

# How to manage HGMD database in MySQL

## Contents

- 1 Current maintenance
- 2 Raw data from HGMD (Jan 2021)
- 3 Summary statistics after pre-processing
- 4 Link to processed data
- 5 How to process HGMD raw data
- 6 How to create and drop a MySQL database
- 7 How to create and delete a table in an existing MySQL database
- 8 How to import .sql file to an existing MySQL database
- 9 How to export a MySQL table to a .csv file
- 10 Basic commands to navigate MySQL on cbsuhy01

## Current maintenance

- Administrator: Yingying Zhang (yz2296@cornell.edu)

## Raw data from HGMD (Jan 2021)

- Phenotypes:
  - HGMD\_Phenbase\_2021\_1.sql
- Mutations:
  - hg38: 2021\_1\_hg38\_fullinfo.vcf
  - hg19: 2021\_1\_hg19\_fullinfo.vcf

## Summary statistics after pre-processing

- Missense mutations (148,114):
  - Mapped to UniProt: 145,332
  - Unmapped: 2,782
- Nonsense mutations (34,405):
  - Mapped to UniProt: 33,168
  - Unmapped: 1,237
- Coding indels (64,392):
  - Mapped to UniProt: 61,491
  - Unmapped: 2,901
- Non-coding mutations (34,945)
- Dropped mutations with non-compliant format (5,288)

## Link to processed data

- Indels in coding regions mapped to UniProt: File:2021\_2\_mapped\_codingindel.xlsx
- Missense mutations mapped to UniProt: File:2021\_2\_mapped\_missense.xlsx
- Nonsense mutations mapped to UniProt: File:2021\_2\_mapped\_nonsense.xlsx
- Non-coding mutations: File:2021\_2\_noncoding.xlsx
- Indels in coding regions which could not be mapped to UniProt: File:2021\_2\_unmapped\_codingindel.xlsx
- Missense mutations which could not be mapped to UniProt: File:2021\_2\_unmapped\_missense.xlsx
- Nonsense mutations which could not be mapped to UniProt: File:2021\_2\_unmapped\_nonsense.xlsx
- Raw data from HGMD: File:HGMD\_concept\_def.xlsx
- Raw data from HGMD: File:HGMD\_mutation.xlsx
- Raw data from HGMD: File:HGMD\_phenotype.xlsx

## How to process HGMD raw data

- Create a MySQL database called HGMD2021\_1 (if not existed). Please find how to do so in #How to create a MySQL database.
- Import HGMD\_Phenbase\_2021\_1.sql to HGMD2021\_1. Please find how to do so in #How to import .sql file to an existing MySQL database. Add "HGMD\_" as prefix to all imported tables. This is to distinguish original HGMD tables from the ones we generated.
- Process 2021\_1\_hg38\_fullinfo.vcf:
  1. The 7th column (index from 0) contains the information about each mutation. Parse the 7th column into 8 columns: Variant\_class, Gene, Strand, dbSNP, Phenotype, Rank\_score, RefSeq\_DNA\_record, and RefSeq\_protein\_record. If the 7th column does not contain RefSeq protein record (for example noncoding variants), set the value in RefSeq\_protein\_record to nan. For example, if the content in the 7th column is "CLASS=DM?;MUT=ALT;GENE=SAMD11;STRAND=+;DNA=NM\_152486.4%3Ac.133A>G;PROT=NP\_689699.3%3Ap.K45E;DB=rs903331232;PHEN="Retinitis\_pigmentosa, 0.21, NM\_152486.4:c.133A>G, NP\_689699.3:p.K45E. Please notice that I replaced "%3A" with ":" for the values in RefSeq\_DNA\_record and RefSeq\_protein\_record. The values in RefSeq\_protein\_record stands for the amino acid change caused by the mutation. In this example, NP\_689699.3:p.K45E stands for an amino acid substitution from K to E at position 45 in the protein NP\_689699.3.
  2. Add one more column POS\_hg19 based on 2021\_1\_hg19\_fullinfo.vcf. Map hg19 coordinates to hg38 coordinates using the HGMD ID of each variant.
  3. Rename the columns and rearrange them so that you will get a dataframe with the following columns: CHR, POS, HGMD\_ID, REF\_N, ALT\_N, POS\_hg19, Variant\_class, Gene, Strand, dbSNP, Phenotype, Rank\_score, RefSeq\_DNA\_record, RefSeq\_protein\_record.

