

 Open access • Posted Content • DOI:10.1101/2020.04.24.050534

## **ACE2 Homo-dimerization, Human Genomic variants and Interaction of Host Proteins Explain High Population Specific Differences in Outcomes of COVID19**

— [Source link](#) 

Swarkar Sharma, Inderpal Singh, Shazia Haider, Md. Zubbair Malik ...+2 more authors

**Institutions:** Shri Mata Vaishno Devi University, Jaypee Institute of Information Technology, Jawaharlal Nehru University

**Published on:** 24 Apr 2020 - bioRxiv (Cold Spring Harbor Laboratory)

**Topics:** Virus and Expression quantitative trait loci

Related papers:

- [A pneumonia outbreak associated with a new coronavirus of probable bat origin](#)
- [Comparative genetic analysis of the novel coronavirus \(2019-nCoV/SARS-CoV-2\) receptor ACE2 in different populations](#)
- [SARS-CoV-2 Cell Entry Depends on ACE2 and TMPRSS2 and Is Blocked by a Clinically Proven Protease Inhibitor](#)
- [Global Spread of SARS-CoV-2 Subtype with Spike Protein Mutation D614G is Shaped by Human Genomic Variations that Regulate Expression of TMPRSS2 and MX1 Genes](#)
- [A Novel Coronavirus from Patients with Pneumonia in China, 2019.](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/ace2-homo-dimerization-human-genomic-variants-and-4jpu3reefr>

## ACE2 Homo-dimerization, Human Genomic variants and Interaction of Host Proteins Explain High Population Specific Differences in Outcomes of COVID19

Swarkar Sharma<sup>1\*</sup>, Inderpal Singh<sup>1,2†</sup>, Shazia Haider<sup>3†</sup>, Md. Zubair Malik<sup>4†</sup>, Kalaiarasan Ponnusamy<sup>5†</sup>, Ekta Rai<sup>1</sup>

1. Human Genetics Research Group, School of Biotechnology, Shri Mata Vaishno Devi University, Katra, Jammu and Kashmir, India.
2. Bioinfores Pvt. Ltd., R. S. Pura, Jammu, Jammu and Kashmir, India.
3. Department of Biotechnology, Jaypee Institute of Information Technology, Noida, Sector-62, Uttar Pradesh, India
4. School of Computational and Integrative Sciences, Jawaharlal Nehru University, New Delhi, India.
5. School of Biotechnology, Jawaharlal Nehru University, New Delhi, India.

† - All authors contributed Equally

### Corresponding Author:

Dr. Swarkar Sharma

Coordinator, Human Genetics Research Group,

School of Biotechnology, Shri Mata Vaishno Devi University, Katra, J&K, India

Email: [swarkar.sharma@smvdu.ac.in](mailto:swarkar.sharma@smvdu.ac.in)

Mobile: +91-9419955636; Ph: +91-1991-285535//285525 Ext. 2385

### ABSTRACT

Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) is a positive single-stranded RNA virus that causes a highly contagious Corona Virus Disease (COVID19). Entry of SARS-CoV-2 in human cells depends on binding of the viral spike (S) proteins to cellular receptor Angiotensin-converting enzyme 2 (ACE2) and on S-protein priming by host cell serine protease TMPRSS2. Recently, COVID19 has been declared pandemic by World Health Organization (WHO) yet high differences in disease outcomes across countries have been seen. We provide evidences to explain these population-level differences. One of the key factors of entry of the virus in host cells presumably is because of differential interaction of viral proteins with host cell proteins due to different genetic backgrounds. Based on our findings, we conclude that a higher expression of *ACE2* is facilitated by natural variations, acting as Expression quantitative trait loci (eQTLs), with different frequencies in different populations. We suggest that high expression of ACE2 results in homo-dimerization, proving disadvantageous for TMPRSS2 mediated cleavage of ACE2; whereas, the monomeric ACE2 has higher preferential binding with SARS-CoV-2 S-Protein vis-a-vis its dimerized counterpart. Further, eQTLs in *TMPRSS2* and natural structural variations in the gene may also result in differential outcomes towards priming of viral S-protein, a critical step for entry of the Virus in host cells. In addition, we suggest that several key host genes, like *SLC6A19*, *ADAM17*, *RPS6*, *HNRNPA1*, *SUMO1*, *NACA*, *BTF3* and some other proteases as Cathepsins, might have a critical role. To conclude, understanding population specific differences in these genes may help in developing appropriate management strategies for COVID19 with better therapeutic interventions.

**Also read at:** <https://science.sciencemag.org/content/367/6485/1444/tab-e-letters>

**RE: ACE2 Homodimerization Affects Binding of SARS-CoV-2 Spike Protein**

## Introduction

The recent emergence of corona virus disease (COVID19) caused by SARS-CoV-2 has resulted in >4.1 Million infections and >285 thousands deaths worldwide so far, and the numbers are increasing exponentially [<https://covid19.who.int/>]. SARS-CoV-2 is reported to be originated in bats (Zhou et al., 2020) and transmitted to humans via unknown intermediate host. However, its origin is still being questioned time and again. With the declaration of SARS-CoV-2 as pandemic by WHO, extensive research worldwide has been carried out. It has been established that Human ACE2 mediates SARS-CoV-2 entry into cells through its Spike (S) Protein, which primarily makes entry to host body through respiratory tract with nasal epithelial cells as potential initial infection site (Sungnak et al., 2020). ACE2 is a functional receptor on human cells for this newly originated coronavirus (Walls et al., 2020) with a higher affinity than the severe acute respiratory syndrome coronavirus (SARS-CoV) originated in 2002 (Wan et al., 2020). However, no substantial evidence exists about the higher expression of ACE2 being primarily associated with the degree of infection (Kuster et al., 2020). Additionally, COVID19 lethality is mostly driven by the extent of underlying lung injury; whereas, a negative correlation has been reported between ACE2 expression and lung injury (Imai et al., 2005). A recent report also suggests an inhibition of SARS-CoV-2 by human recombinant soluble (hrs) ACE2 (Monteil et al., 2020), making it an interesting question to explore.

High differences in clinical outcomes across countries have been [<https://covid19.who.int/>] noted which demonstrate that neither all people who are exposed to SARS-CoV-2 develop infection nor all infected patients end up in severe respiratory illness (Guan et al., 2020), which cannot be explained by immunity alone (Shi et al., 2020). This leads one to hypothesize about differential genetic susceptibility to COVID19 and virulence of SARS-CoV-2 in different populations (Kaiser, 2020). Efforts are being made to gain a better understanding of this disease and even propose prophylactic-hypothetical role of Bacillus Calmette–Guérin (BCG) vaccine, a vaccine primarily used against tuberculosis (TB), for reduced morbidity and mortality for COVID19 in human population (Miller et al., 2020). Due to extensive contemporary research, evidences in favour (Stawiski et al., 2020) as well as against (Cao et al., 2020) the existence of SARS-CoV-2 S-protein binding-resistant ACE2 natural variants, in different populations, have been found recently. In addition, studies highlighting role of eQTLs in *ACE2* expression, resulting in potential differential COVID19 fatality (Cao et al., 2020; Chen et al., 2020a), are pouring in. Where most such studies are targeting only natural variations in *ACE2* gene as SARS-CoV-2 differential susceptibility factor, recent evidences suggest that additional host proteins like cellular serine protease TMPRSS2 act as co-factors and are critical for efficient cellular infection by SARS-CoV-2 (Hoffmann et al., 2020). Further, it is equally important to consider that rare functional variants, with uncertain consequences, may not explain large-scale population level differential clinical outcomes.

This highlights the importance of identification of other potential co-factors and underlying mechanisms these genes could be involved in. At the same time, understanding the interactions of these host proteins with SARS-CoV-2 along with ACE2 may explain many of the unanswered questions. Further evaluation of these host genes, and exploring their natural occurring variants, along with their expression patterns, may also shed some light on better conceptual framework of differential susceptibility to COVID19 and the virulence of SARS-CoV-2. In the present study, we have tried to exploit existing literature to understand mechanisms of correlations, of several relevant host genes, including ACE2, which interact specifically with some prominent viral proteins, using *in-silico* approaches to understand the role of these as factors responsible for population level differences.

## **Methodology**

As SARS-CoV-2 is primarily a respiratory pathogen, the present study was conducted to understand mechanisms of its entry to cells in lungs, and we believe it could be extrapolated to other respiratory tract tissues, to explore potential factors that are responsible for differential outcomes in respiratory illness. To begin with, we started with the most studied host protein ACE2 and tried to understand its interaction with SARS-CoV-2 S-Protein followed by adding other interacting proteins. We further added layers of other methods to have a better understanding of the potential causes and outcomes.

## **Molecular Dynamics (MD) Simulations**

Recently submitted experimental structure (6M17.pdb) of ACE2 dimer bound to two B0AT1 (coded by *SLC6A19*) and two receptor-binding domains (RBD) of S1 of SARS-CoV-2 Spike protein (S-protein), with overall resolution of 2.90 Angstrom and SARS-CoV2 Spike Glycoprotein with RBD domain in Up conformation (6VSB.pdb) were downloaded from the Protein Data Bank for structural visualization and dynamics analysis (Wrapp et al., 2020; Yan and Zhang, 2020). This structure was corrected for missing amino acids, side chains, missing hydrogen atoms, disulphide bonds etc. It was also optimized for pKa corresponding to physiological pH 7.2 and energy minimized to correct steric clashes using the Protein Preparation Wizard of Schrodinger Maestro (Madhavi Sastry et al., 2013). Predefined solvation model TIP3P was used and overall neutrality of the system was maintained by addition of Na<sup>+</sup> and Cl<sup>-</sup> counter ions (Mark and Nilsson, 2001). Physiological salt concentration of 0.15 Molar was generated through addition of NaCl. Periodic boundary condition of 10 Angstrom was set using the System Builder Tool of Desmond software(2006). Total two MD simulations, each 10 nanosecond long, were conducted. In one of the simulations, six chains i.e. ACE2 dimer with each monomer bound to a B0AT1 and RBD domain was simulated, where the resulting solvated system consisted of 424,847 atoms. For the other simulation, monomeric ACE2 bound to viral RBD consisted of 174,900 atoms. Root mean square fluctuation (RMSF) and principal component analysis were done through R based BIO3D module (Grant et al., 2006). MMGBSA free energy of binding between proteins was

calculated through Prime Software of Schrodinger Suite (Jacobson et al., 2002). Structural visualizations and images were traced using pyMOL and VMD (Jacobson et al., 2002) (Humphrey et al., 1996).

### **Gene Expression in Alternate Transcripts and Regulation Analyses**

Genotype-Tissue Expression (GTEx) portal at <https://gtexportal.org/> was used to explore various gene expression profiles. By Tissue, Multigene Query as well as transcript browser was used to understand expression of splice variants and exons in Lungs. eQTLs were viewed by GTEx IGV Browser as well as GTEx Locus Browser was used to plot gene specific eQTLs. eQTLs and other gene regulatory information, splice variants, and ESTs were also explored through UCSC genome browser, <https://genome.ucsc.edu/> with build GRCh37/hg19. Table browser was used to interact with various datasets to look for overlaps and filter outcomes. Retrieved information was saved in files as well as plotted on UCSC browser as custom tracks. Comparison of RNA expression and Protein expression, mainly in lungs, was done at the Human Protein Atlas (Uhlen et al., 2015) at [www.proteinatlas.org](http://www.proteinatlas.org). HaploReg v4.1 (<https://pubs.broadinstitute.org/mammals/haploreg/haploreg.php>) was used to annotate effect of noncoding genome variants of regulation motifs.

### **Analyses of Allele frequency distribution of Variants in different population Groups**

Genetic data for global population groups was explored through the GnomAD portal at <https://gnomad.broadinstitute.org/>, as well as dbSNP database of NCBI at <https://www.ncbi.nlm.nih.gov/snp>. Gene specific SNPs data was also retrieved from 1000 genomes (1000G) Phase 3 data set available through <https://www.ncbi.nlm.nih.gov/variation/tools/1000genomes> as per Homo sapiens:GRCh37.p13 (GCF\_000001405.25). Data for the population groups belonging to five super population groups [African (AFR), Ad mixed American (AMR), East Asian (EAS), European (EUR) and South Asian (SAS)] was analysed. The following sub population groups were studied. AFR: African Caribbeans in Barbados (ACB), Americans of African Ancestry in SW USA (ASW), Esan in Nigeria (ESN), Gambian in Western Divisions in the Gambia (GWD), Luhya in Webeye, Kenya (LWK), Mende in Sierra Leone (MSL), Yoruba in Ibadan, Nigeria (YRI). AMR: Colombians from Medellin, Colombia (CLM), Mexican Ancestry from Los Angeles USA (MXL), Peruvians from Lima, Peru (PEL), Puerto Ricans from Puerto Rico (PUR). EAS: Chinese Dai in Xishuangbanna, China (CDX), Han Chinese in Beijing, China (CHB), Southern Han Chinese (CHS), Japanese in Tokyo, Japan (JPT), Kinh in Ho Chi Minh City, Vietnam (KHV). EUR: Utah Residents (CEPH) with Northern and Western European Ancestry (CEU), Finnish in Finland (FIN), British in England and Scotland (GBR), Iberian Population in Spain (IBS), Toscani in Italia (TSI). SAS: Bengali from Bangladesh (BEB), Gujarati Indian from Houston, Texas (GIH), Indian Telugu from the UK (ITU), Punjabi from Lahore, Pakistan (PJI) and Sri Lankan Tamil from the UK (STU).

LDproxy tools of LDlink 4.0.3 web-based suite (<https://ldlink.nci.nih.gov/>) was used to map proxy variants in strong LD and with putatively functional role. LDpop tool of the LDlink 4.0.3 was used for geographically annotating allele frequencies in 1000G populations ; alternatively also done by a web-based tool Data wrapper (<https://www.datawrapper.de/>) in some of the figures.

### **Computational structural analysis on SARS-CoV-2 S-protein and TMPRSS2**

TMPRSS2 has been recently characterized as a critical component for cell entry by SARS-CoV-2 (Hoffmann et al., 2020). To understand its interactions, the protein sequence of surface glycoprotein (YP\_009724390.1) of SARS CoV-2 and two transcripts of TMPRSS2 (NP\_005647.3 and NP\_001128571.1) protein of homo sapiens, were retrieved from the NCBI-Protein database. Pairwise sequence alignment of TMPRSS2 isoforms was carried out by Clustal Omega tool (Sievers et al., 2011). The protein domain information and transcript variation were retrieved from UniProt (UniProt, 2019), Prosite (Sigrist et al., 2013), Pfam (El-Gebali et al., 2019) and ENSEMBL (Chen et al., 2010), respectively. The homology model of the S-protein and TMPRSS2 were constructed using Swiss model (Waterhouse et al., 2018), whereas 3D structure of the ACE2 was retrieved from the PDB database (PDB ID: 6M17). These structures were energy minimized by the Chiron energy minimization server (Ramachandran et al., 2011). The binding site residues of the proteins retrieved from Uniport and literature. The mutant structure of the TMPRSS2 protein was generated using WHATIF server (China et al., 1995) and energy minimized. The effect of the mutant was analysed using HOPE (Venselaar et al., 2010) and I-mutant (Capriotti et al., 2006). The I-mutant method allows to predict stability of the protein due to mutation. The docking studies for wild and mutant TMPRSS2 with S-protein and ACE2 were carried out using HADDOCK (Dominguez et al., 2003).

### **Network analysis on SARS-CoV-2 with Human Proteins**

We downloaded the SARS-CoV-2 genome (Accession number: MT121215) from the NCBI database. In order to find the Host-Pathogen Interactions (HPIs), the SARS-CoV-2 protein sequences were subjected to Host-Pathogen Interaction Database (HPIDB 3.0) (Ammari et al., 2016; Kumar and Nanduri, 2010). In addition to other host genes, we added two more proteins (TMPRSS2 and SLC6A19) which were found to have an important role in the mechanism of viral entry (Hoffmann et al., 2020). The protein-protein interaction (PPI) and transcription factor regulation of human proteins were retrieved from GeneMANIA (Warde-Farley et al., 2010) and literature (Barros et al., 2012; David et al., 2010; Maitland et al., 2011; Tumer et al., 2013; Yu et al., 2010; Zhang et al., 2009), respectively. The Host Pathogen Interaction Network (HPIN) was visualized using Cytoscape (Shannon et al., 2003) which includes the information collected from HPIDB, PPIs and TFs. Modules were defined as the set of statistics and functionally significant interacting genes (Reichardt and Bornholdt, 2006) which was constructed by using MCODE (Bader and Hogue, 2003). Further, Hubs were identified using Network analyzer (Assenov et al., 2008), the plugin of Cytoscape v3.2.1. We also studied the hub proteins' association with diseases, using DisGeNET (Menche et al., 2015).

The topological properties of the network and their functional significance was studied within the formalism of network theory which could predict important organizational and regulatory candidates in the network. We used network theoretical concepts, namely probability of degree distribution, clustering coefficient and average neighborhood connectivity, to characterize the structural and organizational features of the network (Details in supplementary Information). Removal of high degree nodes (hubs) in biological networks may cause lethality to the corresponding organism; the phenomenon has been referred to as the centrality–lethality rule (Jeong et al., 2001), verified by various genomic investigations. This idea of understanding enabled us to determine the topological features of the network architecture, important functionally related modules and hubs (Malik et al., 2019; Nafis et al., 2016). The topological properties were calculated after each subsequent removal, using the Network Analyzer plugin in Cytoscape version 3.7.1. The calculated topological properties of knock-out experiment were compared with those of the main network and the change in the properties helped in understanding the role of leading hubs in the network.

## Results and Discussion

SARS-CoV-2 uses ACE2 for entry, and its S-protein priming by the serine protease TMPRSS2 is a key factor. Despite the availability of other proteases, like cathepsin B and L, in the host cells, yet only TMPRSS2 activity is essential for pathogenesis and spread (Hoffmann et al., 2020).

### Dynamics and binding analysis of ACE2 with SARS CoV-2 RBD in dimeric and monomeric state

Mutagenesis experiments have reported the importance of Arginine and Lysine residues in ACE2, positioned between 697 to 716 as an important recognition site for TMPRSS2 mediated cleavage (**Figure S1a**) of ACE2. Mapping these residues on the recently resolved structure (61M7.pdb) suggests that this region lies in symmetric dimerization interface between ACE2 homodimer; and is also masked from outside by two BoAT1 in an overall 1:2:1 hetero-tetramer of B0AT1 and ACE2 (**Figure S1a**) (Yan and Zhang, 2020). In a recently reported study, the overexpression of *ACE2* has been reported to have an overall protective role in viral infection among patients (Vaduganathan et al., 2020). Increased density of ACE2 on the cell surface has been observed to protect from lung injury (Imai et al., 2005); also observed in many other instances (Kuster et al., 2020). We propose that an increased expression of ACE2 due to various factors, including naturally occurring variations influencing expression of the gene, may increase collisions and hence result in homo-dimerization of this protein. Binding of B0AT1 to ACE2 is not dependent on ACE2 homodimerization, but its presence along with dimerised ACE2 effectively shields the latter from interaction with TMPRSS2, with a possible prevention or reduction in its efficiency to cleave ACE2, a hypothesis based on visualization of the available experimentally resolved structure 61M7.PDB (Yan and Zhang, 2020). This observation is significant since the cleaved ACE2 interacts more efficiently with the SARS CoV-2 RBD domain.

