# Nonsynonymous mutation and deletion analysis: A fast and accurate pipeline

Md. Shaminur Rahman

BS (Hon's) and MS, Department of Microbiology, University of Dhaka.

**1. Pyfasta** (https://github.com/brentp/pyfasta): For large dataset (more than 10,000 sequences), we need to differentiate the data into two or more dataset as MAFFT alignment tools has a maximum limit $\leq 10,000$.

**Usages (Linux):**

**2. MAFFT online alignment tools** (https://mafft.cbrc.jp/alignment/server/)**:** MAFFT is a very user friendly multiple alignment program for amino acid or nucleotide sequences. For SARS-CoV-2, specialize version (https://mafft.cbrc.jp/alignment/server/add_fragments.html?frommanual) has been launched where the maximum limit was 10,000. Here, existing alignment or reference sequences can be selected in a menu and other menu also available for new alignment. Here, in my experiment, I have used first as reference and the second one were for aligning sequence upload.

**3. Separate ORFs (S for SARS-CoV-2) from the alignment:** I have downloaded the complete S (for example) from reference (https://www.ncbi.nlm.nih.gov/nuccore/1798174254) from NCBI and open it in a text editor (Notepad/Notepad++). Open the alignment in MEGA 7 or other version. Copy the first few nucleotide (at least 20) from S and the go to MEGA. Press, ctrl F, then ctrl V and finally press enter. In this, we can find the first section of S Protein in the alignment. Click on the upstream of the of the S starting and press, shift home, this will select the upstream part of S protein, then delete the upstream part. This will remove the upstream of the S protein. For downstream removal do the same but press, shift end, then delete the downstream part. Insertional sequence must be checked and delete them. If it has importance, preserve the sequence for further analysis.
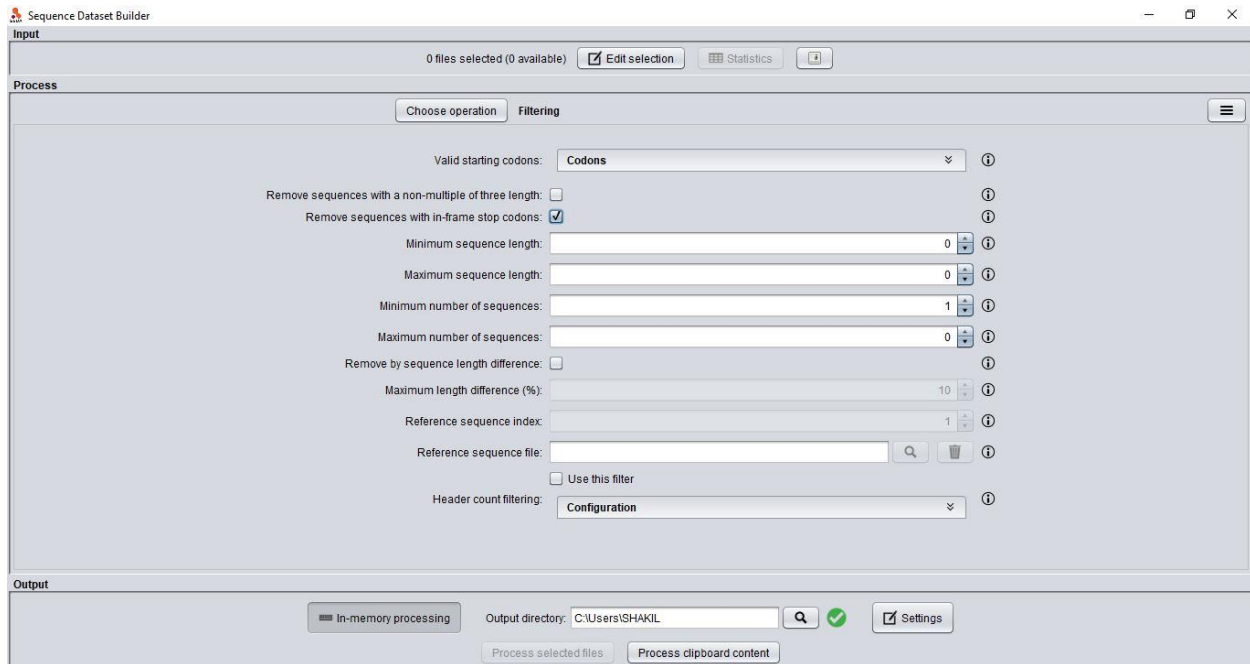


