
esgemme

Release 1.3.0

Mustafa Tekpinar

Jun 22, 2023

CONTENTS:

1	Using ESGEMME via Docker	1
1.1	Requirements	1
1.2	Getting the example input data	1
1.3	Single point mutation calculations	2
1.4	Multiple point mutation calculations	3
1.5	Running several jobs using docker	3
2	Analyzing and Modifying the ESGEMME Output	5
2.1	Raw ESGEMME Scores and Their Interpretation	5
2.2	Entire Single Point Mutation Landscape Calculations	5
2.3	Selected Single Point or Multiple Point Mutation Calculations	6
3	Preparing Your Own Input	7
3.1	Preparing Your Input MSA and PDB with Colabfold	7
4	Installation	9
4.1	Installing the dependencies:	9
4.2	Preparation of the environment and installation of ESGEMME	9
4.3	Configuring default.conf file	10
5	Indices and tables	13

USING ESGEMME VIA DOCKER

1.1 Requirements

You need to have docker installed on your machine. You can consult the following page for this: <https://docs.docker.com/get-docker/>

I am assuming some basic familiarity with Linux/Unix/macOS terminal commands.

Let's start our favorite terminal app.

You must create a folder called docker-tutorial and go to that empty folder:

```
mkdir docker-tutorial
cd docker-tutorial
```

1.2 Getting the example input data

Let's download the sample data provided in the ESGEMME repository for this exercise. First, we will download the multiple sequence alignment file in fasta format:

```
wget http://gitlab.lcqb.upmc.fr/tekpinar/ESGEMME/raw/
↳8d766d4d11af0e93c9da8fc2c5cc1bfc457d2936/data/aliBLAT.fasta
```

If you don't have wget, you can try the same command with curl:

```
curl http://gitlab.lcqb.upmc.fr/tekpinar/ESGEMME/raw/
↳8d766d4d11af0e93c9da8fc2c5cc1bfc457d2936/data/aliBLAT.fasta >aliBLAT.fasta
```

Please verify that the aliBLAT.fasta file is in the folder.

Now, we will download the PDB (Protein Databank) file for BLAT:

```
wget http://gitlab.lcqb.upmc.fr/tekpinar/ESGEMME/raw/
↳8d766d4d11af0e93c9da8fc2c5cc1bfc457d2936/data/blat-af2.pdb
```

1.3 Single point mutation calculations

In order to make sure that the docker is installed:

```
sudo docker -h
```

If it shows you a list of options, you are on a good track. On MacOS, you may not need ‘sudo’ word before the docker command at all.

```
sudo docker run -ti --rm --mount type=bind,source=$PWD,target=/home/tekpinar/research/
↳ myexample \
tekpinar/esgemme-docker:v1.3.0
```

You are in the container (your virtual operating system) now. You created a folder called myexample in your container with the previous command. Let’s change to that folder.

```
cd ../myexample/
```

When you check the data in that folder with ‘ls’ command, you are supposed to see aliBLAT.fasta and blat-af2.pdb files. Basically, your docker-tutorial folder on the host system and myexample folder on the docker container are pointing to the same place.

1.3.1 Obtaining the entire single point mutation landscape

In this step, we will use only evolutionary information from an MSA file:

```
esgemme aliBLAT.fasta -r input -f aliBLAT.fasta
```

After a few minutes of calculation, you must see at least two files named BLAT_normPred_evolCombi.txt and BLAT_normPred_evolCombi.png. You have the entire single point mutational landscape of BLAT protein in these files.

If you want to utilize structural information (highly recommended) as well as evolutionary information:

```
esgemme aliBLAT.fasta -r input -f aliBLAT.fasta \
--pdbfile blat-af2.pdb \
--normweightmode sstjetormax
```

1.3.2 Predicting the effect of a subset of single point mutations

If you are interested in only a bunch of single point mutations, you have to prepare a mut file. The format is a simple text file and each line contains a single point mutation such as D26A.... Fortunately, we have an example mut in data folder of ESGEMME repository.

```
wget http://gitlab.lcqb.upmc.fr/tekpinar/ESGEMME/raw/master/data/Stiffler_2015_BLAT_
↳ ECOLX.mut
```

Similar to the previous step, there are two possible ways to do the calculations: with or without structural information. First, let’s do it without structural information:

```
esgemme aliBLAT.fasta -r input -f aliBLAT.fasta \
-m Stiffler_2015_BLAT_ECOLX.mut
```

You can include structural information in the following way: `.. code:: bash`

```
esgemme aliBLAT.fasta -r input -f aliBLAT.fasta --pdbfile blat-af2.pdb --normweightmode sstjetormax -m  
Stiffler_2015_BLAT_ECOLX.mut
```

You will have BLAT_normPred_evolCombi.txt file in your folder. However, the output format is completely different from the entire mutational landscape scanning file. Each line of this file is a mutation and its predicted effect separated by a space. In addition, you won't have a png file like in the previous case.