Additionally, we also explored if the protease domains (PD) of each of the homodimerized ACE2 could independently bind with whole of the Spike Glycoprotein without affecting each other or any steric hindrance (Yan and Zhang, 2020). In the absence of any experimentally resolved structure available in literature with ACE2 homodimer and complete SARS-CoV2 Spike Glycoproteins binding to these, we aligned the RBD domains of the 6M17.pdb with the prefusion SARS-CoV2 Spike Glycoprotein in the optimal ACE2 binding conformation (Wrapp et al., 2020; Yan and Zhang, 2020) with “up CTD1” and “open S1 subunit” (Song et al., 2018) (**Figure 1’**). Both Spike Glycoproteins (6VSB.pdb) aligned well with viral RBD of (6M17.pdb). Further these were also observed without any steric clash with the ACE2 or partner Spike Glycoprotein in modelled structure. Both of these protruded radially outward from ACE2 binding site at PDs but inwards, at anchorage sites on virus membrane, with distance of approximately 27 nanometers (nm) (**Figure 1’**). However, considering the reported inter spike distance, between 13-15 nm (Neuman et al., 2006) and relative positions of Spike proteins on the viral shell, given almost similar sizes of SARS-CoV (Goldsmith et al., 2004; Neuman et al., 2006) and SARS-CoV-2 (Chen et al., 2020b), these are anticipated to be either parallel or protruding radially outwards from its anchorage site, contrasting to the findings in modelled structure (**Figure 1**). With this finding of inverse orientation of the Spike protein from virus shell in modelled structure, we propose that in case of ACE2 dimerization, only one of the dimerized ACE2 can bind with the Spike protein due to steric hinderances. Thus, the other unbound ACE2 partner may participate in its physiological role and continue to protect from lung injury.

To understand the mechanistic effect of dimerization, we conceived and simulated two situations: one, where ACE2 dimer is bound to two B0AT1 and two RBD domains of SARS CoV-2; and second, where only monomeric ACE2 is bound to one viral RBD. From the RMSF analysis of the simulated trajectories, we observed a higher order fluctuation in monomeric ACE2 homo-dimerization interface; while the rest of the ACE2 displayed a similar fluctuation profile, including residues interacting with RBD (**Figure S1’’**). PCA analysis and amino acid residue loadings on the PCs (PC1 in this case) reported slightly higher away ward conformational dynamics w.r.t RBD in dimeric ACE2 (**Figure S1’’**). This contrasting essential dynamics in dimeric ACE2 is explained well through the MM-GBSA based free energy of binding calculation between the 10th nanosecond final conformation of ACE2-RBD complex. From both simulations through MM-GBSA method, 1.5 fold strong binding of monomeric ACE2-RBD complex compared to dimeric ACE2-RBD, was observed (*i.e.*  $\Delta G = -75.58$  for Dimeric ACE2-RBD vs  $\Delta G = -116.98$  for Monomeric ACE2-RBD Complex). This observation further explains the protective role of overexpressed ACE2, resulting in its dimerization and reduced affinity for the RBD domain of SARS CoV-2 S1 protein. However, we are also not clear whether the presence of B0AT1 has any allosteric effect on this observation, in addition to already observed masking of ACE2 from TMPRSS2 mediated cleavage. Additionally, another metalloprotease ADAM17 has been reported to compete with the TMPRSS2 and cleave ACE2 in a way that only cleavage by TMPRSS2 was reported to drive the SARS CoV entry inside the cell (Heurich et al., 2014), which is yet to be explored in relation to



SARS CoV-2. Recent literature also indicates the potential role of other proteases, like cathepsin B/L that can functionally replace TMPRSS2, which needs to be evaluated extensively. (Sungnak et al., 2020) All these hypotheses warrant further *in-silico* analyses with better computational resources, as well as experimental study designs; and for these we are open to seek collaborations.

Yet, the observations made are of high importance and emphasise on evaluation of expression of the ACE2 and BoAT1 along with TMPRSS2 among patients with varied clinical response to SARS CoV-2 infection, differential outcomes, and correlation with observed mortality. Therefore, a differential expression of these genes among patients displaying varied responses from asymptomatic to acute symptoms is worth an exploration, which may act as biomarkers if proven to predict severity and susceptibility to COVID19.

### **Differential Genomics backgrounds derived expression of *ACE2* and *TMPRSS2***

The observations, from our MD analyses (**Figures S1''**) show that higher expression of ACE2 gene may promote ACE2 homo-dimerization, rendering less binding affinity to SARS CoV-2 RBD as well as masking TMPRSS2 cleavage site (**Figure 1**). Based on evidences appearing in literature of eQTLs and ACE2 higher expression in East Asians (Chen et al., 2020a) as well as reported protective role of ACE2 in lung injury (Imai et al., 2005; Kuster et al., 2020; Vaduganathan et al., 2020), we explored more about ACE2 expression and its differential genomic backgrounds in different population groups of the world.

### ***ACE2* Expression and functional SNPs related to its expression**

Tissue Specific evaluation of ACE2 gene in GTEx portal indicated low level expression of the ACE2 in Lungs (**Figure S2a and S3**). Further, it was observed that not all the transcripts and exons of the gene express in Lungs (**Figure S2b and c**). As expression of *ACE2* is less in Lungs thus, GTEx portal did not return cis-eQTLs for *ACE2* gene in lungs, in all probabilities. However, eQTLs for the gene were found in other tissues (**Figure 2a and b**) in GTEx portal. These eQTLs remained the same across tissue sets and showed similar effect in expression patterns (**Figure 2b**). Yet to correlate the expression of the gene with genomic variations, regions of the gene were explored for potential regulatory elements and overlapped with common variations through table browser tool of UCSC genome browser with an assumption that the variants having population level effect should be common variants (**Figure S2d and S4**). Variations data was retrieved from 1000 Genomes Phase 3 dataset and filtered for variants with at least 10% frequency of the alternate allele in global population. Overlapping the variants from both the exercises resulted in shortlisting of 2 SNPs rs1978124 and rs2106809, which were observed with differential frequency distribution in different populations of the world (**Figure S5**) at both super-population group and sub-population group level. Interestingly, EAS and EUR sub population groups were observed to show relatively uniform frequency distribution within group. However, in AFR, AMR and SAS sub population groups, intra population group differences were observed to be higher,

indicating diversity in the gene pool of population groups. HaploReg annotations indicated that these SNPs are in a region with Enhancer histone marks and DNAase activity. The same were observed through multiple regulatory elements tracks in UCSC Genome browser (**Figure S2d and S4**) including GH0XJ015596 enhancer marked by GeneHancer database. (Fishilevich et al., 2017) Linkage Disequilibrium (LD) values as ( $r^2$ ) with other putative proxy functional variants were also explored and mapped with UCSC genome browser. Interestingly, it was observed that SNPs rs1978124 and rs2106809 have a strong LD block of >100kb with various SNPs in absolute LD but upstream of ACE2 gene across population groups (**Figure S6**). This indicates a strong enhancer activity from the region that may affect higher expression of ACE2 gene in lungs, which may facilitate ACE2 homo-dimerization (**Figure 1**). However, it also requires experimental validation.

### ***TMPRSS2* Expression and functional SNPs in the gene with potential effect**

GTEx portal indicated relatively high levels of *TMPRSS2* expression in Lungs (**Figure S2a and S7c**), differential transcription (**Figure S7c**) but no protein expression in lungs (**Figure S7a**). Evaluation of the Protein atlas portal, source of **Figure S7a**, indicated that both antibodies used for detection *TMPRSS2* (not shown in Ms) were restricted to target near N terminal of the protein (either cytoplasmic domain or proximal extracellular domain near membrane). Evaluation of the GTEx portal for cis-eQTLs for *TMPRSS2* gene in lungs, returned a large number of eQTLs but with a peculiar feature (**Figure 3a and b**). The eQTLs in lungs were different as compared to other tissues and had following features: found to be concentrated in region of the gene with potential alternate transcripts (**Figure a,b and S8**), towards end of the gene in relatively high expressing exons (**Figure 4b**) and coding for the amino acid sequence that has putative functional role in protein as serine protease domain (**Figure 4a**) critical for ACE2 cleavage and SARS CoV-2 S-Protein priming. Variations data was retrieved from 1000 Genomes Phase 3 dataset and filtered for variants with at least 10% frequency of the alternate allele in global population and overlapped with cis-eQTLs data, resulting in 10 SNPs (rs463727, rs55964536, rs4818239, rs734056, rs4290734, rs2276205, rs34783969, rs11702475, rs62217531 and rs383510). Annotation of the SNPs indicated that all the SNPs clustered together and were in strong LD block (**Figure S8 and S9a**). HaploReg annotations indicated Enhancer histone marks and DNAase activity in the regions overlapping these SNPs. Amongst these, rs4818239 showed a prominent putative functional role (**Figure S9b**) which requires experimental validation. However, the frequency distribution of the variant rs4818239 in different populations groups of 1000G showed an interesting differential pattern (**Figure S9c**). We also explored if there were any alternate functional variations, yet common in populations, by screening *TMPRSS2* gene through genomAD browser and filtered for only missense variations. The search returned rs75603675 (NP\_001128571.1:p.Gly8Val or G8V) and rs12329760 (NP\_001128571.1:p.Val197Met or V197M), also observed with differential frequency distribution in 1000G populations (**Figure S9d and e**). The observations indicate *TMPRSS2* variants may influence interaction with *ACE2* as well as SARS CoV-2 (**Figure 1**) resulting in population specific differential outcomes.

### **The variation in TMPRSS2 could inhibit the ingress of SARS-CoV-2**

We further opted to explore the structure and function of TMPRSS2 protein. Pairwise sequence alignment of two isoforms of TMPRSS2 suggested that Isoform-2 lacked 37 residues at N-terminal compared to Isoform-1, which was the longest transcript and coded for 529 amino acids (**Figure S10**). Since the human TMPRSS2 protein structure was not available in the PDB database, we generated a computational protein model. The model was built using Serine protease hepsin (PDB ID: 5CE1) of homo sapiens. The protein 3D structure modelled from 146-491 residues of TMPRSS2 with sequence identity, GMQE and QMEAN of 33.82%, 0.53 and -1.43 values, respectively. It showed that the model was constructed with high confidence and best quality. In a similar manner 3D structure of S-Protein of SARS-CoV-2 was built with sequence identity, GMQE and QMEAN of 99.26%, 0.72 and -2.81 values, respectively. The protein stability analysis showed that rs12329760 (p.Val197Met or V197M) (Isoform-1) variation could decrease the stability of TMPRSS2 protein with  $\Delta\Delta G$  value of -1.51 kcal/mol. The HOPE results suggested, V197M variation is located within a domain, SRCR (GO Term: Scavenger Receptor Activity as annotated in UniProt) and introduces an amino acid with different properties, which could disturb this domain and abolish its function. Analyses of another variation rs75603675 (p.Gly8Val or G8V) showed that wild-type residue, glycine, providing flexibility might be necessary for the protein function and an alteration in this position could abolish this function, as the observed torsion angles for this residue were unusual. It could be speculated that only glycine was flexible enough to make these torsion angles, and a change at the location into another residue would force the local backbone into an incorrect conformation, disturbing the local structure. The variant residue was also observed to be more hydrophobic than the wild-type residue. Since sequence similarity search did not find any significant template at the N-terminal of this protein, we could not generate a quality model for isoform-2, which contains G8V variation; hence we carried out sequence-based stability analysis. The results suggested that G8V variation could increase the stability of the TMPRSS2 protein with  $\Delta\Delta G$  value of -0.10 kcal/mol. Further, it is known that Arginine and lysine residues within amino acids 697 to 716 are essential for efficient ACE2 cleavage by TMPRSS2 (Heurich et al., 2014) and recent studies have shown that SARS-CoV-2 uses the ACE2 for entry and the serine protease TMPRSS2 for S-protein priming. Based on information from previous studies, we docked the TMPRSS2 p.Val197Met wild-type and variant protein with ACE2 and S-protein of SARS-CoV-2. The docking results suggest that the variant V197M protein could promote the binding to ACE2 and inhibit the binding with S-protein (**Figure S11**). However, these observations need critical reevaluation as well as experimental work to understand these interactions better.

***SLC6A19* (B0AT1) expression naturally in Lungs and other respiratory tract cells, may provide protection.**

One of the interesting gene, *SLC6A19*, that codes for protein B0AT1, expresses in a very limited number of tissues and is reported to be absent in Lungs (**Figure S2a and S12a,b,c**). Our findings indicated its protective role with competitive hinderance in binding to TMPRSS2. As its expression was observed to be very low in Lungs, we resolved to look for indirect signatures that may have putative role in providing differential susceptibility. We noticed several variations clustering together in a potential enhancer region (**Figure S12d and e**) and with differential frequencies in different populations. eQTL analyses at GTEx portal also showed a huge list of potential SNPs upregulating or down regulating *SLC6A19* expression in Pancreas, liver, and whole blood cells, overlapping with potential Transcription factor binding sites (**Figure S12f**). We hypothesize a potential chance of some natural occurring variation/s that could induce expression of BoAT1 in respiratory tract thus, providing a protective role in this scenario (**Figure 1**). However, this requires extensive computational exploration as well as experiential validations.

### Host-Pathogen Interaction Modelling

We also believed that there are additional genes and factors which might be playing an additional role in providing differential susceptibility to COVID19. Thus, we carried out a systems biology study to identify the novel key regulators that may influence the Human and SARS CoV-2 interaction. A detailed Host Pathogen Interaction Network (HPIN) was created. The constructed HPIN contained 163 interactions, involving 31 nodes, which included 4 viral proteins, 27 human proteins and 8 Transcription Factors (TF) (**Figure 5a**). From the HPIN, we identified hubs, namely *RPS6*, *NACA*, *HNRNPA1*, *BTF3* and *SUMO1* with 19, 18, 17, 16 & 12 degrees, respectively in the network (**Figure 5a**). This indicated the affinity to attract a large number of low degree nodes towards each hub, which is a strong evidence of controlling the topological properties of the network by these few hubs (Good et al., 2011). Interestingly, out of five significant hubs, four hubs (*RPS6*, *NACA*, *HNRNPA1* and *BTF3*) present in module (Nodes: 16, Edges: 118, Score: 15.73) and one hub *SUMO1* was present at motif level which is considered one of the most important regulating motif of biological network at a fundamental level (**Figure 5b**). From our prediction; we found, four viral proteins (*S*, *N*, *ORF1a* & *ORF1ab*) target five host protein groups (*ACE2*, *SUMO1*, *HNRNPA1*, *RPS6* & *ATP6V1G1*). *SUMO1* & *HNRNPA1* are targeted by same viral proteins (N). The hub protein, *RPS6* (highest degree) directly interacted with one of the important protein *TMPRSS2* which further propagated the signal to *ACE2* & *SLC6A19*. The Transcription factor HIF1A and BCL6, STAT5, YBX1 inhibits *ACE2* and *SUMO1*; whereas MYC, AR and HNF4a, HNF1a activates *HNRNPA1*, *TMPRSS2* and *SLC6A19* respectively. From the gene diseases association study few diseases were highly associated with these important hub proteins (**Figure S13a**); *HNRNPA1* (36%) followed by *RPS6* (25%), *SUMO1* (21%), *BTF3* (9%) and *NACA* (9%) (**Figure S13b**). Clinically, patients with COVID-19 present with respiratory symptoms, Anoxia, fatigue, heart failure etc could be associated with these hub proteins, mainly *HNRNPA1* & *SUMO1*.

By hub removal methodology; we tried to understand the effect of hub removal and calculated the topological properties of the HPIN as a control. The probability of degree distribution (P(k)), Clustering coefficient (C(k))

and Neighbourhood connectivity ( $C_N(k)$ ) showed that the HPIN followed a power law scaling behavior. The power law behaviour was also checked and confirmed by using statistical test for power law fitting (p-value  $\geq 0.1$ ) (Clauset et al., 2009) in the hub-removal process. The removal of *RPS6*, *HNRNPA1*, *SUMO1*, *NACA* & *BTF3* from the HPIN brings significant variations in the topological properties of the HPIN where, degree distribution ( $\alpha$  and  $\beta$ ) change significantly (**Figure S14a**). Similarly, the variations in the measurements of the exponents of clustering coefficient and neighbourhood connectivity ( $\lambda$  and  $\mu$ ) also showed significant (Figure S14a). The knockout experiment of hub could able to highlight the local perturbations driven by these hubs and their effect on global network properties. The result suggest removal of *NACA* and *RPS6* hubs could turn down the degree distribution, so it could be crucial for communicating the signals (**Figure S14b**). In case of clustering coefficient, perturbation of *BTF3* shows minimum  $\gamma$ , indicating removal of *BTF3* makes the network less compact, which may lead to the delay in flow of signal. The HPIN perturbation increases in case of *BTF3* hub removal. In Neighbourhood connectivity increase in  $\mu$  indicates that the information processing in the network becomes faster when *SUMO1* hubs are removed, which means that local perturbations due to removal of hubs are strong enough to cause significant change in global scenario (Canright and Engø-Monsen, 2004) (**Figure S14b**). This indicates hubs are not robust, but it helps in the stability of the network.

Overall network analysis showed *ACE2* is not only the key molecule for entry and survival of SARS-CoV-2 virus, the hub proteins like *RPS6*, *HNRNPA1*, *SUMO1*, *NACA* & *BTF3* might also play a vital role. Analysing these interactions could provide further important understanding for the underlying biological mechanism of SARS-CoV-2 virus infection and identifying putative drug targets.

To conclude (**Figure 1**), higher expression of *ACE2* is facilitated by natural variations with different frequencies in different populations and is functionally associated with expression of the gene. The higher expression of *ACE2* facilitates homo-dimerization resulting in hindrance to TMPRSS2 mediated cleavage of *ACE2*. It becomes more difficult in presence of B0AT1 that usually does not express in Lungs. We also propose that the monomeric *ACE2* has higher preferential binding with SARS-CoV-2 S-Protein vis-a-vis its dimerized counterpart. Further, natural variations in *TMPRSS2*, with potential functional role, and their differential frequencies may also result in differential outcomes towards interaction with *ACE2* or priming of viral S-protein, a critical step for entry of virus in host cells. In addition, we have identified some other potential key host genes like *ADAMI7*, *RPS6*, *HNRNPA1*, *SUMO1*, *NACA* and *BTF3*, that might have a critical role. With all this background, it is anticipated that in populations like Indian populations, with highly diverse gene pool, a great variation in clinical outcomes is expected and that could be population/region specific, but primarily due to gene pool structure of the region, despite similar exposure levels to SARS CoV-2 and resources. Understanding these population specific differences may help in developing appropriate management strategies.

## **Acknowledgment**

All the authors acknowledge Prof. RNK Bamezai (Padamshri), former Vice Chancellor Shri Mata Vaishno Devi University, Katra and former Professor and Coordinator, National Centre for Applied Human Genetics, School of Life Sciences, Jawaharlal Nehru University, New Delhi, for his critical suggestions through various virtual discussion rounds, resulting in present study design. Authors also acknowledge Prof. Gyaneshwar Chaubey, Banaras Hindu University, UP, India for inputs, suggestions and assistance with geographical mapping of the allele frequencies from 1000G dataset.

## **Author Contributions**

SS conceived the concept. IP carried out molecular modelling simulations of ACE2 and related analyses. KP carried out TMPRSS2 modelling and related analyses. SS and ER carried out Human Population data screening and related work. SH and MZM carried out network analyses to identify additional key genes. SS, IP, SH, MZM designed, analysed, executed and wrote their part of manuscript. SS interpreted the results together, ER assisted SS in compilation of figures and SS compiled the overall MS.