**4. Sequence cleaner** (https://github.com/metageni/Sequence-Cleaner)**:** Sequence cleaner will assist to remove all sequences which contain ambiguous ('M', 'D', 'R', 'N', 'K', 'Y', 'S', 'B', 'H', '-', 'V', 'W') characters. -ml 3822 means minimum length less than remove and -mn 0 means remove all N containing sequences.
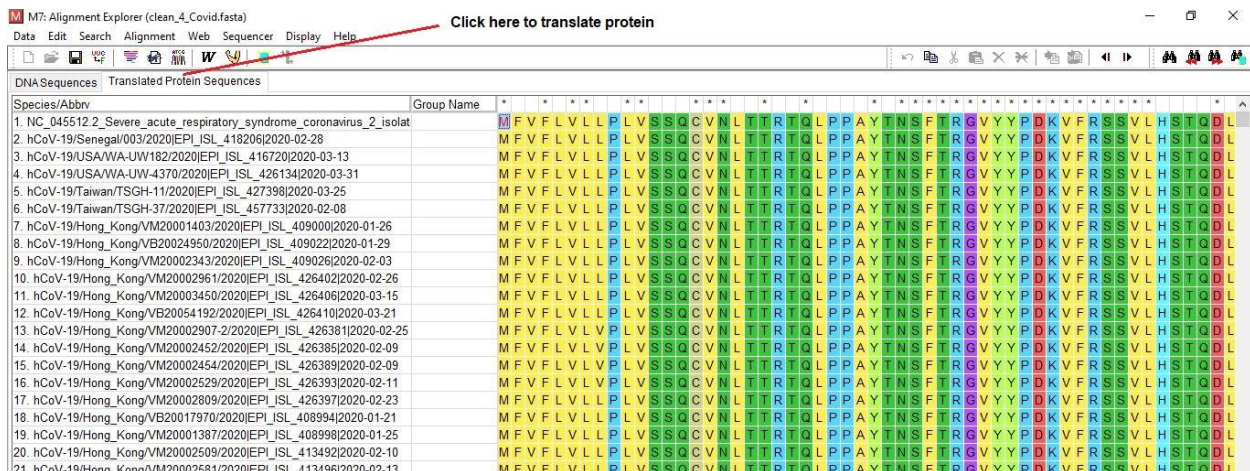
**Usages (Linux)**

**5. Remove internal stop codon using SEDA (**http://www.sing-group.org/seda/download.html**):** Internal stop codon must remove before mutation analysis and SEDA v1.1 (https://www.sing-group.org/seda/manual/operations.html#id5) will assist to remove this. Windows 10 defender takes the file as corrupted, so before using it just make sure off the real-time protection.



**6. DNA to Protein in MEGA** (https://www.megasoftware.net/dload_win_gui): After removal of all low quality and in frame stop codon sequences, we need to translate the DNA to protein. MEGA 7 or other version will aid us to do it. Just open the multiple sequence alignment and press the translate to protein options. Then export the protein alignment to fasta format.

**7. Pairwise mutation analysis** (https://github.com/SShaminur/Mutation-Analysis): Here the mutation results come from pairwise alignment with respect to reference genome (reference amino acid:position:strain amino acid). Before running, we must to ensure that there are no stop codon (*) in the last position of the ORF. We must delete the last stop codon in MEGA analysis.

*Multiple sequence alignment (Top one will be selected as reference)*

```
        1 2    4    6    8    10    13
A   [ 1 M D I I F W V S T F V L F
B   [ 1 M D I E F W F S T F V L F
C   [ 1 M D I I F W F S T W V L F
```

*Results (Reference A)*

```
B I4E
B V7F
C V7F
C F10W
```

**Usages (Linux):** (Requirements: Python ≥3.7, Biopython)




**8. Deletion analysis** (https://bioinf.shenwei.me/seqkit/usage/)**:** SeqKit tools used to arrest all the sequences containing gap (-). From there in frame deletion should be carefully find out. To ensure the in frame deletion with MEGA, remove the reference strain, then remove the gap of the strain/s and translate it to protein. Export the protein sequence/s in fasta format. Usually this translation finds the stop codon in last position. If not, then there are sequencing error (if virus not evolve too much). Then again align the protein sequence to the reference genome protein sequences. We can find the deleted amino acid position. Triplet codon deletion should be screen for deletion analysis. Then finally visualize the deletion data in a suitable software (Jalview, Unipro-UGENE, BioEdit etc.)

**Usages (Linux):**

*DNA deletion*



*Respective protein deletion*



## 9. Mutational analysis using Microsoft Excel:

I.   **Text to column:** Select the column, Then, **Data >Text to Columns > Delimited > Next > Space, Other (/) > Next and Finish.** This will differentiate country name Mutation and others. If we want to separate accession number, the just put Other (|).

**II.** **Flash Fill:** This will separate the mutation (ref:position:strain) into three different columns. From there, we can sort largest to smallest or vice versa. For Flash Fill, first we need to fill up at least two raw then, **Data >** select the 1ˢᵗ desired column > **Flash Fill.** For rest of the two columns, do the same thing.





**III.** **Remove Duplicates:** Unique mutation, Unique position mutation can be found by removing duplicates. Select the column (where duplicates need to remove) then, **Data > Remove Duplicates.**

**IV.** **Frequency count:** For frequency count, 1$^{st}$ column contains the original data and second column contains the duplicate removal data. Then select a cell in third column and then, =COUNTIFS(A1:A17,B1).



**10. Arrest sequence through sequence ID:**

**Usages (Linux):**