

SARS-CoV-2 Protein Analysis

Author: Sushmit Dutta

The virus 'Severe Acute Respiratory Syndrome Coronavirus 2' (SARS-CoV-2) is defined as the agent that causes the novel coronavirus disease which started spreading in

1. The virus was first reported in China when there was a rise in cases of pneumonia around the country. The spread of the disease around the world has affected the livelihood of billions of people. In March of 2020, the World Health Organisation declared the disease to be a pandemic.

Around the world, 120 million people have contracted the disease and out of that 2.66 million have died. Stopping this virus from spreading is the goal of all world governments, doctors, and scientists. There have been a few vaccines that have helped flatten the curve and save countless lives. However, as time passes, newer and stronger evolutions of the strand are spreading around the globe. The key to stopping another wave of the virus from spreading is scientists making vaccines as effective as possible.

The key method to accomplish this goal is to analyse the proteins in the genetic material of the virus. By altering the proteins or inhibiting them from attaching to other cells, scientists can limit the spread. This notebook goes through the process of finding those strands of proteins.

Key Terms

1. Nucleotides - They serve as the monomeric units of DNA and RNA. They are organic molecules consisting of a nucleoside and a phosphate. There are four possible nucleotides - A, T, G, C. A and T are analogous, and G and C are analogous.
2. Codons - A sequences of three nucleotides that together form a unit of genetic code in DNA and RNA. There are two special types of codons - Start and Stop Codons. These sequences dictate where a sequence starts and ends.
3. Open Reading Frames (ORF) - In molecular genetics, an open reading frame is the part of a reading frame that has the ability to be translated. An ORF is a continues stretch of codons that begins with a start codon and ends with a stop codon.

Function Breakdown

The code below is the output of the find_genes function which aims to find all the proteins in the novel coronavirus strain. The function below runs with 1500 shuffles in the noncoding_orf_threshold function.

In [1]:

```
%load_ext autoreload
%autoreload 2

from gene_finder import *

path = "data/NC_045512.2.fa"
find_genes(path)
```