1.4 Multiple point mutation calculations

Sometimes, we need to see effects of double or triple mutations. ESGEMME can perform calculations if you provide a mut file. In this case, the mut file must have the following format:

E26D:Y44R E56N:A77F:H94V

The first line of the text file is impact of a double mutation and the second line is the impact of the triple mutations. As you can see, the mutations are separated by a colon(:) character. The output file will be in a similar format. Each line will contain the multiple mutation and its predicted effect, separated by a space.

1.5 Running several jobs using docker

ANALYZING AND MODIFYING THE ESGEMME OUTPUT

2.1 Raw ESGEMME Scores and Their Interpretation

There is not a hardcoded limit for raw ESGEMME scores. However, the values range between $[-12, 2]$ generally. The lower values mean the mutations is impactful, while values close to 0 means the mutation does not have any significant impact.

We should note that most of the ‘impactful’ mutations are deleterious but it is not always the case.

2.2 Entire Single Point Mutation Landscape Calculations

By default, ESGEMME will only output the combined (independent and epistatic) scores*:

There are three output files:

1. myProt_normPred_evolCombi.txt

myProt is the short name in the MSA file for your protein. The ‘myProt_normPred_evolCombi.txt’ file contains 20 rows (for 20 amino acids in alphabetical order) and L columns, where L is the number of amino acids in your protein of interest. Since this file is horizontal, it is easy to read it in R or Python but difficult to find the mutations you are interested.

2. myProt_normPred_evolCombiTransposed.txt

As the name implies, this is the transposed version of the combined results. It is easier to find the mutations you are interested in this file. Just check the row corresponding to the mutation.

3. myProt_normPred_evolCombi.png

This is the image file of the combined results. It selects ‘Oranges_r’ matplotlib color map by default. You can change it by adding ‘--colormap turbo_r’ for a more fancy look during the esgemme call. It --colormap argument accepts all the color maps in matplotlib.

Note*: If you want to see epistatic and independent contributions as well, you should add ‘--verbose true’ argument while calling esgemme.

2.3 Selected Single Point or Multiple Point Mutation Calculations

Since you used a mutation file to predict mutations, you will have your combined (epistatic+independent) results in the same format, such as :

*. **myProt_normPred_evolCombi.txt:**

A2C -6.23 A2D -1.23 . . . D286F -0.23

PREPARING YOUR OWN INPUT

3.1 Preparing Your Input MSA and PDB with Colabfold

You have a fasta file for your protein of interest and you want to understand impact of (certain) mutations. Before starting, please make sure that your fasta file does not contain a gap. The quickest method to obtain both multiple sequence alignment and a protein structure is to use Colabfold. Let's do this step by step:

1. Let's go the Colabfold web site:

<https://colab.research.google.com/github/sokrypton/ColabFold/blob/main/AlphaFold2.ipynb>

Sign in using your gmail account.

2. Click on the 'Connect' button on the top right hand side.
3. Clean 'query_sequence' box and paste your sequence to the 'query_sequence' box. For me, I selected adenylate kinase (AKE) as my example fasta sequence (<https://www.rcsb.org/fasta/entry/4AKE/display>).
4. Change the 'jobname' to something that makes more sense to you.
5. Go to the menu bar of your 'AlphaFold2.ipynb' notebook, where 'File, Edit, View, Insert, Runtime, Tools, Help' are listed. Click on the Runtime and select 'Run all'.
6. This process make take from a few minutes to a few hours depending on your protein size. It will give you an a3m file and up to 5 PDB models. Put these files in a clean folder and change the directory to that folder in your terminal.
7. Unfortunately, a3m file is not in fasta format and it contains gap columns. We have to clean those gaps. We can do that with a GUI program like Ugene or Jalview. However, it is a labor intensive procedure. Here, I will use a small tool that I developed and added to the ESGEMME docker image that I created.
8. Start the docker image with the following command:

```
sudo docker run -ti --rm --mount type=bind,source=$PWD,target=/home/tekpinar/  
→research/myexample \  
tekpinar/esgemme-docker:v1.3.0
```

9. Now, change the directory to myexample folder.

```
cd ../myexample/  
ls -l
```

We are supposed to see our a3m and pdb files in this folder.

10. Let's use a small script from hhsuite to convert a3m file to fasta format.

```
reformat.pl a3m fas AKE.a3m AKE.fasta
```

11. Final step and we are there:

```
demust removegaps -i AKE.fasta -o AKE_nogaps.fasta
```

There is one last step to reach our goal. ID and description parts of the a3m and fasta files are too long. We have to shorten them. We can do that with

```
awk 'BEGIN{FS=" "}{{if(NF>1) {printf(">%s\n", $1)}else{print $0}}}' AKE_nogaps.  
↪ fasta > AKE_nogaps_short_names.fasta  
# Recheck this command if you can remove extra >
```

Congratulations! Now, you have all the input files required for ESGEMME: I. An input MSA: AKE_nogaps_short_names.fasta II. An input PDB: myprotein.pdb

INSTALLATION

ESGEMME is implemented in Python 3 and R. It has been tested only on Linux. Since ESGEMME has many dependencies, we recommend using our web site or our docker image. If you are a determined user, here comes the steps required to install it from the source.