## **Competing Interests**

The authors declare no competing interests associated with MS. For declaration purposes, SS is founder, chief scientific advisor of a startup “Biodroid Innovations Pvt Ltd” and IP is director of “Bioinfores Pvt. Ltd.”.

## REFERENCES

- (2006). Proceedings of the 2006 ACM/IEEE conference on Supercomputing (Tampa, Florida: Association for Computing Machinery).
- Ammari, M.G., Gresham, C.R., McCarthy, F.M., and Nanduri, B. (2016). HPIDB 2.0: a curated database for host-pathogen interactions. *Database (Oxford)* 2016.
- Assenov, Y., Ramirez, F., Schelhorn, S.E., Lengauer, T., and Albrecht, M. (2008). Computing topological parameters of biological networks. *Bioinformatics* 24, 282-284.
- Bader, G.D., and Hogue, C.W. (2003). An automated method for finding molecular complexes in large protein interaction networks. *BMC Bioinformatics* 4, 2.
- Barros, P., Lam, E.W., Jordan, P., and Matos, P. (2012). Rac1 signalling modulates a STAT5/BCL-6 transcriptional switch on cell-cycle-associated target gene promoters. *Nucleic Acids Res* 40, 7776-7787.
- Borgatti, S., Carley, K., and Krackhardt, D. (2006). On the Robustness of Centrality Measures Under Conditions of Imperfect Data. *Social Networks* 28, 124-136.
- Canright, G., and Engø-Monsen, K. (2004). Roles in networks. *Science of Computer Programming* 53, 195-214.
- Cao, Y., Li, L., Feng, Z., Wan, S., Huang, P., Sun, X., Wen, F., Huang, X., Ning, G., and Wang, W. (2020). Comparative genetic analysis of the novel coronavirus (2019-nCoV/SARS-CoV-2) receptor ACE2 in different populations. *Cell Discov* 6, 11.
- Capriotti, E., Calabrese, R., and Casadio, R. (2006). Predicting the insurgence of human genetic diseases associated to single point protein mutations with support vector machines and evolutionary information. *Bioinformatics* 22, 2729-2734.
- Chen, J., Jiang, Q., Xia, X., Liu, K., Yu, Z., Tao, W., Gong, W., and Han, J.J. (2020a). Individual Variation of the SARS-CoV2 Receptor ACE2 Gene Expression and Regulation. Pre prints, preprints.org > 202003.200191.v202001.
- Chen, N., Zhou, M., Dong, X., Qu, J., Gong, F., Han, Y., Qiu, Y., Wang, J., Liu, Y., Wei, Y., *et al.* (2020b). Epidemiological and clinical characteristics of 99 cases of 2019 novel coronavirus pneumonia in Wuhan, China: a descriptive study. *Lancet* 395, 507-513.
- Chen, Y., Cunningham, F., Rios, D., McLaren, W.M., Smith, J., Pritchard, B., Spudich, G.M., Brent, S., Kulesha, E., Marin-Garcia, P., *et al.* (2010). Ensembl variation resources. *BMC Genomics* 11, 293.
- China, G., Padron, G., Hooft, R.W., Sander, C., and Vriend, G. (1995). The use of position-specific rotamers in model building by homology. *Proteins* 23, 415-421.
- Clauset, A., Shalizi, C.R., and Newman, M.E.J. (2009). Power-Law Distributions in Empirical Data. *SIAM Review* 51, 661-703.
- David, C.J., Chen, M., Assanah, M., Canoll, P., and Manley, J.L. (2010). HnRNP proteins controlled by c-Myc deregulate pyruvate kinase mRNA splicing in cancer. *Nature* 463, 364-368.
- Dominguez, C., Boelens, R., and Bonvin, A.M. (2003). HADDOCK: a protein-protein docking approach based on biochemical or biophysical information. *J Am Chem Soc* 125, 1731-1737.
- El-Gebali, S., Mistry, J., Bateman, A., Eddy, S.R., Luciani, A., Potter, S.C., Qureshi, M., Richardson, L.J., Salazar, G.A., Smart, A., *et al.* (2019). The Pfam protein families database in 2019. *Nucleic Acids Res* 47, D427-D432.
- Fishilevich, S., Nudel, R., Rappaport, N., Hadar, R., Plaschkes, I., Iny Stein, T., Rosen, N., Kohn, A., Twik, M., Safran, M., *et al.* (2017). GeneHancer: genome-wide integration of enhancers and target genes in GeneCards. *Database (Oxford)* 2017.
- Goldsmith, C.S., Tatti, K.M., Ksiazek, T.G., Rollin, P.E., Comer, J.A., Lee, W.W., Rota, P.A., Bankamp, B., Bellini, W.J., and Zaki, S.R. (2004). Ultrastructural characterization of SARS coronavirus. *Emerg Infect Dis* 10, 320-326.
- Good, M.C., Zalatan, J.G., and Lim, W.A. (2011). Scaffold proteins: hubs for controlling the flow of cellular information. *Science* 332, 680-686.
- Grant, B.J., Rodrigues, A.P., ElSawy, K.M., McCammon, J.A., and Caves, L.S. (2006). Bio3d: an R package for the comparative analysis of protein structures. *Bioinformatics (Oxford, England)* 22, 2695-2696.
- Guan, W.J., Ni, Z.Y., Hu, Y., Liang, W.H., Ou, C.Q., He, J.X., Liu, L., Shan, H., Lei, C.L., Hui, D.S.C., *et al.* (2020). Clinical Characteristics of Coronavirus Disease 2019 in China. *N Engl J Med*.

- Heurich, A., Hofmann-Winkler, H., Gierer, S., Liepold, T., Jahn, O., and Pohlmann, S. (2014). TMPRSS2 and ADAM17 cleave ACE2 differentially and only proteolysis by TMPRSS2 augments entry driven by the severe acute respiratory syndrome coronavirus spike protein. *J Virol* *88*, 1293-1307.
- Hoffmann, M., Kleine-Weber, H., Schroeder, S., Kruger, N., Herrler, T., Erichsen, S., Schiergens, T.S., Herrler, G., Wu, N.H., Nitsche, A., *et al.* (2020). SARS-CoV-2 Cell Entry Depends on ACE2 and TMPRSS2 and Is Blocked by a Clinically Proven Protease Inhibitor. *Cell*.
- Humphrey, W., Dalke, A., and Schulten, K. (1996). VMD: Visual molecular dynamics. *Journal of molecular graphics* *14*, 33-38.
- Imai, Y., Kuba, K., Rao, S., Huan, Y., Guo, F., Guan, B., Yang, P., Sarao, R., Wada, T., Leong-Poi, H., *et al.* (2005). Angiotensin-converting enzyme 2 protects from severe acute lung failure. *Nature* *436*, 112-116.
- Jacobson, M.P., Friesner, R.A., Xiang, Z., and Honig, B. (2002). On the Role of the Crystal Environment in Determining Protein Side-chain Conformations. *Journal of molecular biology* *320*, 597-608.
- Jeong, H., Mason, S.P., Barabasi, A.L., and Oltvai, Z.N. (2001). Lethality and centrality in protein networks. *Nature* *411*, 41-42.
- Kaiser, J. (2020). How sick will the coronavirus make you? The answer may be in your genes. Health, Coronavirus, Sciencemagorg.
- Kumar, R., and Nanduri, B. (2010). HPIDB--a unified resource for host-pathogen interactions. *BMC Bioinformatics* *11 Suppl 6*, S16.
- Kuster, G.M., Pfister, O., Burkard, T., Zhou, Q., Twerenbold, R., Haaf, P., Widmer, A.F., and Osswald, S. (2020). SARS-CoV2: should inhibitors of the renin-angiotensin system be withdrawn in patients with COVID-19? *Eur Heart J*.
- Madhavi Sastry, G., Adzhigirey, M., Day, T., Annabhimoju, R., and Sherman, W. (2013). Protein and ligand preparation: parameters, protocols, and influence on virtual screening enrichments. *Journal of computer-aided molecular design* *27*, 221-234.
- Maitland, N.J., Frame, F.M., Polson, E.S., Lewis, J.L., and Collins, A.T. (2011). Prostate cancer stem cells: do they have a basal or luminal phenotype? *Horm Cancer* *2*, 47-61.
- Malik, M.Z., Chirom, K., Ali, S., Ishrat, R., Somvanshi, P., and Singh, R.K.B. (2019). Methodology of predicting novel key regulators in ovarian cancer network: a network theoretical approach. *BMC Cancer* *19*, 1129.
- Mark, P., and Nilsson, L. (2001). Structure and Dynamics of the TIP3P, SPC, and SPC/E Water Models at 298 K. *The Journal of Physical Chemistry A* *105*, 9954-9960.
- Menche, J., Sharma, A., Kitsak, M., Ghiassian, S.D., Vidal, M., Loscalzo, J., and Barabasi, A.L. (2015). Disease networks. Uncovering disease-disease relationships through the incomplete interactome. *Science* *347*, 1257601.
- Miller, A., Reandelar, M.J., Fasciglione, K., Roumenova, V., Li, Y., and Otazu, G.H. (2020). Correlation between universal BCG vaccination policy and reduced morbidity and mortality for COVID-19: an epidemiological study. medRxiv, 2020.2003.2024.20042937.
- Monteil, V., Kwon, H., Prado, P., Hagelkruys, A., Wimmer, R.A., Stahl, M., Leopoldi, A., Garreta, E., Hurtado Del Pozo, C., Prosper, F., *et al.* (2020). Inhibition of SARS-CoV-2 Infections in Engineered Human Tissues Using Clinical-Grade Soluble Human ACE2. *Cell*.
- Nafis, S., Ponnusamy, K., Husain, M., Singh, R.K., and Bamezai, R.N. (2016). Identification of key regulators and their controlling mechanism in a combinatorial apoptosis network: a systems biology approach. *Mol Biosyst* *12*, 3357-3369.
- Neuman, B.W., Adair, B.D., Yoshioka, C., Quispe, J.D., Orca, G., Kuhn, P., Milligan, R.A., Yeager, M., and Buchmeier, M.J. (2006). Supramolecular architecture of severe acute respiratory syndrome coronavirus revealed by electron cryomicroscopy. *Journal of virology* *80*, 7918-7928.
- Ramachandran, S., Kota, P., Ding, F., and Dokholyan, N.V. (2011). Automated minimization of steric clashes in protein structures. *Proteins* *79*, 261-270.
- Reichardt, J., and Bornholdt, S. (2006). Statistical mechanics of community detection. *Physical Review E* *74*, 016110.
- Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin, N., Schwikowski, B., and Ideker, T. (2003). Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res* *13*, 2498-2504.



- Shi, Y., Wang, Y., Shao, C., Huang, J., Gan, J., Huang, X., Bucci, E., Piacentini, M., Ippolito, G., and Melino, G. (2020). COVID-19 infection: the perspectives on immune responses. *Cell Death Differ.*
- Sievers, F., Wilm, A., Dineen, D., Gibson, T.J., Karplus, K., Li, W., Lopez, R., McWilliam, H., Remmert, M., Soding, J., *et al.* (2011). Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega. *Mol Syst Biol* 7, 539.
- Sigrist, C.J., de Castro, E., Cerutti, L., Cuče, B.A., Hulo, N., Bridge, A., Bougueleret, L., and Xenarios, I. (2013). New and continuing developments at PROSITE. *Nucleic Acids Res* 41, D344-347.
- Song, W., Gui, M., Wang, X., and Xiang, Y. (2018). Cryo-EM structure of the SARS coronavirus spike glycoprotein in complex with its host cell receptor ACE2. *PLoS Pathog* 14, e1007236.
- Stawiski, E.W., Diwanji, D., Suryamohan, K., Gupta, R., Fellouse, F.A., Sathirapongsasuti, J.F., Liu, J., Jiang, Y.-P., Ratan, A., Mis, M., *et al.* (2020). Human ACE2 receptor polymorphisms predict SARS-CoV-2 susceptibility. *bioRxiv*, 2020.2004.2007.024752.
- Sungnak, W., Huang, N., Bécavin, C., Berg, M., Queen, R., Litvinukova, M., Talavera-López, C., Maatz, H., Reichart, D., Sampaziotis, F., *et al.* (2020). SARS-CoV-2 entry factors are highly expressed in nasal epithelial cells together with innate immune genes. *Nature Medicine*.
- Tumer, E., Broer, A., Balkrishna, S., Julich, T., and Broer, S. (2013). Enterocyte-specific regulation of the apical nutrient transporter SLC6A19 (B(0)AT1) by transcriptional and epigenetic networks. *J Biol Chem* 288, 33813-33823.
- Uhlen, M., Fagerberg, L., Hallstrom, B.M., Lindskog, C., Oksvold, P., Mardinoglu, A., Sivertsson, A., Kampf, C., Sjostedt, E., Asplund, A., *et al.* (2015). Proteomics. Tissue-based map of the human proteome. *Science* 347, 1260419.
- UniProt, C. (2019). UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res* 47, D506-D515.
- Vaduganathan, M., Vardeny, O., Michel, T., McMurray, J.J.V., Pfeffer, M.A., and Solomon, S.D. (2020). Renin-Angiotensin-Aldosterone System Inhibitors in Patients with Covid-19. *N Engl J Med* 382, 1653-1659.
- Venselaar, H., Te Beek, T.A., Kuipers, R.K., Hekkelman, M.L., and Vriend, G. (2010). Protein structure analysis of mutations causing inheritable diseases. An e-Science approach with life scientist friendly interfaces. *BMC Bioinformatics* 11, 548.
- Walls, A.C., Park, Y.J., Tortorici, M.A., Wall, A., McGuire, A.T., and Velesler, D. (2020). Structure, Function, and Antigenicity of the SARS-CoV-2 Spike Glycoprotein. *Cell*.
- Wan, Y., Shang, J., Graham, R., Baric, R.S., and Li, F. (2020). Receptor Recognition by the Novel Coronavirus from Wuhan: an Analysis Based on Decade-Long Structural Studies of SARS Coronavirus. *Journal of Virology* 94, e00127-00120.
- Warde-Farley, D., Donaldson, S.L., Comes, O., Zuberi, K., Badrawi, R., Chao, P., Franz, M., Grouios, C., Kazi, F., Lopes, C.T., *et al.* (2010). The GeneMANIA prediction server: biological network integration for gene prioritization and predicting gene function. *Nucleic Acids Res* 38, W214-220.
- Waterhouse, A., Bertoni, M., Bienert, S., Studer, G., Tauriello, G., Gumienny, R., Heer, F.T., de Beer, T.A.P., Rempfer, C., Bordoli, L., *et al.* (2018). SWISS-MODEL: homology modelling of protein structures and complexes. *Nucleic Acids Res* 46, W296-W303.
- Wrapp, D., Wang, N., Corbett, K.S., Goldsmith, J.A., Hsieh, C.-L., Abiona, O., Graham, B.S., and McLellan, J.S. (2020). Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation. *Science (New York, NY)* 367, 1260-1263.
- Yan, R., and Zhang, Y. (2020). Structural basis for the recognition of SARS-CoV-2 by full-length human ACE2. *367*, 1444-1448.
- Yu, Y.N., Yip, G.W., Tan, P.H., Thike, A.A., Matsumoto, K., Tsujimoto, M., and Bay, B.H. (2010). Y-box binding protein 1 is up-regulated in proliferative breast cancer and its inhibition deregulates the cell cycle. *Int J Oncol* 37, 483-492.
- Zhang, R., Wu, Y., Zhao, M., Liu, C., Zhou, L., Shen, S., Liao, S., Yang, K., Li, Q., and Wan, H. (2009). Role of HIF-1alpha in the regulation ACE and ACE2 expression in hypoxic human pulmonary artery smooth muscle cells. *Am J Physiol Lung Cell Mol Physiol* 297, L631-640.
- Zhou, P., Yang, X.L., Wang, X.G., Hu, B., Zhang, L., Zhang, W., Si, H.R., Zhu, Y., Li, B., Huang, C.L., *et al.* (2020). A pneumonia outbreak associated with a new coronavirus of probable bat origin. *Nature* 579, 270-273.

## FIGURE LEGENDS

### **Figure 1. The overall interaction of host proteins and SARS-CoV-2 entry to cell**

The figure summarises various factors involved that can influence the entry of Virus in host cell. It depicts brief mechanisms that may arise due to natural occurring variations influencing expression of host genes or structural changes affecting interactions within host proteins or viral proteins resulting in effect on efficiency of virus entry in host cells which could be key factor is providing differential clinical outcomes in different population groups. The figure depicts, higher expression of ACE2 is facilitated by natural variations with different frequencies in different populations and functionally associated with expression of the gene. The higher expression of ACE2 facilitates homo-dimerization resulting in hindrance to TMPRSS2 mediated cleavage of ACE2. Monomeric ACE2 has higher preferential binding with SARS-CoV-2 S-Protein vis-a-vis its dimerized counterpart. It becomes more difficult in presence of B0AT1 that usually does not express in Lungs. Further, natural variations in TMPRSS2, with potential functional role, and their differential frequencies may also result in differential outcomes towards interaction with ACE2 or priming of viral S-protein, a critical step for entry of Virus in host cells. In addition, other potential key host genes like ADAM17, RPS6, HNRNPA1, SUMO1, NACA and BTF3 might have a critical role.

### **Figure 1'. Complex of homo dimer of ACE2 with complete SARS-CoV-2 Spike Glycoprotein**

Modelled complex of homo dimer of ACE2 colored as brick red (6M17.pdb) with complete Spike Glycoprotein in open state coloured as sky blue (6VSB.pdb) by aligning the complete Spike Glycoprotein with RBD colored as purple blue from 6M17.pdb. The head of the Spike protein can bind the dimer ACE2, protruding radially outward with RBD domains 6nm apart but inwards, at anchorage sites on virus membrane, with distance of approximately 27 nm.

### **Figure 2. The screenshot from GTEx portal depicting eQTLs of ACE2 Gene**

eQTLs in different tissues were plotted through (a) GTEx IGV Browser as well as (b) GTEx Locus Browser. (a) Red dots indicate eQTLs that showed up in the region against the query. Size of the dot in (b) indicate level of significance (as negative p values) whereas colour depicts positive or negative correlation with Normalized effect size (NES) of the eQTL from -1 to 0 to 1. Red color shades represent upregulation whereas blue color shades show downregulation. (b) also depicts Linkage disequilibrium (LD) with value range from 0 to 1 as white to black shades with 1(dark) as absolute LD.

### **Figure 3. The screenshot from GTEx portal depicting eQTLs of TMPRSS2 gene**

eQTLs in different tissues were plotted through (a) GTEx IGV Browser as well as (b) GTEx Locus Browser. (a) Red dots indicate eQTLs that showed up in the region against the query. Size of the dot in (b) indicate level of significance (as negative p values) whereas colour depicts positive or negative correlation with Normalized effect size (NES) of the eQTL from -1 to 0 to 1. Red color shades represent upregulation whereas blue color

shades show downregulation. (b) also depicts Linkage disequilibrium (LD) with value range from 0 to 1 as white to black shades with 1(dark) as absolute LD.

It was noted [also indicated by blue arrow in (a)] that where in other tissues eQTLs are mainly towards 5' UTR of the gene, in Lungs the eQTLs are towards end of the gene extending towards 3'UTR.