Out[1]:

```
[ 'MVPHISRQRLTKYTMADLVYALRHFDENCDTLKEILVTYNCCDDDYFNKKDWYDFVENPDILRVYANLGERVRQ
ALLKTVQFCDAMRNAGIVGVLTLDNQDLNGNWDYDFGDFIQTTPGSGVPVVDSSYSSLMPILTLTRALTAESHVDTDL
TKPYIKWDLKLYDFTEERLKLFDYFYWDQTYHPNCVNCLEDDRCILHCANFNVLSTVFPPTSFGPLVRKIFVDGV
PFVSTGYHFRELGVVHNQDVNLHSSRLSFKELLVYAADPAMHAASGNLLLDKRTTCFSVAALTNNVAFQTVKPGNF
NKDFYDFAVSKGFFKEGSSVELKHFFFAQDGNAAISDYDYRYNLPTMCDIRQLLFVVEVVDKYFDCYDGGCINANQ
VIVNNLDKSAGFPFNKGKARLYYDSMSYEDQDALFAYTKRNVIPITITQMNLYAISAKNRARTVAGVSICSTMTNR
QFHQKLLKSI AATRGATVVIGTSKFYGGWHNMLKTVYSDVENPHLMGWDYPKCDRAMPNMLRIMASLVLARKHTTCC
SLSHRFYRLANECQVLSEVMVMCGGSLYVKPGGTSSGDATTAYANSVFNICQAVTANVNALLSTDGKNIADKYVRNL
QHRLYECLYRNRDVTDFVNEFYAYLRKHFSMMILSDDAVVCFNSTYASQGLVASIKNFKSVLYYQNNVFMSEAKCW
TETDLTKGPHEFCSQHTMLVKQGGDYVYLPYDPSPRILGAGCFVDDIVKTDGTLMIERFVSLAIDAYPLTKHPNQEY
ADVHLYLQYIRKLHDELTHGMLDMYSVMLTNDNTSRYWEPEFYEAMYPHTVLQAVGACVLCNSQTSRLRCGACIRR
PFLCCKCCYDHVISTSHKLVL SVNYPVCNAPGCDVTDVTQLYLGMSYICKSHKPPISFPLCANGQVFGLYKNTCVG
SDNVTDFNAIATCDWTNAGDYILANTCTERLKLFAAETLKATEETFKLSYGIATVREVLSDRELHLSWEVGKPRPPL
NRNYVFTGYRVTKNSKVQIGEYTFEKG DYGDVAVYRGTTTTYKLVNGDYFVLTSHTVMPLSAPTLPQEHYVRITGLY
PTLNISDEFSSNVANYQKVG MQYSTLQGPPTGKSHFAIGLALYPSARIVYTACSHA AVDALCEKALKYLPIDKC
SRIIPARARVECFDKFKVNSTLEQYVFCTVNALPETTADIVVFDEISMATNYDLSV VNARLRAKHVYIGDPAQLPA
PRTLLTKGTLEPEYFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTV SALVYDNKLKAHKDKSAQC FKMFKGVITHDV
SSAINRPQIGVVREFLTRNPAWRKAVFISPYNQNAVASKILGLPTQTV DSSQGSEYDYVIFTQTTETAHSCNVNRF
NVAITRAKVGILCIMSDDRDLYDKLQFTSLEIPRRNVATLQAE NVTGLFKDCSKVITGLHPTQAPTHLSVDTKFKTEG
LCVDIPGIPKDMTYRRLISMMGFKMNYQVNGYPNMFITREEAIRHVR AWIGFDVEGCHATREAVGTNLPLQLGFSTG
VNLVAVPTGYVDTPNNTDFSRVSAKPPPGDQFKHLIPLMYKGLPWNVVR IKIVQMLSDTLKNLSDRVVFVLWAHGFE
LTSMKYFVKIGPERTCCLCDRRATCFSTASDTYACWHHSIGFDYVYNPF MIDVQQWGFTGNLQSNHDLYCQVHGNAH
VASCDAIMTRCLAVHECFVKRVDWTIEYPIIGDELKINAACRKVQHMV VKAALLADKFPVLHDIGNPKAIKCPQAD
VEWKFYDAQPCSDKAYKIEELFYSYATHSDKFTDGVCLFWNCNVD RY PANSIVCRFDTRVLSNLNLP GCDGGSLYVN
KHAFHTPAFDKSAFVNKQLPFFYYSDSPCESHKGQVVS DIDIYVPLKSATCITRCNLGGAVCRHHANEYRLYLDAYN
MMISAGFSLWVYKQFDTYNLWNTFTRLQSL ENVAFNVVNKGHFDGQQGEVPVSIINNTVYTKVDGVDVELFENK TTL
PVNVAFELWAKRNIKPVEVKILNNLGVDIAANTVIWDYKRDA PAHISTIGVCSMTDIAKKPTETICAPLTVFFDGR
VDGQVDLFRNARNGVLITEGSVKGLQPSVGPQKQASLNGVTLIGE AVKTQFNYYKKVDGVVQQLPETYFTQSRNLQEF
KPRSQMEIDFLELAMDEFIERYKLEGYAFEHIVYGDFS HSQLGGLHLLIGLAKRFKESPFEELED FIPMDSTVKNYFI
TDAQTGSSKCVCSVIDL LLLDDFVEI IKSQDLSVSVKVKVTIDY TEISFMLWCKDGHVETFY PKLQSSQAWQPGVAM
PNLYKMQRMLLEKCDLQNYGDSATLPKGIMMNVAKYTQLCQY LNTLT LAVPYNMRV IHFGAGSDKG VAPGTAVLRQW
LPTGTLLVDSDLNDFVSDADSTLIGDCATVHTANKWDLI ISDMYDPKTKNVT KENDSKEGFFTYICGFIQQKLALGG
SVAIKITEHSWNADLYKLMGHFAWWTAFVTNVNASSSEAF LIGCN YLGKPREQIDGYVMHANYIFWRNTNPIQLSSY
SLFDMSKFPLKLRGTAVMSLKEGQINDMILSLLSKGR LIIRENNRVVISSDVLVNN',
'MDLFMRIFTIGTVTLKQGEIKDATPSDFVRATATIPIQASLPFGWLIVGVALLAVFQSASKIITLKKRWQLALSK
GVHFVCNLLLLFVTVYSHLLLVAAGLEAPFLYLYALVYFLQS INFVRIIMRLWLCWKCRSKNP LLYDANYFLCWHTN
CYDYCIPYNSVTSSIVITSGDGTTSPISEHDYQIGGYTEKWESGVKDCVVLHSYFTSDYYQLYSTQLSTDTGVEHVT
FFIYNKIVDEPEEHVQIHTIDGSSGVVNPVMEPIYDEPTTTTSVPL',
'MYSFVSEETGTLIVNSVLLFLAFVVFLLVTLAILTALRLCAYCCNIVN VSLVKPSFYVYSRVKNLNSSRVPDLLV
',
'MKIILFLALITLATCELYHYQE CVRGTTVLLKEPCSSGTYEGNSPFHPLADNKFALT CFSTQFAFACPDGVKH
VYQLRARSVSPKLFIRQEEVQELYSPIFLIVAAIVFITLCFTLKRKTE',
```