4.1 Installing the dependencies:

ESGEMME has the following external dependencies:

- Joint Evolutionary Trees: <http://www.lcqb.upmc.fr/JET2/> and its dependencies:
 - java
 - naccess: <http://www.bioinf.manchester.ac.uk/naccess/>
 - muscle: <https://www.drive5.com/muscle/>
- seqinr R package: <https://cran.r-project.org/web/packages/seqinr/index.html>
- dssp for secondary structure prediction.

These tools should be installed to be able to use ESGEMME.

4.2 Preparation of the environment and installation of ESGEMME

Step by step installation on Ubuntu 22.04

Prepare your environment and install the required packages:

```
sudo apt-get update --fix-missing && \  
sudo apt-get install -y --no-install-recommends apt-utils && \  
sudo apt-get install -y software-properties-common && \  
sudo apt-get install -y autotools-dev && \  
sudo apt-get install -y automake && \  
sudo apt-get install -y build-essential && \  
sudo apt-get install -y python3-dev && \  
sudo apt-get install -y python3-pip && \  
sudo apt-get install -y r-base r-base-core && \  
sudo apt-get install -y muscle && \  
sudo apt-get install -y default-jre && \  
sudo apt-get install -y ncbi-blast+ && \  
sudo apt-get install -y nano && \  

```

(continues on next page)

(continued from previous page)

```
sudo apt-get install -y less && \  
sudo apt-get install -y wget && \  
sudo apt-get install csh && \  
sudo apt-get install -y hammer && \  
sudo apt-get install -y libboost-all-dev && \  
sudo apt-get clean && \  
sudo rm -rf /var/lib/apt/lists/* /tmp/* /var/tmp/*
```

#Dssp installation If you are using Ubuntu 20.04, you can install dssp by the following command .. code:: bash

```
sudo apt-get install dssp
```

Otherwise, you can install it from the source by the following commands. Please note that default dssp in Ubuntu 22.04 is not working properly. .. code:: bash

```
wget https://github.com/cmbi/dssp/archive/refs/heads/master.zip && unzip -o master.zip && cd dssp-  
master/ && ./autogen.sh && ./configure && make && sudo make install && sudo ln -s  
/usr/local/bin/mkdssp /usr/local/bin/dssp && cd ../ && sudo rm -rf dssp-master/ && sudo rm -f master.zip
```

#HHSUITE installation .. code:: bash

```
wget https://github.com/soedinglab/hh-suite/releases/download/v3.3.0/hhsuite-3.3.0-AVX2-Linux.tar.gz  
&& mkdir hhsuite && mv hhsuite-3.3.0-AVX2-Linux.tar.gz hhsuite/ && cd hhsuite && tar xvfz hhsuite-  
3.3.0-AVX2-Linux.tar.gz && rm -f hhsuite-3.3.0-AVX2-Linux.tar.gz
```

#Add it to your path permanently inside .bashrc or .profile or .bash_profile Check the location of hhsuite folder and add it to your path In my case it was in /home/tekpinar/research/lcqb folder. Therefore, I added the following line to my .profile file. PATH=”/home/tekpinar/research/lcqb/hhsuite/bin:/home/tekpinar/research/lcqb/hhsuite/scripts:\$PATH”

Then source ~/.profile

#

```
cd ESGEMME
```

#Download ESGEMME from <http://gitlab.lcqb.upmc.fr/tekpinar/ESGEMME> repository and go inside the ESGEMME folder.! .. code:: bash

```
cd ESGEMME
```

4.3 Configuring default.conf file

Inside ESGEMME/esgemme folder, there is an important file called default.conf. This file contains essential parameters of ESGEMME, such as paths of external parts, default internal parameters. etc. You have to correct the Software section of this file according to your system.

```
pip3 install -e . &&\  
cd ../
```

#Installing the required R packages .. code:: bash

```
sudo Rscript -e 'install.packages(“seqinr”, repos=“http://cran.us.r-project.org”, dependencies=TRUE)'
```

#Installing secondary programs such as ev_couplings to obtain MSA files.

```
wget https://github.com/debbiemarkslab/plmc/archive/refs/heads/master.zip && \  
unzip -o master.zip && \  
cd plmc-master && \  
make all-openmp32 && \  
sudo cp bin/plmc /usr/local/bin/ && \  
cd ../ && \  
rm -rf master.zip plmc-master
```


INDICES AND TABLES

- `genindex`
- `modindex`
- `search`