4. Pick out missense mutations from the parsed dataframe:
  - Missense mutations always have the format of `r'[:-p.[A-Z]d+[A-Z]'` in `RefSeq_protein_record` (this is a python regular expression, please google to see what it means). Pick out all rows in the parsed dataframe where this pattern could be found in `RefSeq_protein_record`, and add two more columns `REF_AA` and `ALT_AA` to these rows, one for reference amino acid and one for altered amino acid. In the aforementioned example, the values in these two columns should be K and E.
  - Map missense mutations to protein positions using the pipeline written by Charles Liang: Genomics and protein coordinate conversion. After this step you should get two text files, **mapped\_missense.txt** and **unmapped\_missense.txt**. The columns in `mapped_missense.txt` should be: `CHR`, `POS`, `HGMD_ID`, `REF_N`, `ALT_N`, `POS_hg19`, `REF_AA`, `ALT_AA`, `UNIPROT`, `UNIPROT_position`, `Variant_class`, `Gene`, `Strand`, `dbSNP`, `Phenotype`, `Rank_score`, `RefSeq_DNA_record`, `RefSeq_protein_record`. These columns should also appear in `unmapped_missense.txt` except for `UNIPROT` and `UNIPROT_position`.
5. Pick out nonsense mutations from the parsed dataframe:
  - Nonsense mutations always have the format of `r'[:-p.[A-Z]d+[*]'` in `RefSeq_protein_record`. Pick out all rows in the parsed dataframe where this pattern could be found in `RefSeq_protein_record`, and repeat the same procedures as what you did when picking out missense mutations. Finally you should get two text files, **mapped\_nonsense.txt** and **unmapped\_nonsense.txt**.
6. Pick out coding indels from the parsed dataframe:
  - From the parsed data frame, pick out variants whose `REF_N` and `ALT_N` have different length and `RefSeq_protein_record` is not nan. Put these variants into `VEP(GRCh38)`: make sure to check the "Protein" and "UniProt" boxes under the "Identifiers" section, and also make sure to check "Transcript biotype" under the "Transcript annotation" section. In the "Filters" section, check the box saying "Return results for variants in coding regions only". In "Restrict results", select "Show one selected consequence per variant".
  - After getting VEP output, filter the output file and retain the rows where the `SWISSPROT` and `TREMBL` columns are not be both empty. There should be only one identifier in `SWISSPROT` and might be multiple identifiers in `TREMBL` (if so separated by ", "). If the `SWISSPROT` or `TREMBL` identifiers are in the format of `XXXXXX.XXX` (for example `P43489.17`), only retain the content before "." (P43489). Then choose the best UniProt among `SWISSPROT` and `TREMBL` identifiers for each row, prioritizing `SWISSPROT` over `TREMBL`, and choosing the longest `TREMBL` identifier if there is no `SWISSPROT`. Naming the column containing the best UniProt as `UNIPROT`. Rename the `Protein_position` column as `UNIPROT_position`.
  - Parsing the `Amino_acids` column in the VEP output file into two columns: `REF_AA` and `ALT_AA`. For example, if the value in `Amino_acids` is "RLTQTV/X", the values in `REF_AA` and `ALT_AA` should be "RLTQTV" and "X" respectively.
  - Retain the six columns in VEP output: `#Uploaded_variation`, `UNIPROT`, `UNIPROT_position`, `REF_AA`, `ALT_AA` and `Consequence`, and map the VEP output to the coding indels using `#Uploaded_variation` as key. Finally you should get a text file named **mapped\_coding\_indel.txt** with 19 columns: `CHR`, `POS`, `HGMD_ID`, `REF_N`, `ALT_N`, `POS_hg19`, `REF_AA`, `ALT_AA`, `UNIPROT`, `UNIPROT_position`, `Variant_class`, `Gene`, `Strand`, `dbSNP`, `Phenotype`, `Rank_score`, `RefSeq_DNA_record`, `RefSeq_protein_record`, `Consequence`. Also you should get **unmapped\_coding\_indel.txt** with 14 columns (missing `REF_AA`, `ALT_AA`, `UNIPROT`, `UNIPROT_position`, and `Consequence`).
7. Pick out noncoding variants from the parsed dataframe:
  - From the parsed data frame, pick out variants whose `RefSeq_protein_record` is nan. Put the selected rows into **noncoding.txt**.
8. Put the remaining variants that do not belong to any of the aforementioned categories into **other\_variants.txt**.
9. Import all text files you generated from the previous steps as tables into `HGMD2021_1`. Please find how to do so in #How to create and delete a table in an existing MySQL database.