'MESLVPGFNEKTHVQLSLPVLQVRDVLVRGFGDSVEEVLSEARQHLKDGTCLVEVEKEGVLPOLEQPYVFIK
RSDARTAPHGHVMVELVAELEGIQYGRSGETLGLVPHVGEIPVAYRKVLLRKNGNKGAGGHSYGADLKSFDLG
DELGTDPYEDFQENWNTKHSSGVTRELMRELNGGAYTRYVDNNFCGPDGYPLECIKDLLARAGKASCTLSEQLD
FIDTKRGVYCCREHEHEIAWYTERSEKSYELQTPFEIKLAKKFDTFNGECPNFVFPPLNSIIKTIQPRVEKKKLD
GFMGRIRSVYPVASPNECNQMCLSTLMKCDHCGETSWQTGDFVKATCEFCGTENLTKEGATTCGYLPQNAVVKI
YCPACHNSEVGPESHSLAEYHNESGLKTI LRKGGRTIAFGGCVFSYVGCHNKAYWVPRASANIGCNHTGVVGE
SEGLNDNLLEILQKEKVNINIVGDFKLNEEIAIILASFSASTSAFVETVKGLDYKAFKQIVESCGNFKVTKGKA
KKGAWNIGE QKSILSPLYAFASEAARVVRISFSRTLETAQNSVRVLQKAAITILDGISQYSLRLIDAMMFTSDL
ATNNLVVMAYITGGVVQLTSQWLTNIFGTVYEKLPVLDWLEEKFKEGVEFLRDGWEIVKFISTCACEIVGGQI
VTCAKEIKESVQTFFKLVNKFLALCADSIIIGGAKLKALNLGETFVTHSKGLYRKCVKSREETGLLMPLKAPKE
IIFLEGETLPTEVLTEEVVLKTGDLQPLEQPTSEAVEAPLVGTPVCINGLMLLEIKDTEKYCALAPNMMVTNNT
FTLKGGA PTKVTFGDDTVIEVQGYKSVNITFELDERIDKVLNEKCSAYTVELGTEVNEFACVVADAVIKTLQPV
SELLTPLGIDLDEWSMATYYLFDSEGEFKLASHMYCSFYPPDEDEEEGDCEEEEFEPSTQY EYGTEDDYQ GKPL
EFGATS AALQPEEEQEEDWLDDDSQQTVGQQDGS EDNQT TTIQTIVEVQPQLEMELTPVVQTI EVNSFSGYLKL
TDNVYIKNADIVEEAKVKPTVVVNAANVYLKHGGGVAGALNKATNNAMQVESDDYIATNGPLKVGGS CVLSGH
NLAKHCLHVVGPNVNKGEDIQLLKSAYENFNQHEVLLAPLLSAGIFGADPIHSLRVCVDTVRTNVYLAVFDKNL
YDKLVSS FLEMKSEKQVEQKIAEIPKEEVKPFITESKPSVEQRKQDDKKIKACVEEVTTTLEETKFLTENLLLY
IDINGNLHPDSATLVSDIDITFLKKDAPYIVGDVVQEGVLTAVVIPTKKAGGTTEMLAKALRKVPTDNYITTY
P
GQGLNGYTVEEAKTVLKKCKSAFYILPSIISNEKQEILGTVSWNLREMLAHAEETRKLMPVCVETKAIVSTIQ
R
KYKGIKIQEGVVDYGARFYFYTSKTTVASLINTLNDLNETLVTMPLGYVTHGLNLEEAARYMRS LKVPATVSVS
SPDAVTAYNGYLTSSSKTPEEHFIETISLAGSYKDWYSYGQSTQLGIEFLKRGDKSVYYTSNPTTFHLDGEVIT
FDNLKTL LSLREVRTIKVFTTVDNINLHTQVVDMSMTY GQQFGPTYLDGADVTKIKPHNSHEGKTFYVLPNDDT
LRVEAFEYYHTTDP SFLGRYMSALNHTKKWKYPQVNGLT SIKWADNNCYLATA LLTLQQIELKFNP PALQDAY
Y
RARAGEAANFCALILAYCNKTVGELGDVRETMSYLFQHANLDSCKRVLNVVCKTCGQQQTTLKGVEAVMYMGTL
SYEQFKKG VQIPCTCGKQATKYL VQQESPFVMM SAPP AQYELKHGTFTCASEYTGNYQC GHYKHITSKETLYCI
DGALLTKSSEYKGPITDV FYKENSYTTTIKPV TYKLDGVVCTEIDPKLDNYYKKD NSYFTEQPIDLVPNQYPN
ASFDNFKFVCDNIKFADDNLQLTGYKKPASRELKVTFFPDLNGDVVAIDYKHYTPSFKKGAKLLHKPIVWHVNN