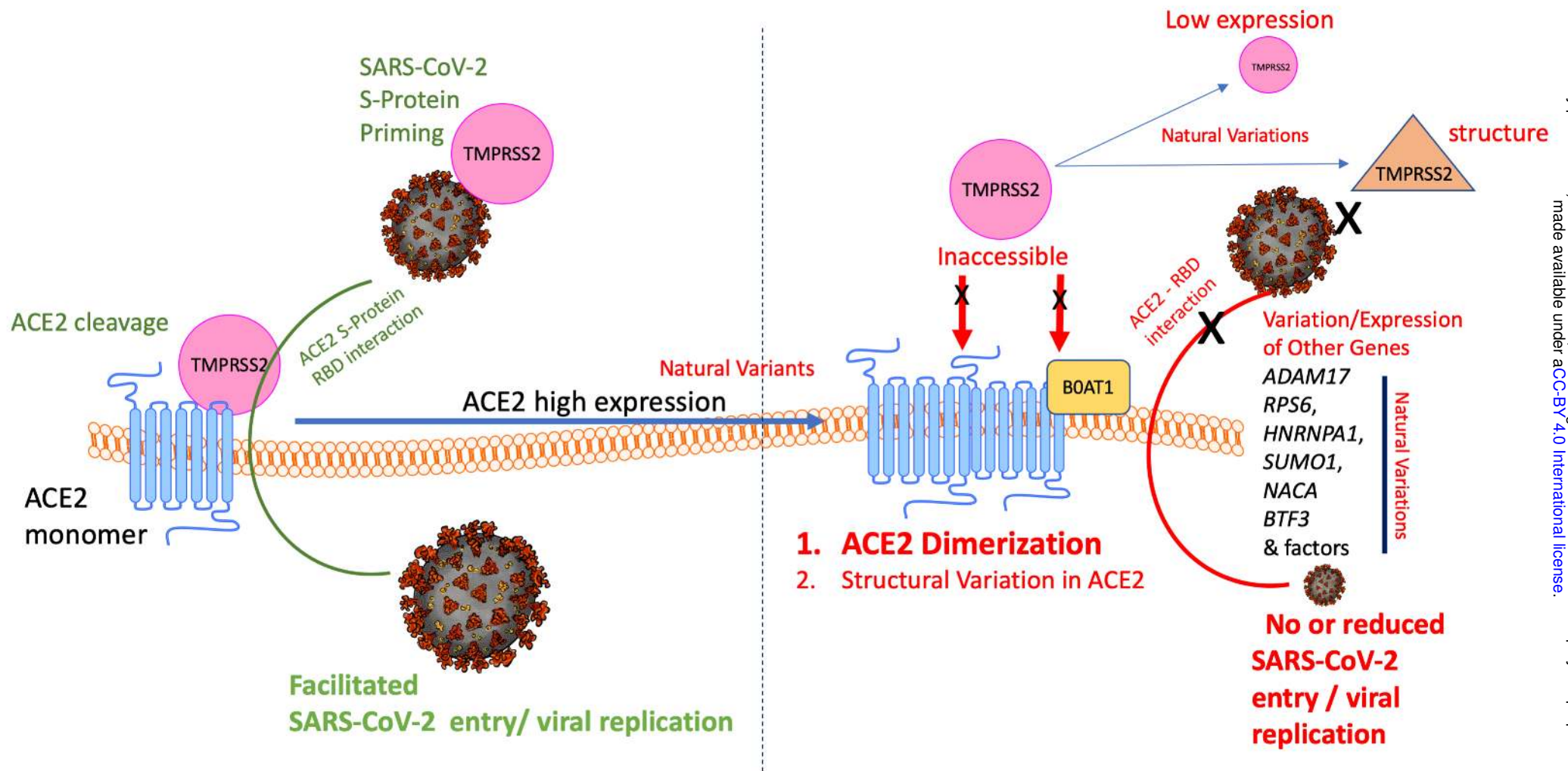
#### **Figure 4. Structure of TMPRSS2 gene, its alternative transcripts and functional domains**

(a) TMPRSS2 canonical transcript is constituted of 14 exons and alternate transcripts have been seen. The coded protein is a transmembrane protein with 1-85 amino acids (aa) forming Cytoplasmic Domain, and 112-492 aa constituting Extracellular domain. (b) Expression levels as median read count per base as shades of blue from low to high values indicating exons expression in different tissues in GTEx portal. Gene models of alternate transcripts are also depicted. Interesting to note exons coding for the extracellular domain have higher expression in Lungs (row marked by blue arrow) then other exons.

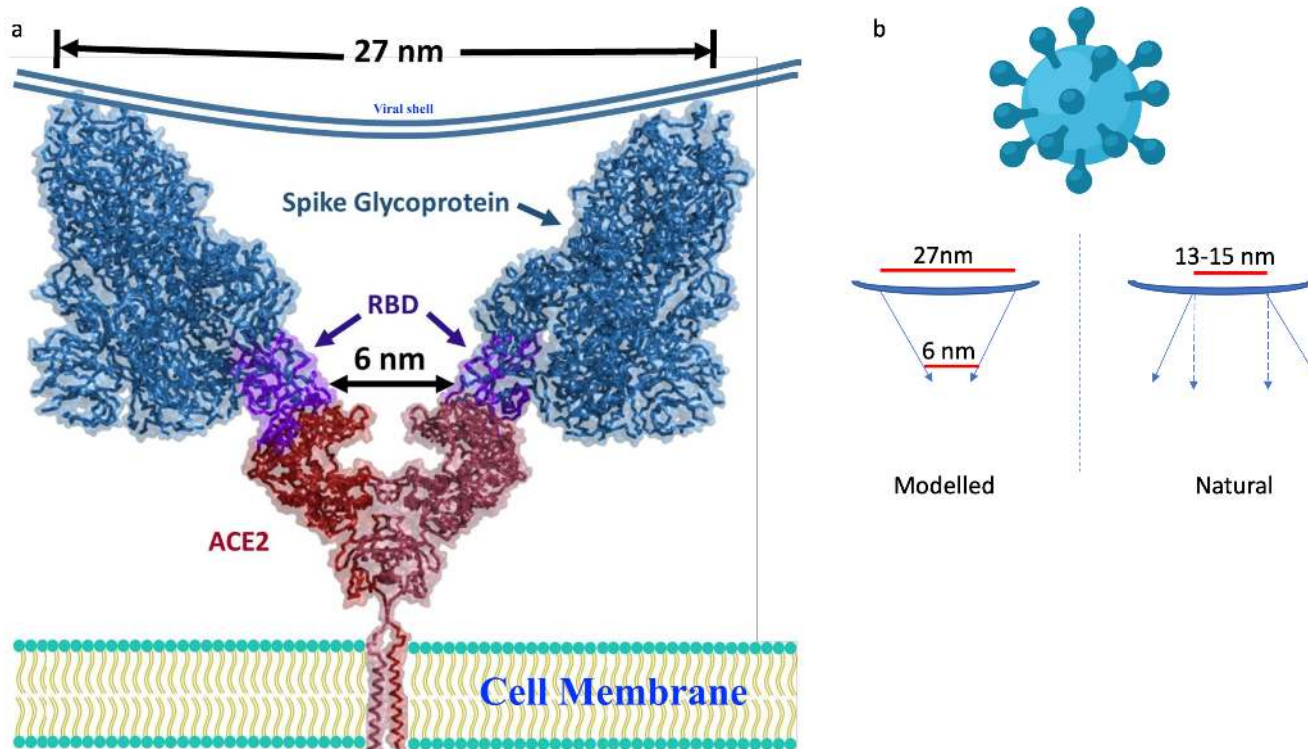
#### **Figure 5. Host-Pathogen Interaction Network and their significant hubs.**

(a) The network view of HPIN imported from Cytoscape. The Viral proteins, Human proteins and TFs are represented as nodes and edges denote the physical interaction. All the nodes of viral proteins (red), human proteins (blue) and TF (green) are filled triangles, circles and V-shaped respectively. The edges between Virus-human proteins are shown in orange-headed arrows, PPIs in grey lines and TF-human protein in pink arrow-headed (activators) and flat-headed (Inhibitors). The significant existence of sparsely distributed few main hub proteins, namely RPS6, NACA, HNPRNPA1, BTF3 and SUMO1 are colored as dark blue in the network were represented in the order of four enlarged sized circles.

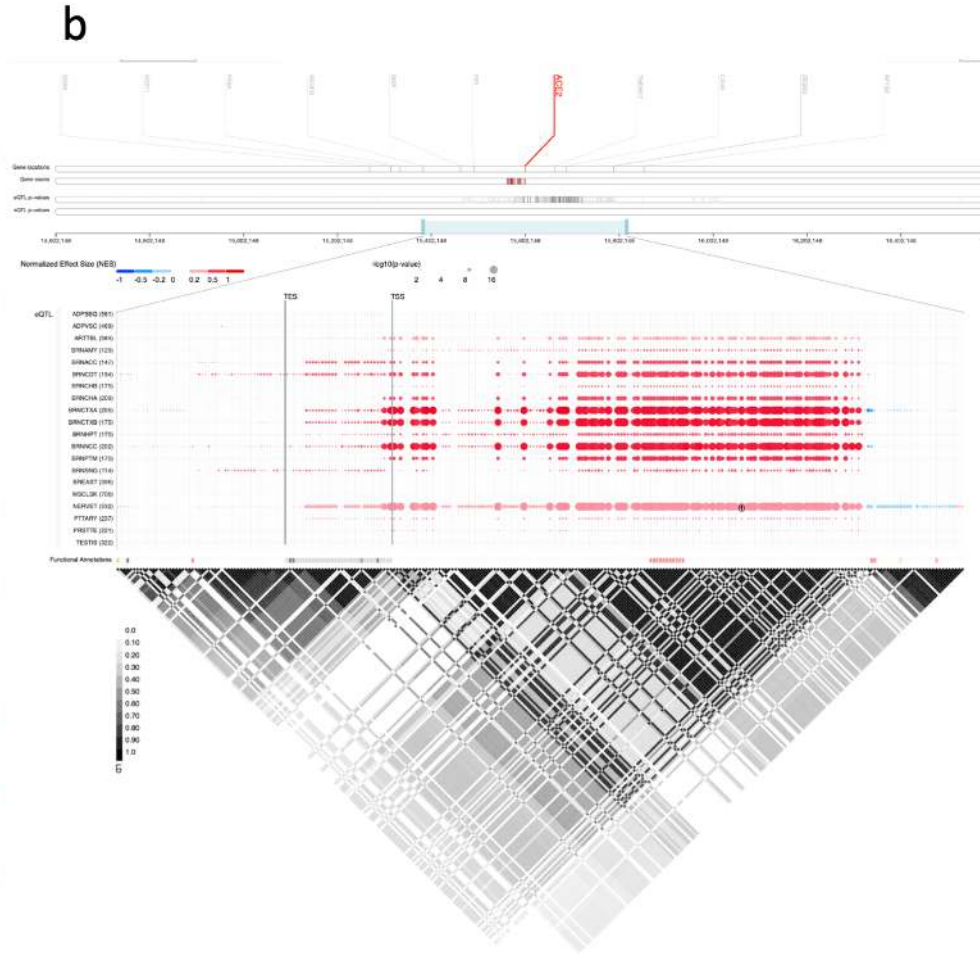
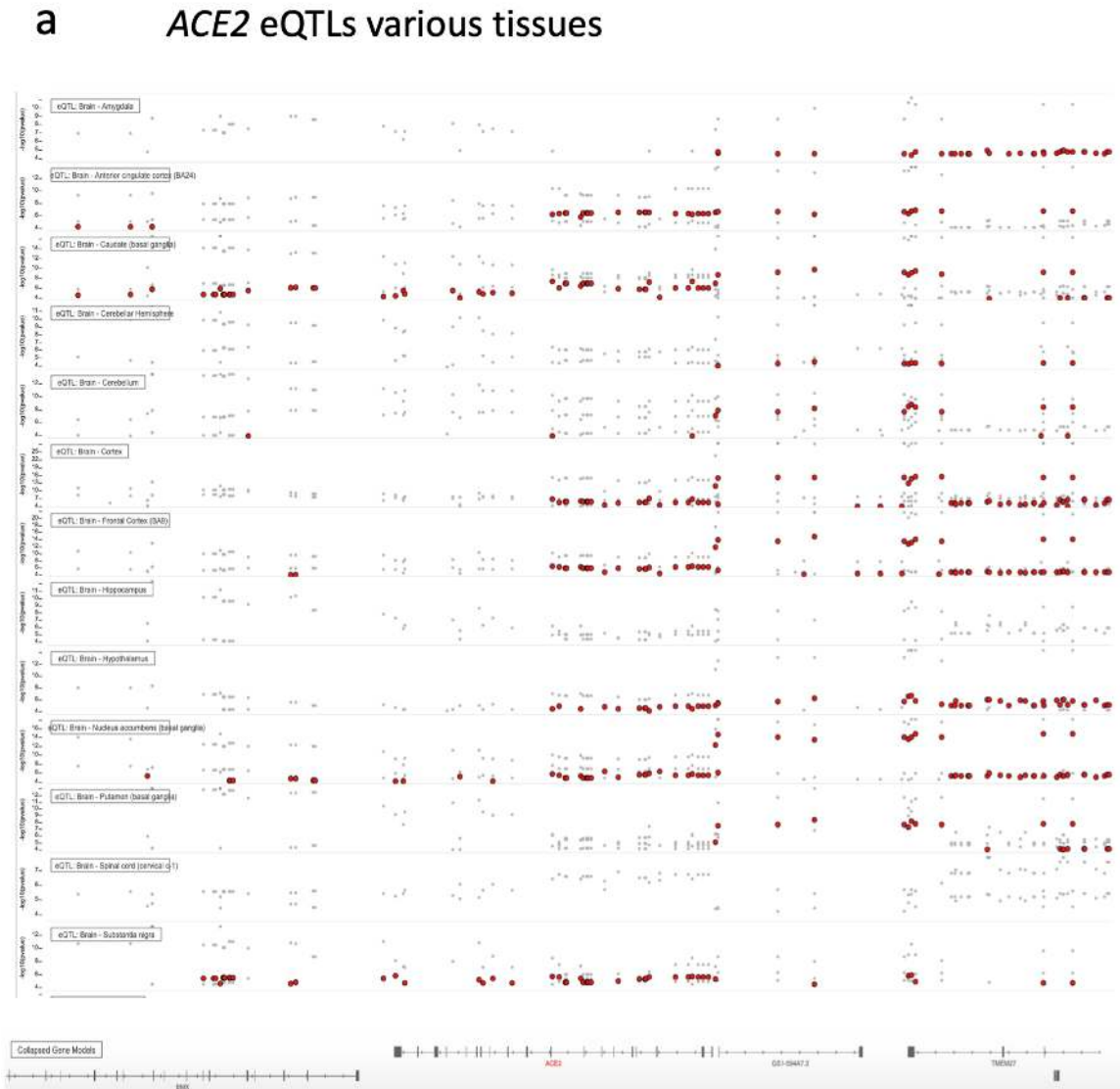
(b) The Module and Motif constructed and analyzed using MCODE. All the nodes and edges of module & motif are in blue and grey color, respectively. The significant hubs, present in module & motif highlighted in dark blue color.



