# Algorithms in Bioinformatics

## BLAST (Basic Local Alignment Search Tool) for COVID 19 Genome Strains

### Group Members

Anmol Kumar (2018382)

Prutyay Gautam (2018403)

Sanskar Sachdeva (2018411)
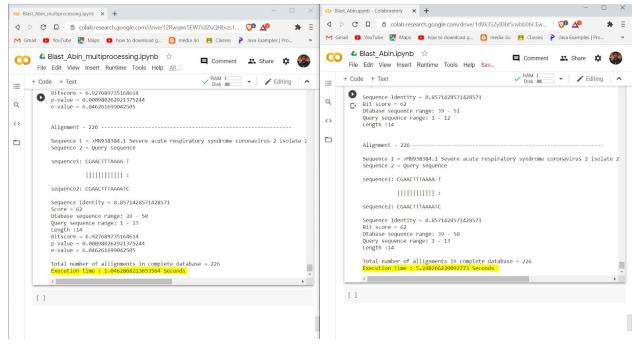
Sarthak Pal (2018412)

### Program execution instructions

To run the following program, execute it with the following command:

```
python3 file_name.py
```

You are prompted to input the following options :
- **Query Sequence:** String input
- **Word Length:** Integer Input
- **HSSP Score Threshold** : Integer Input
- **Extension Threshold**: Integer Input

## Optimizations

- Optimization was done on improving the computation time, by applying multiprocessing. 5 times less time was used by the program using multiprocessing.



- The Smith-Waterman algorithm was applied on a strict sub range of database and query which was developed by identifying the adjacent words. The algorithm was not extended outside of these ranges as they would have led to less score or highly gapped local alignment. Therefore computation time was saved by not calculating such low scoring local alignments. If in case user wanted to get these alignments they can get it by applying a smaller HSSP Thershold.

## References

- Blast: The Algorithm: https://www.youtube.com/watch?v=wfi_KimrNQM
- Heuristics of Blast: https://www.youtube.com/watch?v=jzSIC2UzxZ4
- Basic Local Alignment Search tool: https://www.sciencedirect.com/science/article/pii/S0022283605803602