ATNKATYKPN TW CIRCLWSTKPVETSNSFDVLKSEDAQGM DN LACEDLKPVSEEVVENPTIQKDVLECNVKTTE
VVGDIILK PANNSLKITEEVGHTDLMAAYVDNSSLTIKKPNELSRVLGLKTLATHGLAAVNSVPWDTIANYAKP
FLNKVVSTTTNIVTRCLNRVCTNYMPYFF TLLLQLCTFTRSTNSRIKASMPPTIAKNTVKS V GKFCLEASFNYL
KSPNFSKLINIIWFLLLSVCLGSLIYSTAALGV LMSNLGMPSYCTGYREGYLNSTNVTIATYCTGSIPCSVCL
SGLDSLDTYPSLETIQITISSFKWDLTAFGLVAEWFLAYILFTRFFYVLGLAAIMQLFFSYFAVHFISNSWLMW
LIINLVQMAPISAMVRMYIFFASFYVWKS YVHVVDGCNSSTCM MCYKRNRATRVECTTIVNGVRRSFYVYANG
GKGFC LKHNWNCVNCDTFCAGSTFISDEVARDLSLQFKRPINPTDQSSYIVDSVTVKNGSIHLYFDKAGQKTYE
RHSLSHFVNLDNLRANNTKGSLP INVIVFDGKSKCEESSAKSASVYYSQ LMCQPILLDDQALVSDVGD SAEVAV
KMFDAYVNTFSSTFNVPMEK LKTLVATAEAE LAKNVSLDNVLSTFISAARQGFVDS DVETKDVVECLKL SHQSD
IEVTGDSCNNYMLTYNKVENMTPRDLGACIDCSARHINAQVAKSHNIALIWNVKDFMSLSEQLRKQIRSAAKN
NLPFKLTCATTRQVVNVVTTKIALKGGKIVNNWLKQLIKVTLVFLFVA AIFYLITPVHVM SKHTDFSSEIIGYK
AIDGGVTRDIASDTCTFANKHADFDTWFSQRGGSYTN DKACPLIAAVITREVG FVVPGLPGTILRTTNGDFLHF
LPRVFSAVGNICYTPSKLIEYTD FATSACVLAAECTIFKDASGKVPYCYDTNVLEGSVAYESLRPDTRYV LMD
GSIIQFPNTYLEGSRVVTTFDSEYCRHGTCERSEAGVCVSTSGRWVLNNDYRSLPGVFCGVDAVNLLTNMFT
PLIQPIGALDISASIVAGGIVAI VVTC LAYYFMRFRRAFGEYSHVVAFNTLLFLMSFTVLC LTPVYSFLPGVYS
VIYLYLTFYLTNDVSFLAHIQWMVMFTPLVPFWITIA YIIICISTKH FYWFFSNY LKRRVVFNGVSFSTFEEAAL
CTFLLNKEMYLKLRSDVLLPLTQYNRYLALYNKYKYFS GAMDTTSYREAA CCHLAKALNDFSNSGSDVLYQPPQ
TSITSAVLQSGFRKMAFP SGKVEGCMVQVTCGTTTLNGLWLDDVVYCPRHVICTSEDMLNPNYEDLLIRKSNHN
FLVQAGNVQLRVIGHSMQNCVLKLKVDTANPKTPKYKFVRIQPGQTF SVLACYNGSPSGVYQCAMRPNFTIKGS
FLNGSCGSVGFNIDYDCVSFCYMHMELPTGVHAGTDLEGNFYGPFVDRQTAQAAGTDTTITVNVLAWLYAAVI
NGDRWFLNRFTTTLNDFNLVAMKYNYEPLTQDHVDILGPLSAQTGIAVLDMCASLKELLQNGMNGRTILGSALL
EDEFTPFDVVRQCSGVTFQSAVKRTIKGTHH WLLL TILTSLLVLVQSTQWSLFFFLYENAF LFPFAMGIIAMSAF
AMMFVKHKHAF LCLFLLPSLATVAYFNMVMPASWVMRIMTWLDMVDTSLSGFKLKDCVMYASAVVLLIIMTAR
TVYDDGARRVW TLMNVLT LVYK VYYGNALDQAI SMWALIISVTSNYSGVVTTVMFLARGIVFMCVEYCPIFFIT
GNTLQCIMLVYCF LGYFCTCYFGLFCLLNRYFRLTLGVYDYL VSTQEF RYMNSQGLLPKNSIDAFKLN ILLG
VGGKPCIKVATVQSKMSDVKCTSVVLLSVLQQLRVESSSKLWAQCVQLHNDILLAKDTTEAFEKMSVLLSVLLS
MQGAVDINKLCEEMLDNRATLQAIASEFSS LPSYAAFATAQEAYEQAVANGDSEVVLK LKLSLVAKSEFDRD

