar(48),
Synonyms varchar(1000),
description varchar(4000),
type_of_gene varchar(24),
Modification_date varchar(12));
```

We can use the SQL *insert into* statement to populate our new table from our existing table (takes a couple minutes to run) :

```
insert into gene_info2
select tax_id, GeneID, Symbol, LocusTag, Synonyms, description,
type_of_gene, Modification_date
from gene_info
where GeneID is not null;
```

Note: We are not including rows where there is no GeneID with the statement:
```
where GeneID is not null;
```
Without being able to uniquely identify a record, its useless.

MySQL reference to insert:
**http://dev.mysql.com/doc/refman/8.0/en/insert.html**

We no longer need the original *gene_info* table we used for importing gene data, so we can drop it.

**drop table gene_info;**

*Note: To delete all rows from a table without removing a table:*
**delete from gene_info;**
Drop and create table is generally much faster than delete for large tables.

**Primary keys and indexes [major steps – entire section, but feel free to abbreviate results shown in report]**

Issue the following select statement to retrieve the record for a specific gene, and record the time:
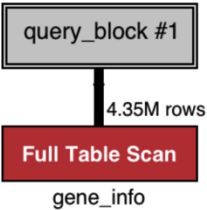
```
select * from gene_info2 where GeneID=4226999;

303, 4226999, 'pWW53_43', 'pWW53_43', '-', 'putative reductase
protein', 'protein-coding', '20080514'

1 row fetch (12.4951s)
```

This next step should be done in the GUI, even if you have been using the command line! **Re-execute the command with the explain option by selecting the *Lighting Bolt* with the *magnifying glass* from the Toolbar.** Explain will display how the table is accessed. You may need to click on the various GUI buttons arraigned vertically on the right side of the results to see this figure and text (view source). ***Record the result.***

```
explain select * from gene_info2 where GeneID=4226999;
```



```
"query_block": {
  "select_id": 1,
  "table": {
    "table_name": "gene_info",
    "access_type": "ALL",
    "rows": 4354758,
    "filtered": 100,
    "attached_condition":
"(`genomics`.`gene_info`.`GeneID` =
4226999)"
  }
}
```

```
Full table scan. I.e., we are scanning every record in the database.
```

This is very *slow* for retrieving a record. Also, we want to ensure that our table is in first normal form. At a minimum, this means we need to be able to address each record/row of a database table uniquely using a single field or some combination of fields. We do this by defining the *primary key* for the table, which is GeneID. When you define a primary key, the DBMS automatically creates an index on this key.

An index reduces search time by creating a tree-like or hash data structure to access rows in a database table. Importing data into a table goes much faster if you define the *primary key* and other indexes *after* the data has been imported. Otherwise the indexes have to be verified and updated with each row insert.

alter table gene_info2 add constraint pk_gene_info primary key(GeneID);

```
select * from gene_info2 where GeneID=4226999;
```

```
303, 4226999, 'pWW53_43', 'pWW53_43', '-', 'putative reductase
protein', 'protein-coding', '20080514'

1 row fetch (0.0004s)
```

That's more like it. *Note: You can create primary keys and indexes as a composite of multiple fields as well.*

We may also want to create indexes on other fields in the database. For example we may be interested in identifying all genes within a specific taxonomy.  So, let's create an index on *tax_id*. Note: the *tax_id* is not distinct.

```
create index tax_id_idx on gene_info2(tax_id);
```

Next let's import the gene2pubmed data. For large data imports, you may get a time-out error in the GUI, which is why we use the command line to import data. The gene2pubmed table allows us to cross-reference Pubmed articles. Format: tax_id, GeneID, PubMed_ID

```
create table gene2pubmed (
tax_id int,
GeneID int,
PubMed_ID int);

Load Data local Infile '/data/gene2pubmed' into table
genomics.Gene2pubmed fields terminated by '\t' lines terminated by
'\n';
```

### ERD
Create an ERD diagram (We learn about this in the 2^nd lecture) of gene_info2, gene2pubmed and the following Pubmed table (we are not going to populate it right now):
```
Create table Pubmed (
Pubmed_ID int,
Text varchar(4000));
```

You can use *desc* to display each table's data types as follows:
```
desc gene2pubmed;
```

**https://en.wikipedia.org/wiki/Entity%E2%80%93relationship_model**

### Deliverables
Complete a lab report containing the following:
- What you've learned from the lab
- How it could be improved
- Observations and query results for each major step of the lab.
- ERD
- Submit to Canvas in PDF format!