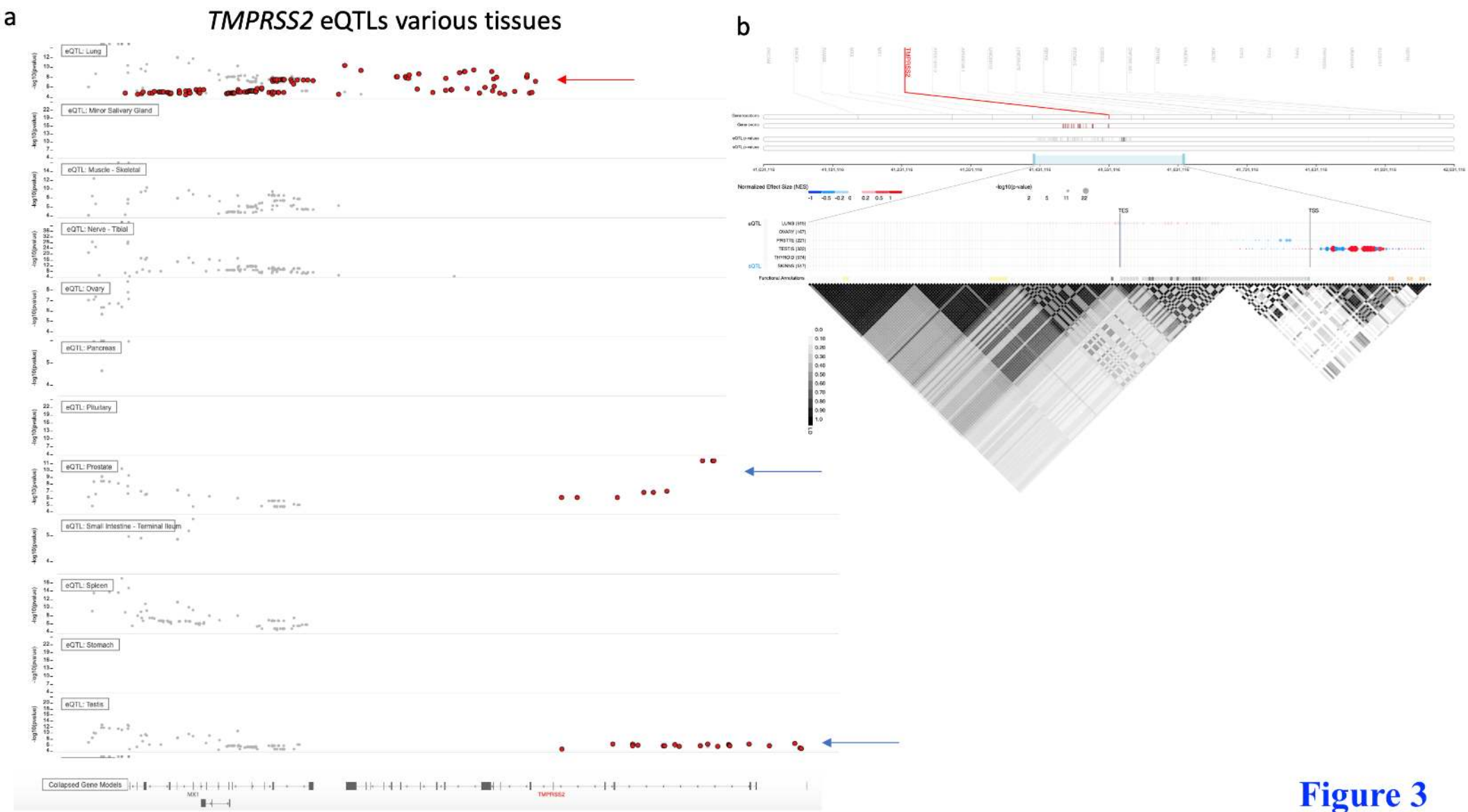
**Figure 1**



**Figure 1'**

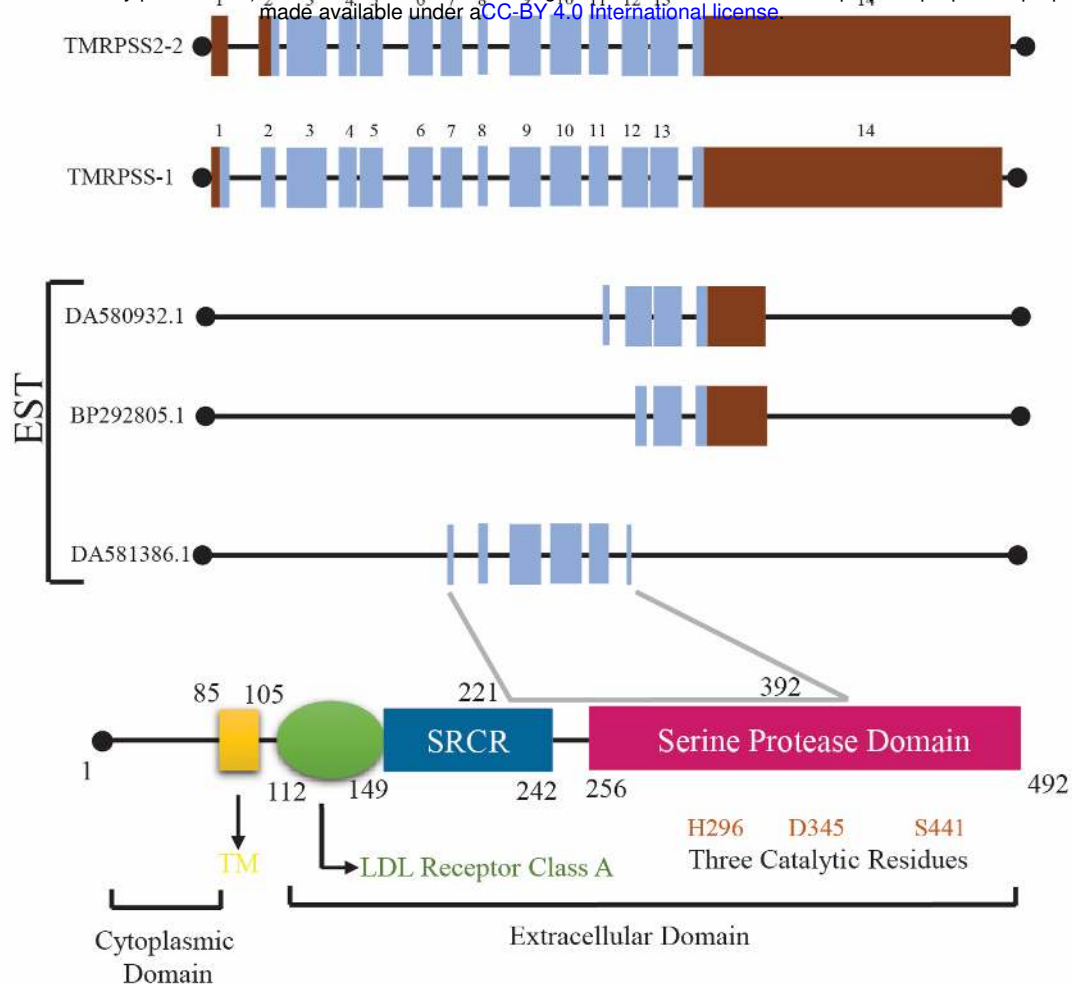


**Figure 2**



**Figure 3**

A.



B.

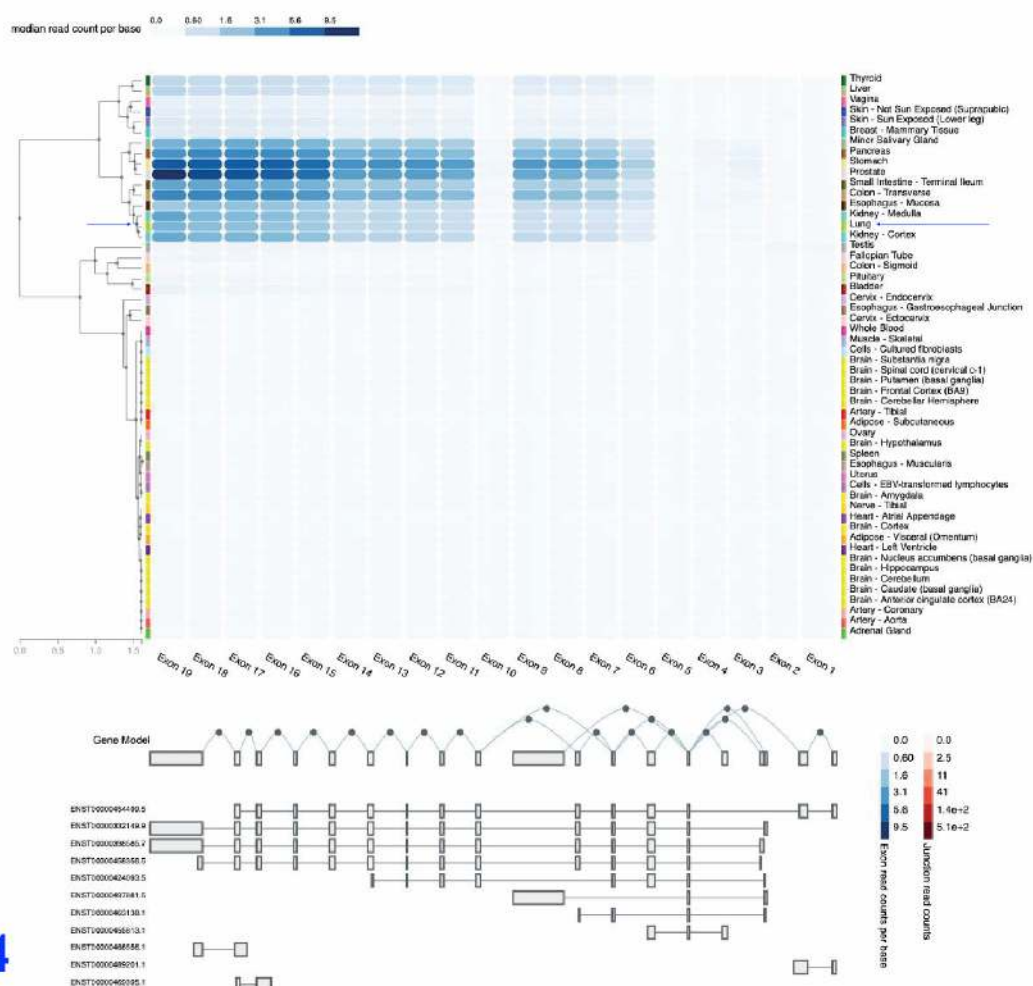
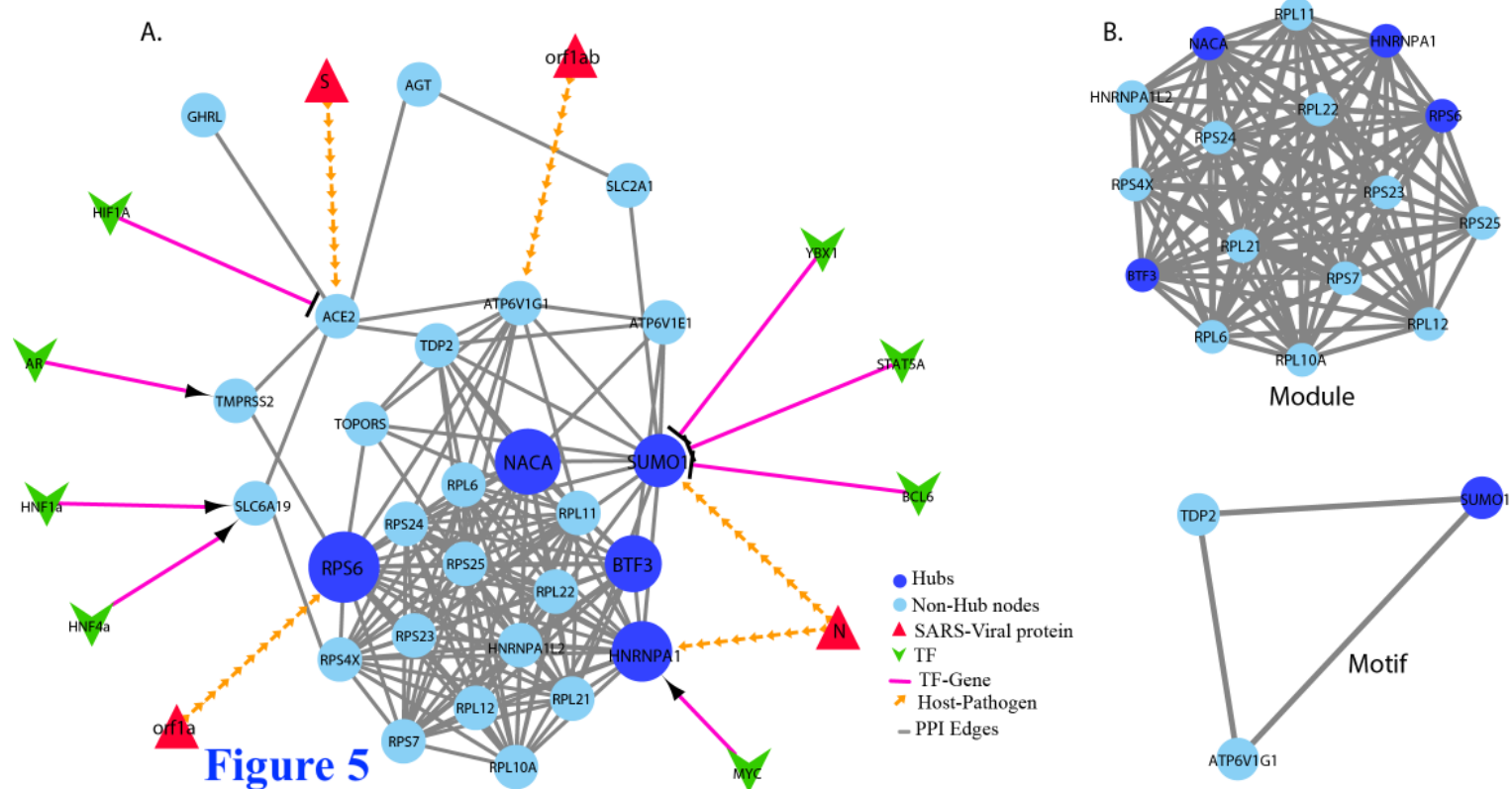


Figure 4





## Supplementary Information

### **ACE2 Homo-dimerization, Human Genomic variants and Interaction of Host Proteins Explain High Population Specific Differences in Outcomes of COVID19**

Swarkar Sharma<sup>1\*</sup>, Inderpal Singh<sup>2†</sup>, Shazia Haider<sup>3†</sup>, Md. Zubair Malik<sup>4†</sup>, Kalaiarasan Ponnusamy<sup>5†</sup>, Ekta Rai<sup>1</sup>

1. Human Genetics Research Group, School of Biotechnology, Shri Mata Vaishno Devi University, Katra, Jammu and Kashmir, India.
2. Bioinfores Pvt. Ltd., R. S. Pura, Jammu, Jammu and Kashmir, India.
3. Department of Biotechnology, Jaypee Institute of Information Technology, Noida, Sector-62, Uttar Pradesh, India
4. School of Computational and Integrative Sciences, Jawaharlal Nehru University, New Delhi, India.
5. School of Biotechnology, Jawaharlal Nehru University, New Delhi, India.

† - All authors contributed Equally

#### **Corresponding Author:**

Dr. Swarkar Sharma

Coordinator, Human Genetics Research Group,

School of Biotechnology, Shri Mata Vaishno Devi University, Katra, J&K, India

Email: [swarkar.sharma@smvdu.ac.in](mailto:swarkar.sharma@smvdu.ac.in)

Mobile: +91-9419955636; Ph: +91-1991-285535//285525 Ext. 2385

## Additional information about knock-out experiment in Network analysis

### Perturbation by leading hub removal analysis

**Degree Distribution ( $P(k)$ ):** Degree  $k$  is the number of interaction a node in the HPIN.  $P(k)$  is the probability of randomly chosen node to have  $k$  interaction with the neighbour. The probability of degree distribution ( $P(k)$ ) of the network is calculated by:

$$P(k) = \frac{n_k}{N}$$

where,  $n_k$  is equal to the number of nodes with degree  $k$  and  $N$  is equal to the size of the network (Albert and Barabási, 2002; Barabasi and Albert, 1999).

**Clustering Coefficient ( $C(k)$ ):** Clustering coefficient defines how strongly the nodes in a network tend to cluster together. Clustering coefficient of the  $i^{th}$  node in undirected network can be obtained by:

$$C(k_i) = \frac{2e_i}{k_i(k_i - 1)}$$

where,  $e_i$  is the number of connected pairs of the nearest-neighbour of the  $i$ -th node and  $k_i$  denotes the degree of the respective node (Ravasz and Barabasi, 2003; Ravasz, et al., 2002).

**Neighborhood Connectivity ( $C_N(k)$ ):** Neighbourhood connectivity of a node is defined as the connectivity of all the neighbours of the node (Maslov and Sneppen, 2002). The average connectivity of the nearest neighbours of a node is given by:

$$C_N(k) = \sum_q qP(q|k)$$

where,  $P(q|k)$  denotes the conditional probability that a link belonging to a node with connectivity  $k$  points to a node with connectivity  $q$ .

For scale free network,  $C_N(k)$  is constant while it follows power law  $C_N(k) = c^{-\mu}$  for hierarchical network with  $\mu$  approximately equal to 0.5 (Pastor-Satorras, et al., 2001). The positive and negative signs in  $\mu$  indicates the assortivity or disassortivity of network respectively (Barrat, et al., 2004).

## References

1. Albert, R. and Barabási, A.-L. Statistical mechanics of complex networks. *Reviews of Modern Physics* 2002;74(1):47-97.
2. Barabasi, A.L. and Albert, R. Emergence of scaling in random networks. *Science* 1999;286(5439):509-512.
3. Barrat, A., et al. The architecture of complex weighted networks. *Proceedings of the National Academy of Sciences of the United States of America* 2004;101(11):3747-3752.
4. Maslov, S. and Sneppen, K. Specificity and stability in topology of protein networks. *Science* 2002;296(5569):910-913.
5. Pastor-Satorras, R., Vazquez, A. and Vespignani, A. Dynamical and correlation properties of the internet. *Phys Rev Lett* 2001;87(25):258701.
6. Ravasz, E. and Barabasi, A.L. Hierarchical organization in complex networks. *Phys Rev E Stat Nonlin Soft Matter Phys* 2003;67(2 Pt 2):026112.

7. Ravasz, E., et al. Hierarchical organization of modularity in metabolic networks. *Science* 2002;297(5586):1551-1555.

## Supplementary figures

### Figure S1. Structural representation and essential dynamics of ACE2

A multi chain complex of dimeric ACE2 is shown (colored as purple and forest green surface representation) bound to B0AT1 (colored as chocolate and orange surface) through its C terminal amino acids and SARS-CoV-2 RBD domain through its N terminal amino acids (shown as black and blue surface), the recognition site of TMPRSS2 is also highlighted as surface against the cartoon background [a]. The cartoon representation along with CPK spheres blue-white-red showing the direction of essential dynamics calculated through PCAs PC1 of ACE2-RBD region from both simulations is shown [b and c].

### Figure S1". Root mean square fluctuation and contribution of ACE2 residues to PC1 and 2 along with proportion of variance

Root mean square fluctuation (RMSF) plots of ACE2 protein from monomeric ACE2-RBD viral complex and dimeric ACE2(single chain)-RBD complex are shown. The localized fluctuation profile of the amino acids across the protein were similar for both simulation conditions except higher fluctuation in the region that constitutes the ACE2 dimerization interface and following C terminal helix which interacts with B0AT1 (B0AT1 was deleted in the monomeric ACE2-RBD simulation [a]. The contribution of the amino acids to the principal component (PC) 1 and 2 (black and red respectively) along with the scree plot reporting the proportion of variance by each PC are shown [b and c]. It can be observed that the N terminal amino acids consisting of the ACE2-RBD binding site contributed more to the overall essential dynamics in dimeric compared monomeric ACE2 simulated in the present study.

### Figure S2. Expression of host genes in Lungs, alternative transcripts of ACE2 and expression data depiction and regulatory elements and other annotations related to ACE2

(a) Expression of various human genes in Human Lungs from GTEx portal. Normalized Expression Values as Transcripts per million (TPM) are plotted as shades from yellow to deep blue in low to high order. (b) Expression levels as median read count per base as shades of blue from low to high values indicating exons expression in different tissues in GTEx portal. Gene models of alternate transcripts are also depicted. (c) Expression levels as median read count per base as shades of red from low to high values indicating exons junctions expression in different tissues in GTEx portal. (d) UCSC genome browser showing various tracks at 5' Untranslated region of ACE2 gene, common dbSNPs 153 intersecting with conserved transcription factors and conserved transcription factor binding sites are also depicted.

### **Figure S3 The screenshot from Human Protein Atlas depicting Human ACE2 expression**

RNA and Protein expression in different tissues is depicted. In lungs, only RNA expression of ACE2 gene is shown with no ACE2 protein expression.

### **Figure S4. Annotation of ACE2 depicting various regulatory elements in UCSC Genome browser**

UCSC genome browser showing SNPs rs1978124 and rs2106809 in intron 1 of ACE2 gene. Various tracks covering first 2 exons and 5' Untranslated region of ACE2 gene are shown. The ACE2 gene is located on negative strand. In the region GH0XJ015596 enhancer (in red colour) can be located in the track Enhancer and promoter from GeneHancer database.

### **Figure S5. Frequency distribution of SNPs rs1978124 and rs2106809 in 1000G population groups**

Frequency distribution of SNPs, rs1978124 and rs2106809 in 1000G populations dataset, also depicted on world map by a web based tool Datawrapper (<https://www.datawrapper.de/>). Derived allele frequencies in both the SNPs are also plotted for comparison, indicating differential frequency distribution in different population groups. SNPs showed relatively uniform frequency distribution trend in sub populations belonging to same super population groups EUR and EAS. However, differences can be seen amongst the AFR, AMR and SAS sub population groups.

### **Figure S6. UCSC Genome browser screenshot depicting proxy SNPs and LD structure in the ACE2 region**

UCSC genome browser showing  $r^2$  plot along with other annotations for SNPs rs1978124 and rs2106809 from the intron 1 of ACE2 gene. GeneHancer regulator elements and interaction between GeneHancer regulatory elements as curved lines can be seen. In addition, it can be seen that all the proxy SNPs (in Brown) are clustered together within potential functional regulatory enhancer element and interaction region (in brown color).

### **Figure S7. Expression of TMPRSS2 in different tissues**

(a) Expression of TMPRSS2 in different tissues as retrieved from GTEx portal. Transcripts per million (TPM) are plotted on Y axis and different tissues on X axis. (b) RNA and Protein expression in different tissues is depicted in Human Protein Atlas online portal. In lungs, only RNA expression of TMPRSS2 gene is noted with no protein expression depiction. (c) Expression levels as read count in shades of purple from low to high values indicating expression of different transcripts in different tissues in GTEx portal. Gene models of alternate transcripts are also depicted.

### **Figure S8. UCSC Genome browser screenshot depicting Lung eQTLs along with SNPs and LD structure in the TMPRSS2 region**

Lung eQTLs that overlapped with the common dbSNPs are shown and along  $r^2$  values depicting LD in the region is plotted. All the SNPs appeared to cluster towards 3' end of the gene. Other annotations include GeneHancer regulator elements and interaction between GeneHancer regulatory elements as curved lines can be seen. (dark blue color for *TMPRSS2* gene). Alternate splicing graph is also indicated.