## How to create and drop a MySQL database

- Please contact Yingying Zhang (yz2296@cornell.edu) for the MySQL account and password
- Download sequel pro to your computer: <http://sequelpro.com/>
- When the first time you launch the sequel pro application, the following window will pop up:

Type in "127.0.0.1" to "MySQL Host".

Type in the MySQL account username and password to "Username" and "Password".

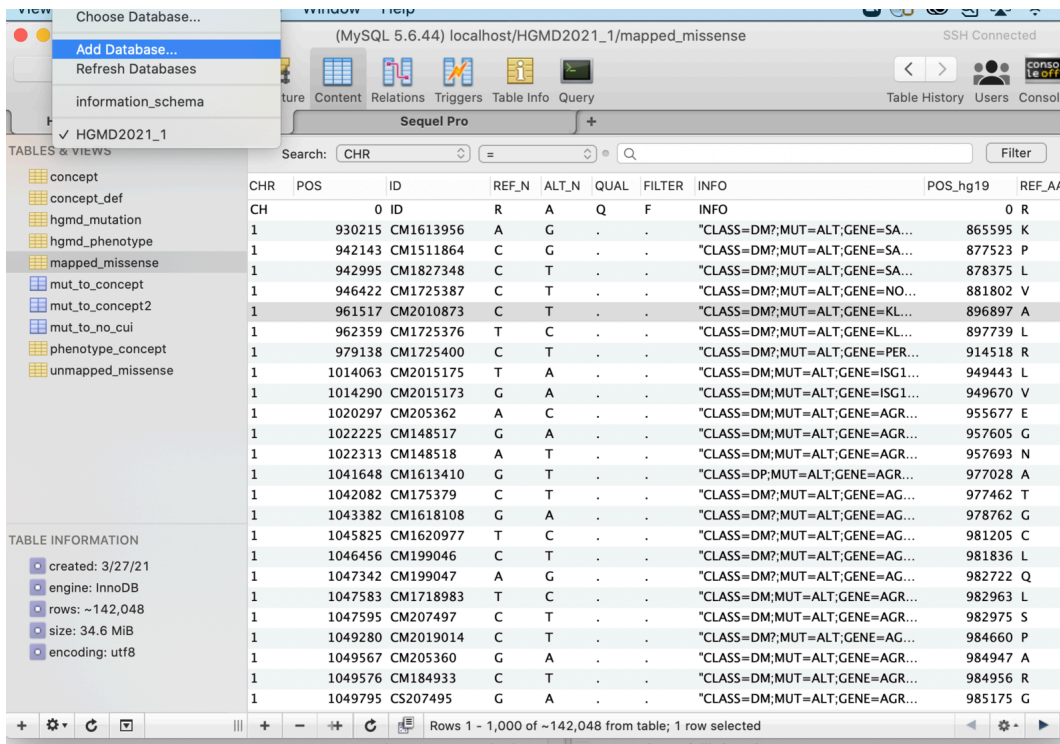
Type in "cbsuhy01.biohpc.cornell.edu" to "SSH Host"

Type in your NetID to "SSH User"

Type in your ssh password to cbsuhy01 to "SSH Password"

Then you should have successfully logged in to MySQL on cbsuhy01

- If you want to create a new database, click the "Choose Database..." button at the upper left corner of the sequel pro interface, and you will see a drop-down list, among which there is an entry called "Add Database...". By clicking it you can create a new database and specify the name.