AAMQQRKLEKMADQAMTQMYKQARSEDKRAKVTSAMQTMLFTMLRKLNDNDALNNIINNARDGCVPLNIIPLTTAA
 KLMVVI PDYNTYKNTCDGTTFTYASALWEIQQVVDADSKIVQLSEISMDNSPNLAWPLIVTALRANSVAVKLQNN
 ELSPVALRQMSCAAGTTQTACTDDNALAYYNTTKGGRFVLALLSDLQDLKWARFPKSDGTGTIYTELEPPCRFV
 TDTPKGPVKYLYFIKGLNNLNRMVGLGSLAATVRLQAGNATEVPANSTVLSFCAFAVDAAKAYKDYLASGGQP
 ITNCVKMLCTHTGTGQAITVTPEANMDQESFGGASCCLYCRCHIDHPNPKGFCDLKGKYVQIPTTCANDPVGFT
 LKNTVCTVCGMWKGYGCSCDQLREPMLQSADAQSFLNGFAV',
 'MLNQVEPHQEMPQLMLLIVFLTFVKLSRPMLMHFYLLMVTKLPISMSAIYNTDFMSVSIEIEMLTQTL',

Interpretation

As you can see above the strains vary in length. The table below provides a breakdown of some of the amino acid chains found above:

Protein	Accession	Query Length
ORF1a	YP_009725 295.1	4405
Nucleocapsid Phosphoprotein	YP_009724 397.2	419
Envelope Protein	YP_009724 392.1	75
Membrane Glycoprotein	YP_009724 393.1	222
Surface Glycoprotein	BCN86353.1	1282

These are not the list of all the strains found within the FATSA file passed but they are some of the important ones that scientists can aim to work around. By experimenting with these prominent proteins, the vaccines can aim to alter the configuration of the polymers to slow them down or make them less harmful.