### **Figure S9. Elucidation of key functional SNPs of *TMPRSS2* and their Frequency distribution in 1000G population groups**

(a) *TMPRSS2* gene eQTLs in lungs appear to cluster in the gene towards 3' end and potentially are associated with expression of alternative transcripts in lungs. (b) Amongst these common eQTLs, HaploReg annotations indicated rs4818239 as an important SNP. (c) Frequency distribution on world map of rs4818239

(d) rs75603675 (p.Gly8Val) (e) rs12329760 (p.Val197Met)

The frequencies of SNPs are based on 1000G populations dataset and depicted on world map by LDpop tool of web based LDlink 4.0.3 suite from National Cancer Institute, USA. Derived allele frequencies from low to high are plotted as shades of white to blue.

### **Figure S10. Sequence alignment and structure of *TMPRSS2***

(a) The pairwise sequence alignment shows the N-terminal difference in two isoforms and highlights two mutation regions. (b) The predicted 3D structure of *TMPRSS2* protein. Two domains of *TMPRSS2* and mutation residue are highlighted in the carton model with different colors.

### **Figure S11. *In silico* protein-protein docking analysis of *TMPRSS2*.**

(A1) Surface Model (A2) Carton Model of *TMPRSS2* (Pink) human protein interaction with ACE2 (Green) of SARS CoV-2. (B1) Surface Model (B2) Carton Model of Docked complex of SARS CoV-2 S-protein protein (Blue) with *TMPRSS2* (Pink). (C) Comparison of docking score between wild-type and mutant *TMPRSS2* with ACE2 and SARS CoV-2 S-protein.

### **Figure S12. Summarised information about *SLC6A19* gene and its variants**

(a) Expression of *TMPRSS2* in different tissues as retrieved from GTEx portal indicates very restricted expression of the gene. Expression levels as read count in shades of purple from low to high values. Gene models of alternate transcripts are also depicted. (b) Expression levels as median read count per base as shades of blue from low to high values indicating exons expression in different tissues in GTEx portal. Gene models of alternate transcripts are also depicted. (c) *SLC6A19* RNA and Protein expression in different tissues is depicted in Human Protein Atlas online portal indicating *SLC6A19* has very restricted expression in tissues and in lungs neither its RNA nor Protein expresses. (d) UCSC genome browser screenshot highlights a prominent Enhancer element GH05J001199 at 5'UTR of the gene overlapping exon 1 and extending in intron 1 of the gene. (e) Screenshot from UCSC genome browser depicting SNPs cluster overlapping and

intersecting enhancer region GH05J001199. It also depicts common variants (with applied filter on track for SNPs with frequency greater than 20 percent in global populations) in dbSNP version 151. (f) eQTLs in different tissues for *SLC6A19* were plotted through GTEx Locus Browser. Size of the dot indicate level of significance (as negative p values) whereas colour depicts positive or negative correlation with Normalized effect size (NES) of the eQTL from -1 to 0 to 1. Red color shades represent upregulation whereas blue color shades show downregulation. It also depicts Linkage disequilibrium (LD) in region with value range from 0 to 1 as white to black shades with 1(dark) as absolute LD. Significant upregulation of the gene by eQTLs was observed in Pancreas whereas downregulation in Liver and no data on Lungs as it is reported not to be expressing in Lungs at GTEx portal.

**Figure S13. Number of diseases associated with hub proteins.**

(a) The network view where nodes represent diseases and hub proteins and edges the association between them. (All the nodes of diseases rectangle (red), diseases highly associated with hub proteins rectangle (light green) and hub proteins circle (cyan) are filled color and edges in lines (grey).

(b) The pie chart graph showing the percentage distribution of each hub proteins associated with diseases.

**Figure S14. Topological characteristics of the HPIN and hub removal.**

(a) The figure illustrates the network properties, such as in-degree, out-degree, clustering coefficient, neighbourhood connectivity in the HPIN. The In-degree ( $P(k)_{in}$ ) and out-degree ( $P(k)_{out}$ ) distribution, Clustering Co-efficient  $C(k)$ , Neighbourhood connectivity  $C_N(k)$  is fitted to the power law distribution with exponent values ( $\alpha$ ,  $\beta$ ,  $\gamma$  &  $\mu$ ).

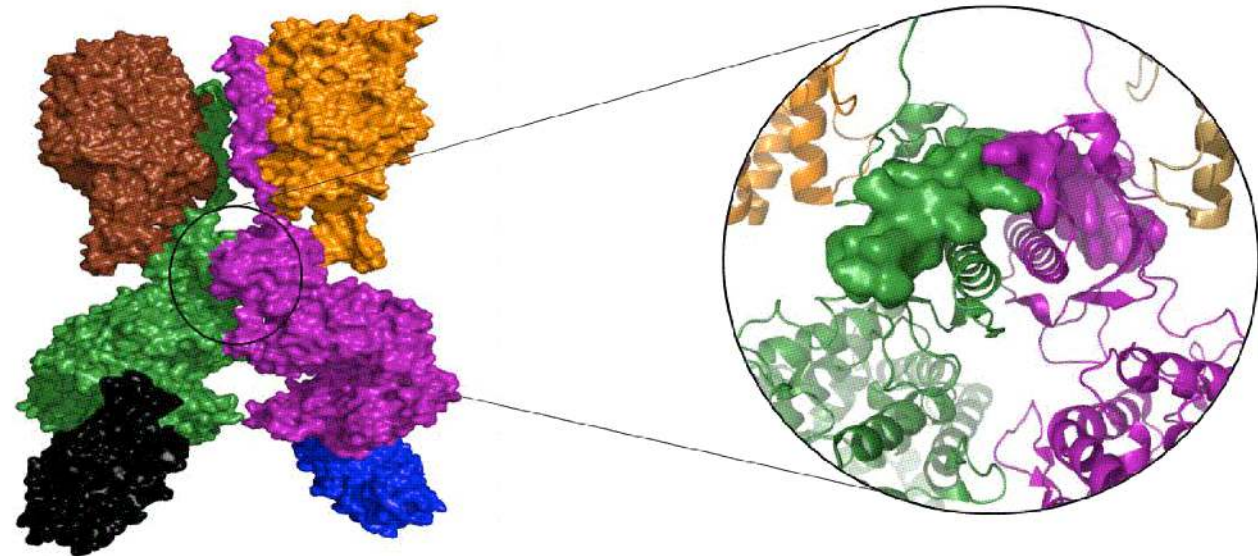
(b) Power law distribution with exponent values  $\alpha$ ,  $\beta$ ,  $\gamma$  &  $\mu$  of HPIN (complete), hub removed proteins network (RPS6, NACA, HNRNPA1 and BTF3).

**A**

**BOAT1**

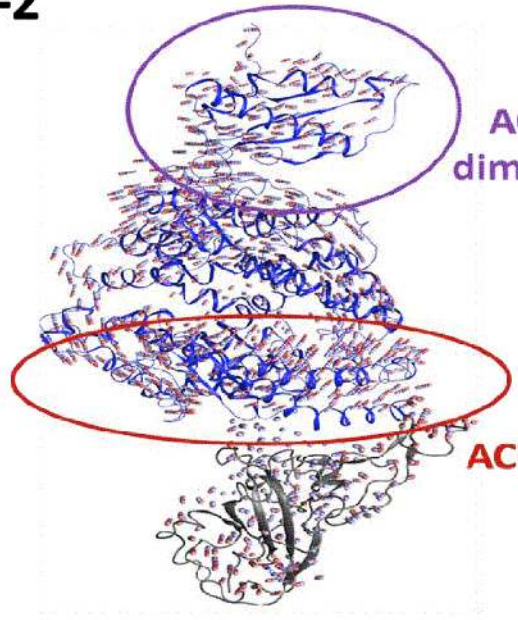
**ACE2**

**RBD  
SARS-CoV-2**



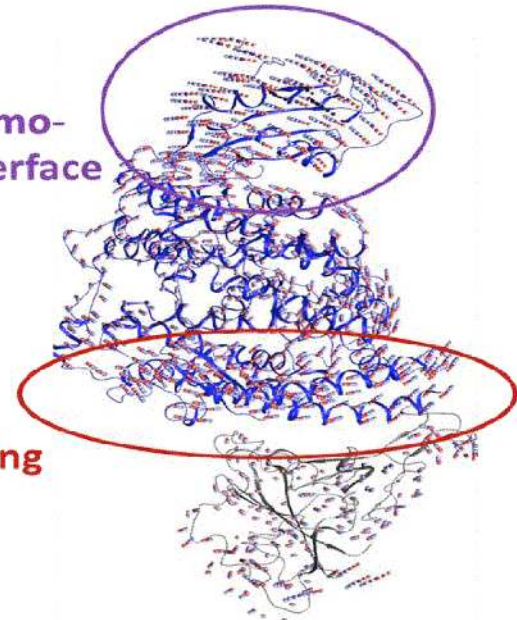
**Residues 697-716 of ACE2 (TMPRSS2 recognition site)**

**B**



**ACE2-RBD from dimeric ACE2 simulation**

**C**



**ACE2-RBD from monomeric ACE2 simulation**

**Figure S1**



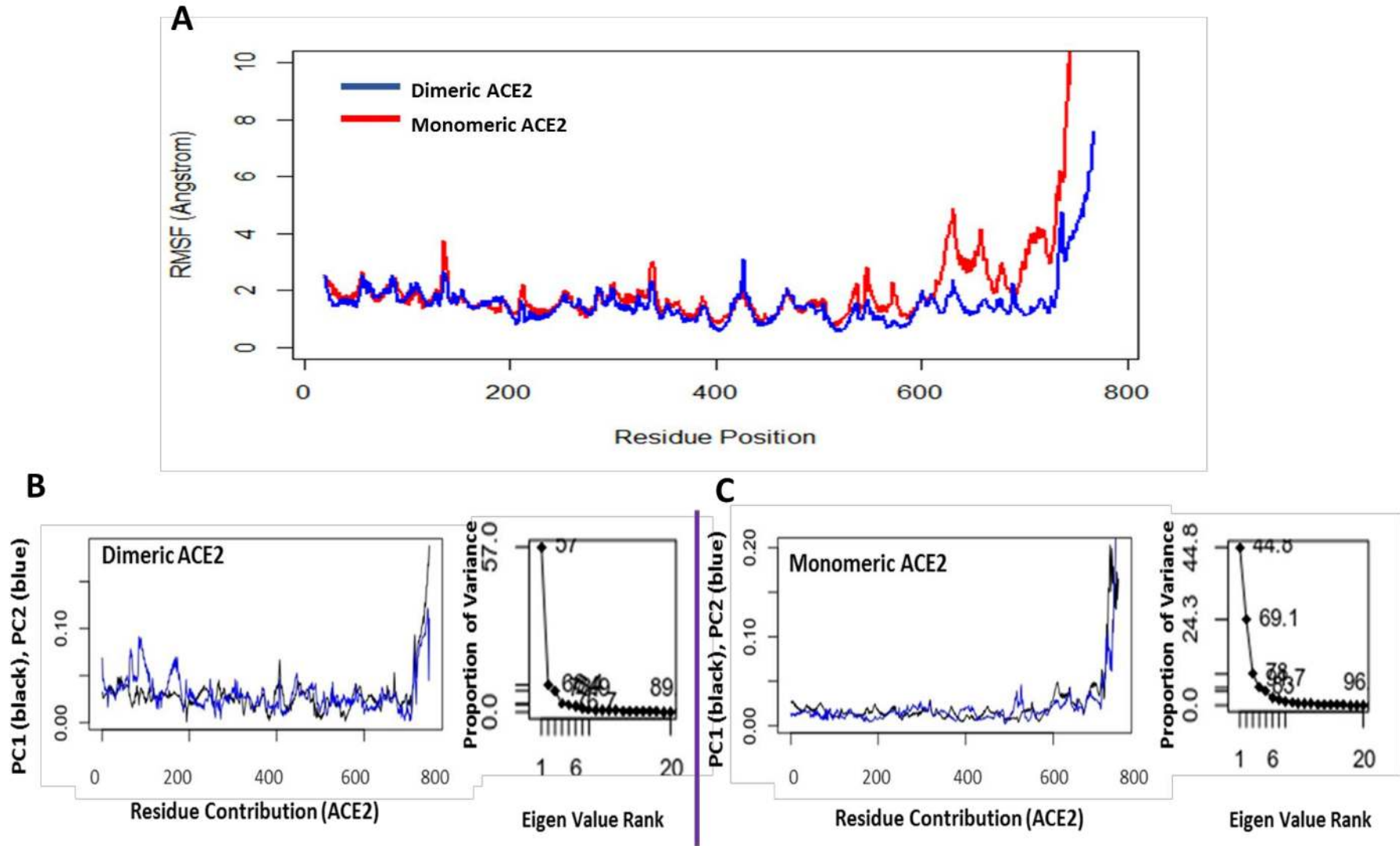


Figure S1''



TISSUE ATLAS

PRIMARY DATA

GENE/PROTEIN

ANTIBODIES AND VALIDATION



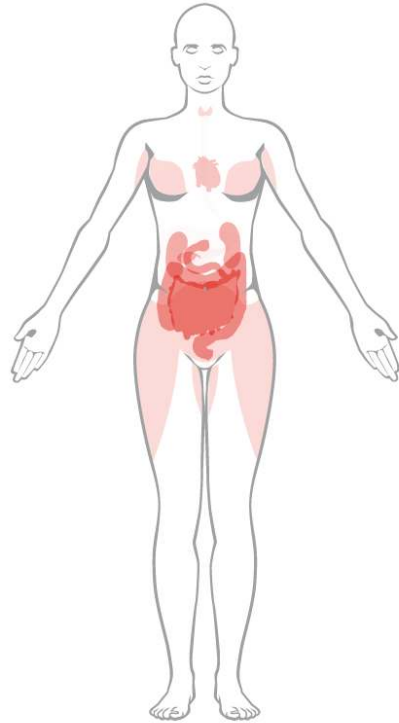
Dictionary



Tissue proteome



RNA AND PROTEIN EXPRESSION SUMMARY<sup>1</sup>



Expression Detection All organs



	RNA expression (NX) <sup>1</sup>	Protein expression (score) <sup>1</sup>
Brain		
Eye		
Endocrine tissues		
Lung		
Proximal digestive tract		
Gastrointestinal tract		
Liver & gallbladder		
Pancreas		
Kidney & urinary bladder		

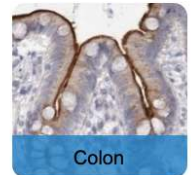


Figure S3

bioRxiv preprint doi: <https://doi.org/10.1101/2020.04.24.050534>; this version posted May 19, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

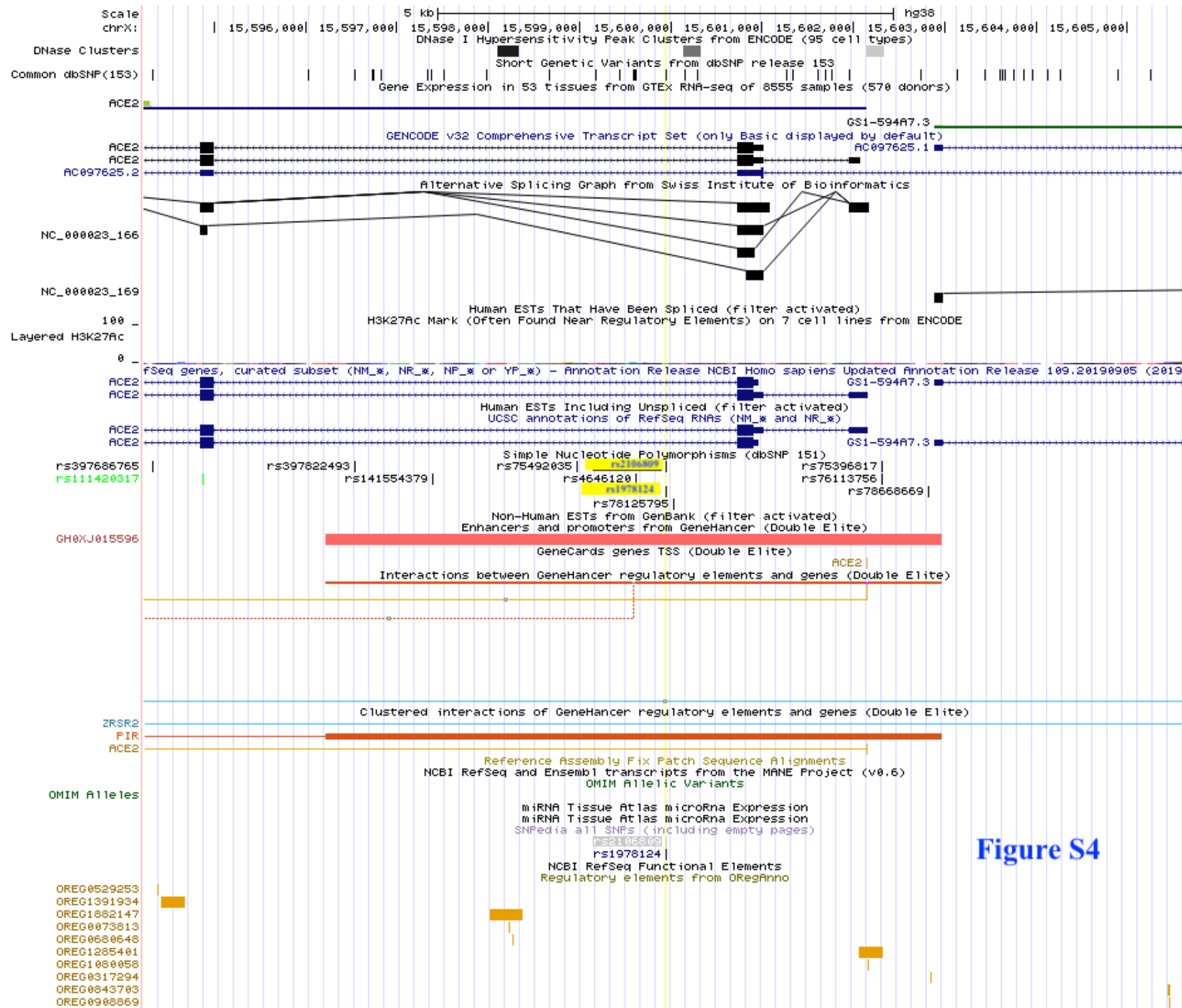
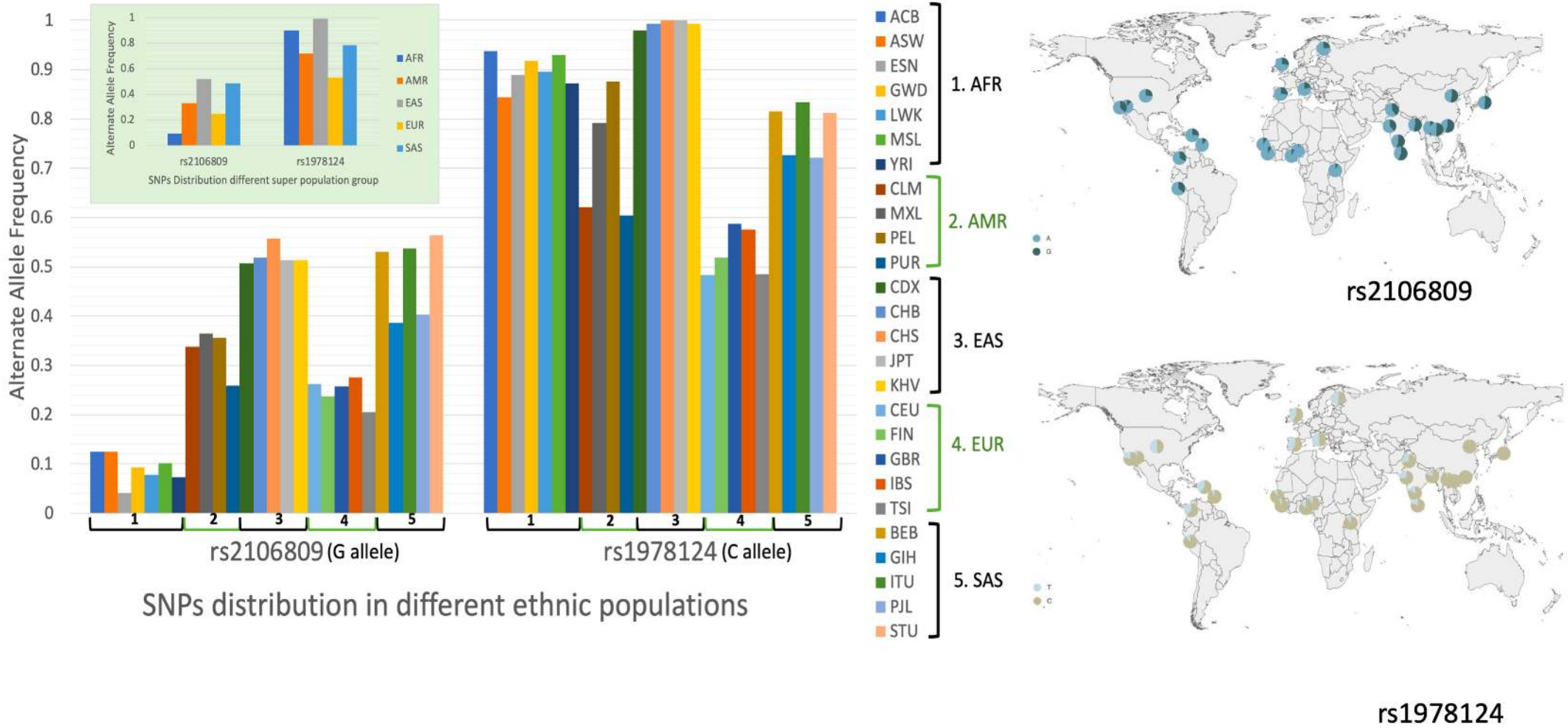