- If you want to drop a database, use the following command:

```
DROP DATABASE database_name;
```

## How to create and delete a table in an existing MySQL database

- If you want to create a new table in an existing database, for example a table "STUDENT\_INFO" composed of two columns "NAME" and "AGE" in a database "STUDENTS\_2021", use the following command:

```
mysql> use STUDENTS_2021;
mysql> CREATE TABLE STUDENT_INFO (NAME varchar(20), AGE Int);
```

- If you want to import a text file to an existing table, for to import student\_info.txt to the table STUDENT\_INFO, use the following command (replace /path/to/student\_info.txt with the path to the text file that you want to import. If your text file does not have a header line, remove "IGNORE 1 LINES" from the following command):

```
mysql> LOAD DATA LOCAL INFILE '/path/to/student_info.txt' INTO TABLE STUDENT_INFO IGNORE 1 LINES;
```

- If you want to permanently delete a table, type the following command (Replace "table\_name" with the real table name):

```
mysql> DROP TABLE table_name;
```

## How to import .sql file to an existing MySQL database

- If you want to import a .sql file to an existing MySQL database, please login to cbsuhy01 through your terminal, go to the folder where the .sql file is stored, and then type in the following command:

```
sed 's/\sDEFINER=[^']*@`[^']*`//g' -i your_sql_file.sql
```

Please replace "your\_sql\_file.sql" with the name of the .sql file that you want to import. The command above is to remove the "DEFINER=.." statement from your .sql file, so that you won't run into any access problem.

Then type in the following command to import the .sql file to MySQL database:

```
mysql -u username -p database_name < your_sql_file.sql
```

Please replace "username", "database\_name" with your MySQL account user name and the existing MySQL database name, respectively. Please also replace "your\_sql\_file.sql" with the name of the .sql file that you want to import. As you've already in the right folder, there is no need to include a file path.

## How to export a MySQL table to a .csv file

```
mysql -B -u username -ppassword database_name -h localhost -e "SELECT * from table_name" | sed "s/'/\"/;s/\\t/\\n/g;s/\\/\"/;s/\\n/\\n/g" > /path/to/your/csv_file
```

Please replace "username" and "password" with the MySQL account username and password from Yingying Zhang. Please notice that there is no space between "-p" and your password [IMPORTANT]. Replace "database\_name" and "table\_name" with the real name of the database and table that you want to export. Replace "/path/to/your/csv\_file" with the path to the output .csv file.

## Basic commands to navigate MySQL on cbsuhy01

- ssh login to cbsuhy01:

```
net_id@cbsuhy01.biohpc.cornell.edu (please replace "net_id" with your own NetID)
```

- Type in your terminal:

```
mysql -u username -p
```

Please replace "username" with your MySQL account user name and then type in your MySQL account password. If your account exists, the following information will pop up:

```
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 26
Server version: 5.6.44 MySQL Community Server (GPL)
Copyright (c) 2000, 2019, Oracle and/or its affiliates. All rights reserved.
Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.
Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.
mysql>
```

To get started on your own database, first check which databases currently exist:

```
mysql> show databases;
```

To the select a database, issue the use command (Replace "database\_name" with the real database name that you want to use):

```
mysql> use database_name;
mysql> show tables;
```

Retrieved from "[http://wiki.yulab.org/index.php?title=How\\_to\\_manage\\_HGMD\\_database\\_in\\_MySQL&oldid=2858](http://wiki.yulab.org/index.php?title=How_to_manage_HGMD_database_in_MySQL&oldid=2858)"

- This page was last edited on 24 May 2021, at 14:25.