Figure S4



**Figure S5**

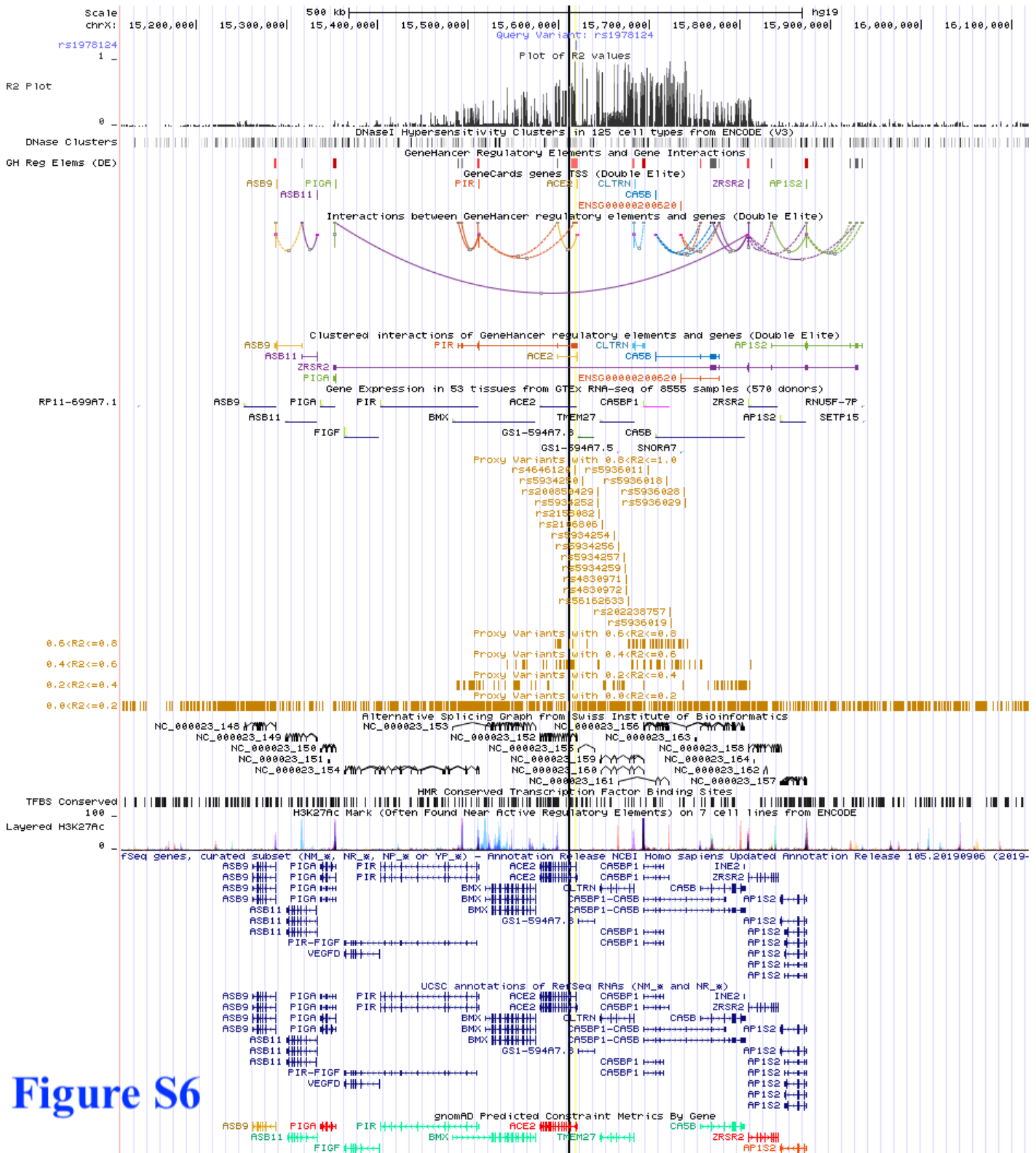
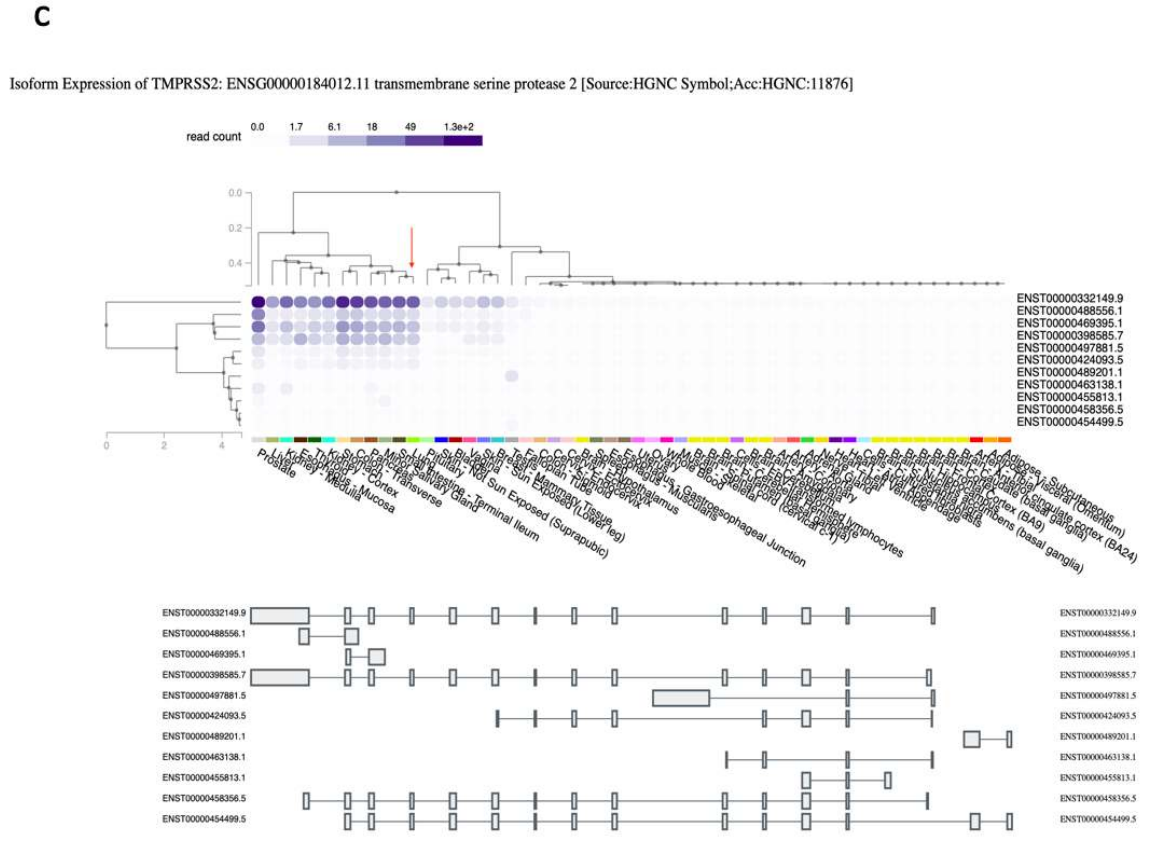
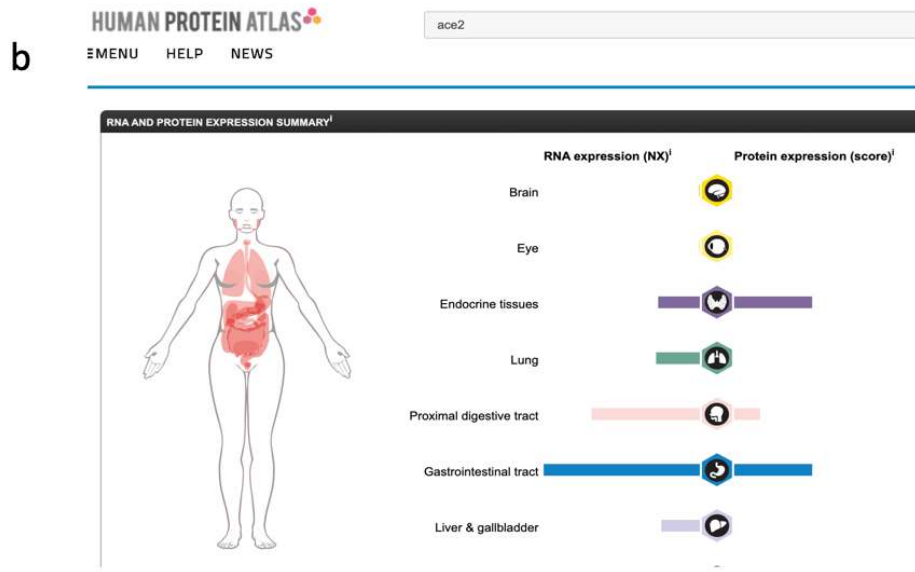
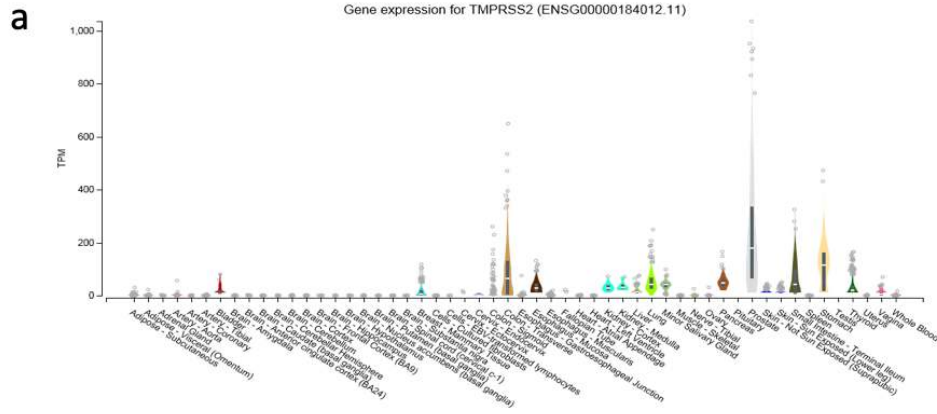
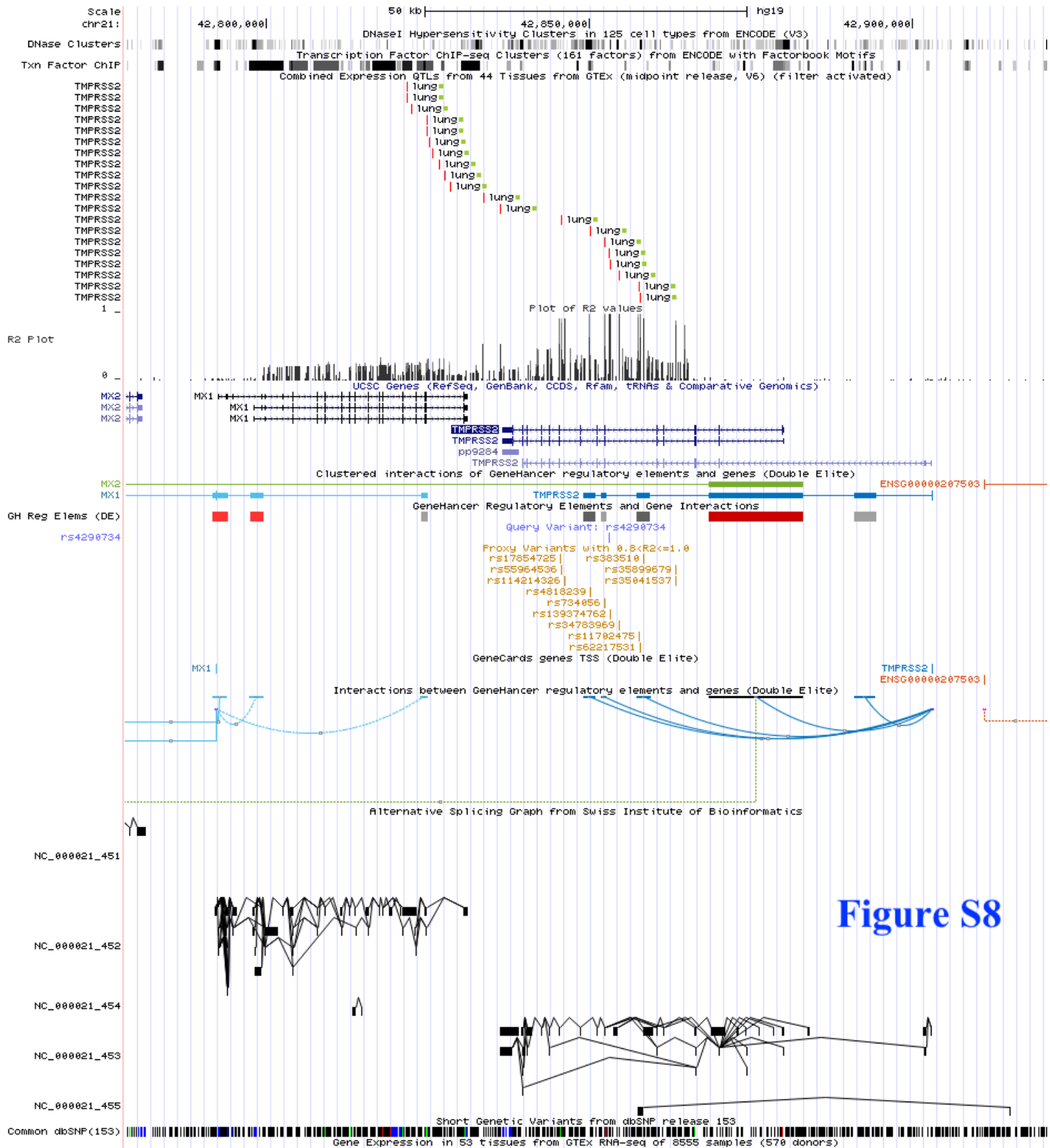


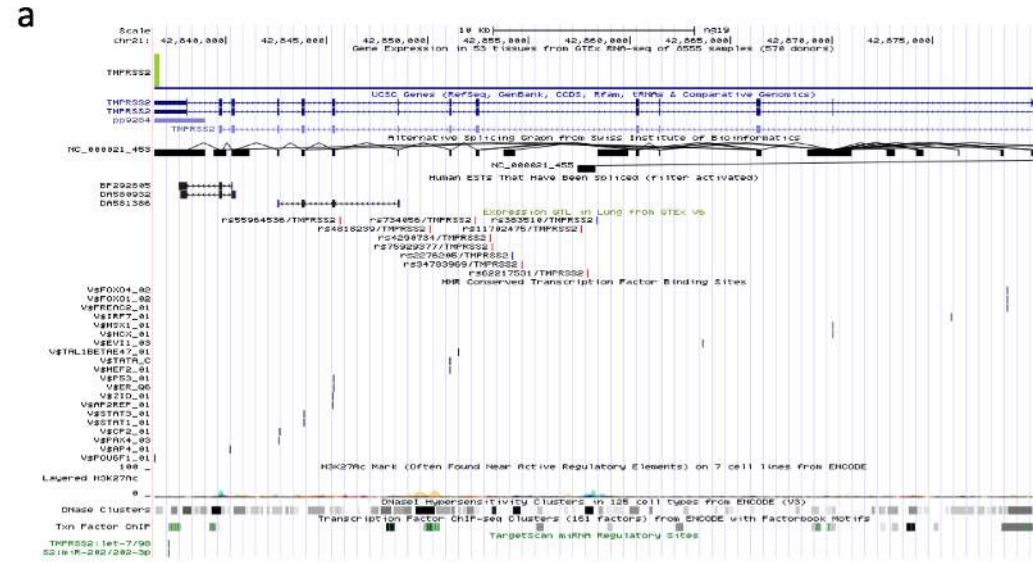
Figure S6



**Figure S7**



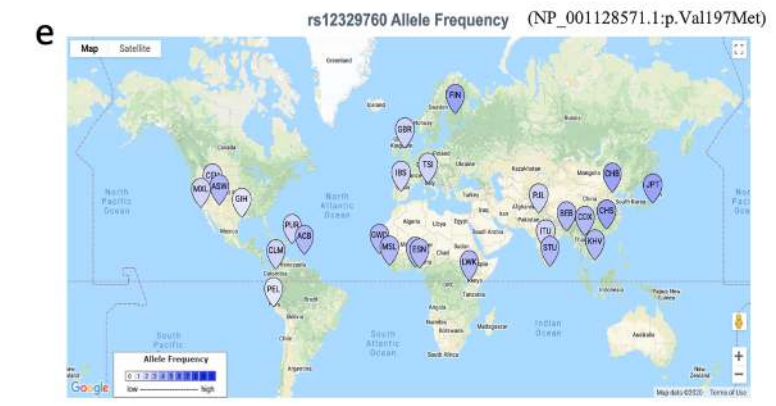
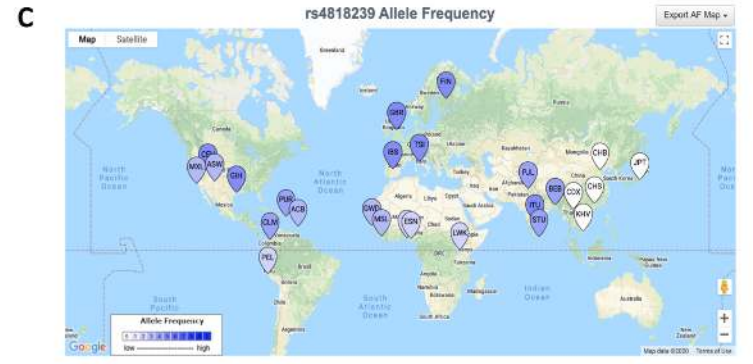




**b**

Query SNP: **rs4818239** and variants with  $r^2 \geq 0.8$

chr	pos (hg38)	LD (r)	LD (D)	variant	Ref	Alt	AFR freq	AMR freq	ASN freq	EUR freq	SIphy cons	Promoter histone marks	Enhancer histone marks	DNAse	Proteins bound	Motifs changed	NHGRI/VEBI GWAS hits	GRASP hits	QTL hits	Selected eQTL hits	GENCODE genes	dbSNP func annot
21	41478148	1	1	<b>rs4818239</b>	T	C	0.25	0.35	0.01	0.48		<b>H3K4, H3K27</b>	<b>21 tissues</b>	<b>18 tissues</b>	<b>9 bound proteins</b>	<b>4 altered motifs</b>			<b>1 hit</b>		TMPS22	intronic
21	41480393	0.93	0.97	<b>rs734056</b>	C	A	0.22	0.33	0.00	0.48		<b>8 tissues</b>	<b>PANC, GI</b>		<b>Gli1</b>			<b>3 hits</b>		TMPS22	intronic	
21	41481156	0.93	0.97	<b>rs4290734</b>	A	G	0.09	0.32	0.00	0.48		<b>4 tissues</b>	<b>ESDR, BRN</b>		<b>5 altered motifs</b>			<b>3 hits</b>		TMPS22	intronic	
21	41481287	0.93	0.97	<b>rs138374782</b>	T	C	0.13	0.32	0.00	0.48		<b>4 tissues</b>	<b>5 tissues</b>		<b>5 altered motifs</b>			<b>1 hit</b>		TMPS22	intronic	
21	41482711	0.92	0.97	<b>rs34783969</b>	A	T	0.14	0.33	0.00	0.48		<b>8 tissues</b>	<b>6 tissues</b>		<b>5 altered motifs</b>			<b>1 hit</b>		TMPS22	intronic	
21	41485974	0.8	0.93	<b>rs6212531</b>	C	T	0.31	0.33	0.00	0.46		<b>19 tissues</b>	<b>8 tissues</b>		<b>MZF1, THAP1, YY1</b>			<b>1 hit</b>		TMPS22	intronic	

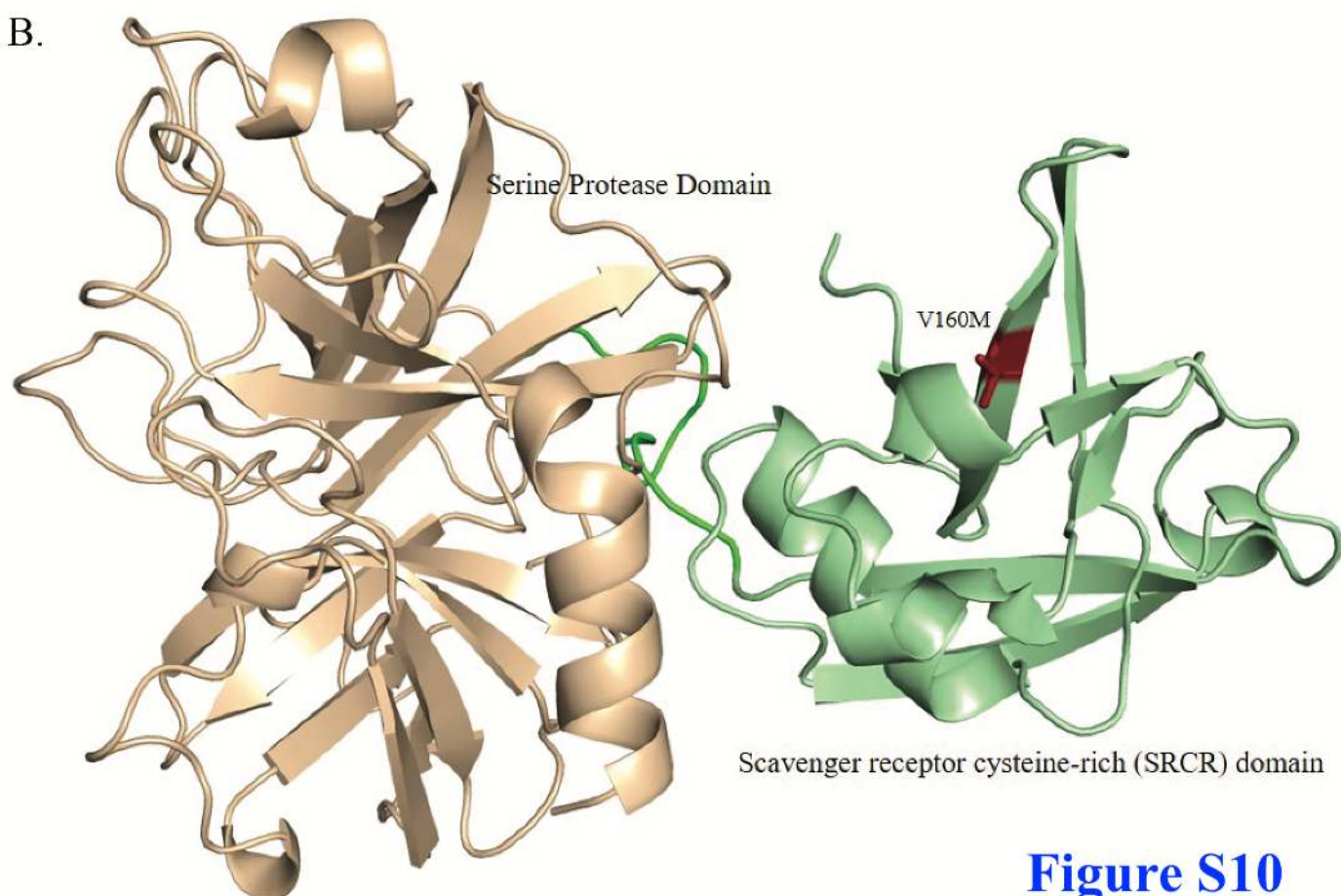


**Figure S9**

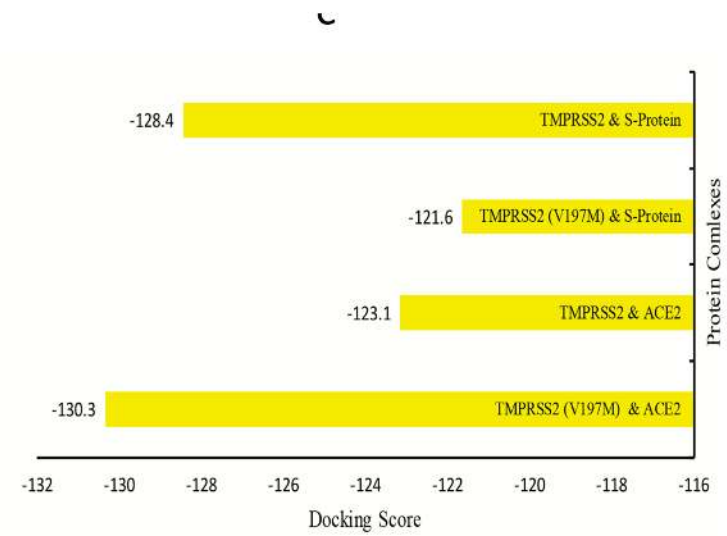
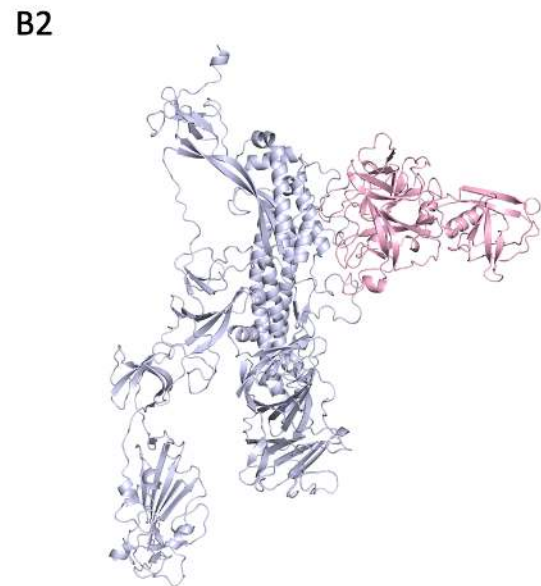
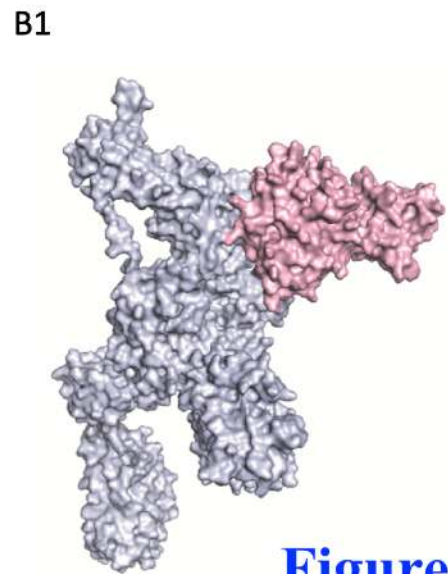
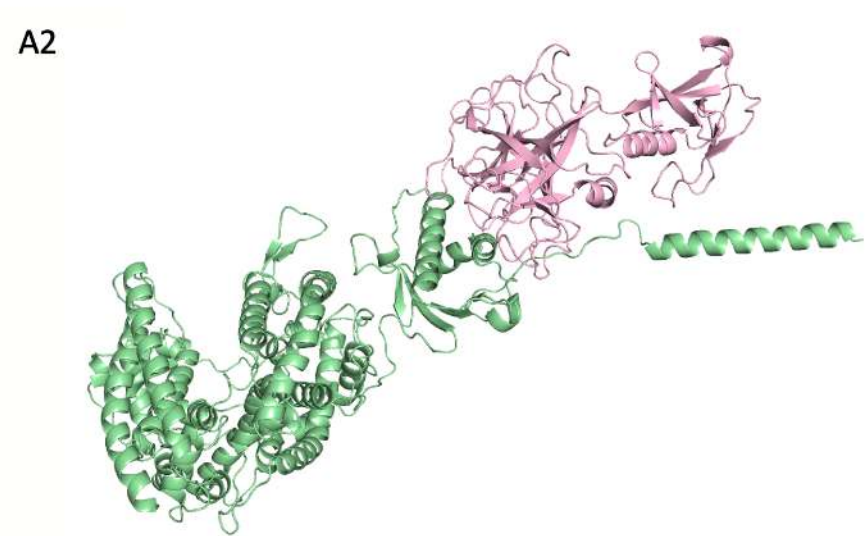
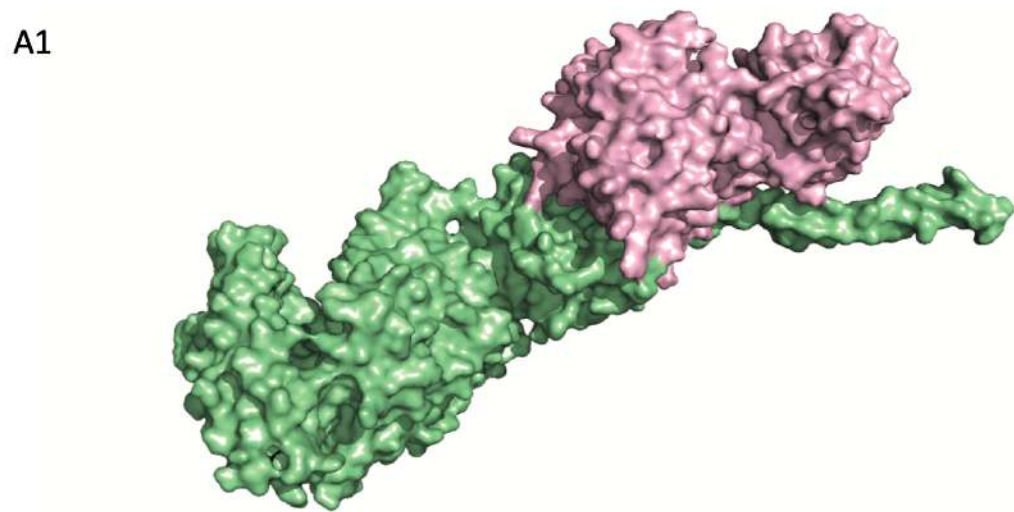
A.

TMPS2-2	MPPAPPGE SGCEERGAAGHIEHSRYLSLLDAVDNSK	60
TMPS2-1	-----MALNSGSPPAIGPYENHGYQPE	23
	*****	
TMPS2-2	NPYPAQPTVVPTVVEVHPAQYYPSPVPQYAPRVL TQASNPVVCTQPKSPSGTVCTSKTKK	120
TMPS2-1	NPYPAQPTVVPTVVEVHPAQYYPSPVPQYAPRVL TQASNPVVCTQPKSPSGTVCTSKTKK	83
	*****	
TMPS2-2	ALCITLTLGTFLVGAALAAAGLLWKFMGSKCSNSGIECDSSGTCINPSNWC DGVSHCPGGE	180
TMPS2-1	ALCITLTLGTFLVGAALAAAGLLWKFMGSKCSNSGIECDSSGTCINPSNWC DGVSHCPGGE	143
	*****	
	V160M ← V197M	
TMPS2-2	DENRCVRLYGPNFILQVYSSQRKSWHPVCQDDWNE NYGRAACRDMGYKNNFYSSQGIVDD	240
TMPS2-1	DENRCVRLYGPNFILQVYSSQRKSWHPVCQDDWNE NYGRAACRDMGYKNNFYSSQGIVDD	203
	*****	
TMPS2-2	SGSTSFMKLNTSAGNVDIYKKLYHSDACSSKAVVSLRCIACGVNLS SRQSRIVGGESAL	300
TMPS2-1	SGSTSFMKLNTSAGNVDIYKKLYHSDACSSKAVVSLRCIACGVNLS SRQSRIVGGESAL	263
	*****	
TMPS2-2	PGAWPWQVSLHVQNVHVC GGSIIITPEWIVTAAHCVEKPLNPNPWHWTAFAGILRQSFMYG	360
TMPS2-1	PGAWPWQVSLHVQNVHVC GGSIIITPEWIVTAAHCVEKPLNPNPWHWTAFAGILRQSFMYG	323
	*****	
TMPS2-2	AGYQVEKVISHPNYDSKTKNNDIALMKLQKPLTFNDLVKPVCLPNPGMMLQPEQLCWI SG	420
TMPS2-1	AGYQVEKVISHPNYDSKTKNNDIALMKLQKPLTFNDLVKPVCLPNPGMMLQPEQLCWI SG	383
	*****	
TMPS2-2	WGATEEKGKTSEVLNAAKVLLIETQRCNSRYVYDNLITPAMICAGFLQGNVDSCQGDSSG	480
TMPS2-1	WGATEEKGKTSEVLNAAKVLLIETQRCNSRYVYDNLITPAMICAGFLQGNVDSCQGDSSG	443
	*****	
TMPS2-2	PLVTSKNNIWWLIGDTSWGS GCAKAYRPGVYGNMVF TDWIYRQMRADG	529
TMPS2-1	PLVTSKNNIWWLIGDTSWGS GCAKAYRPGVYGNMVF TDWIYRQMRADG	492
	*****	

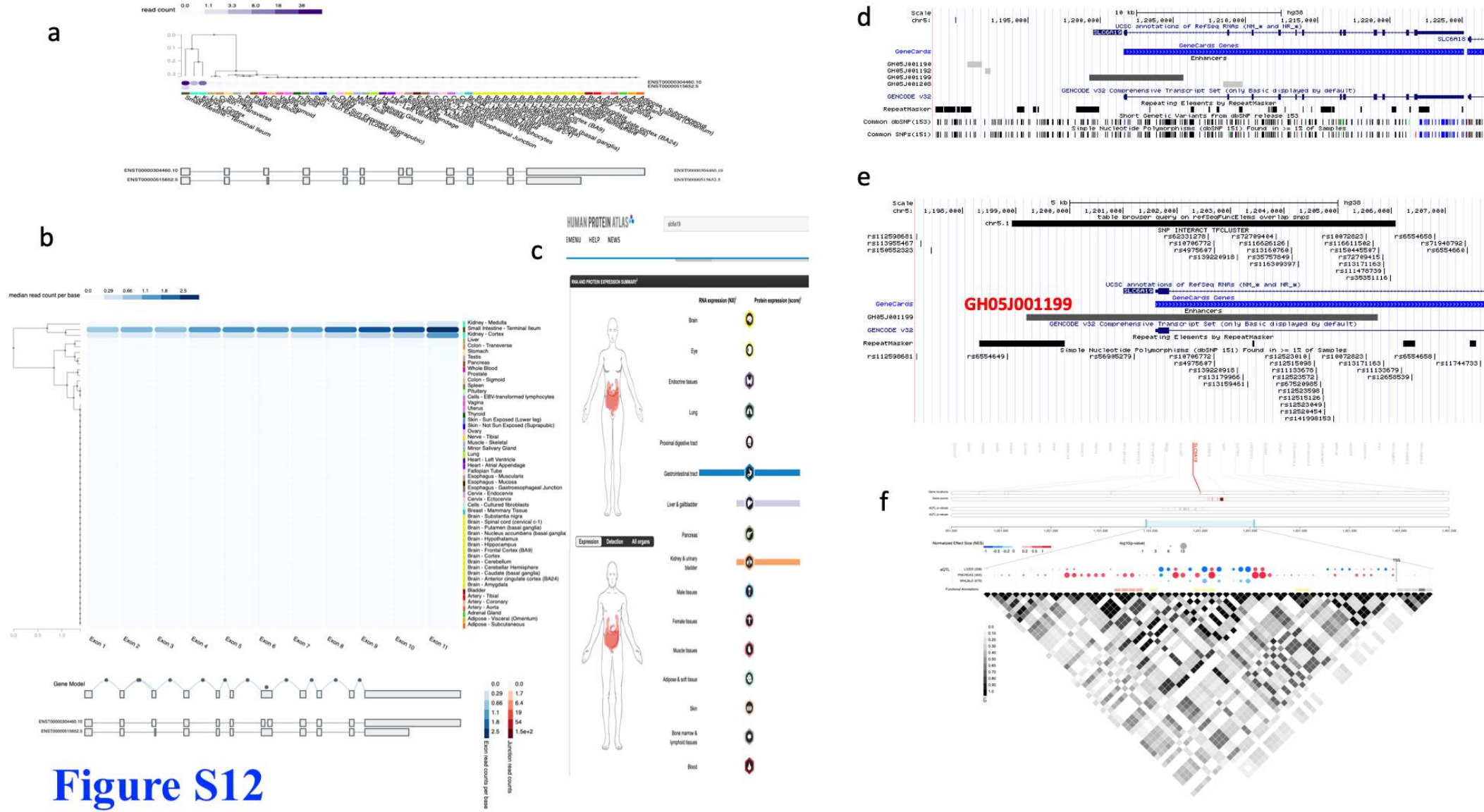
B.



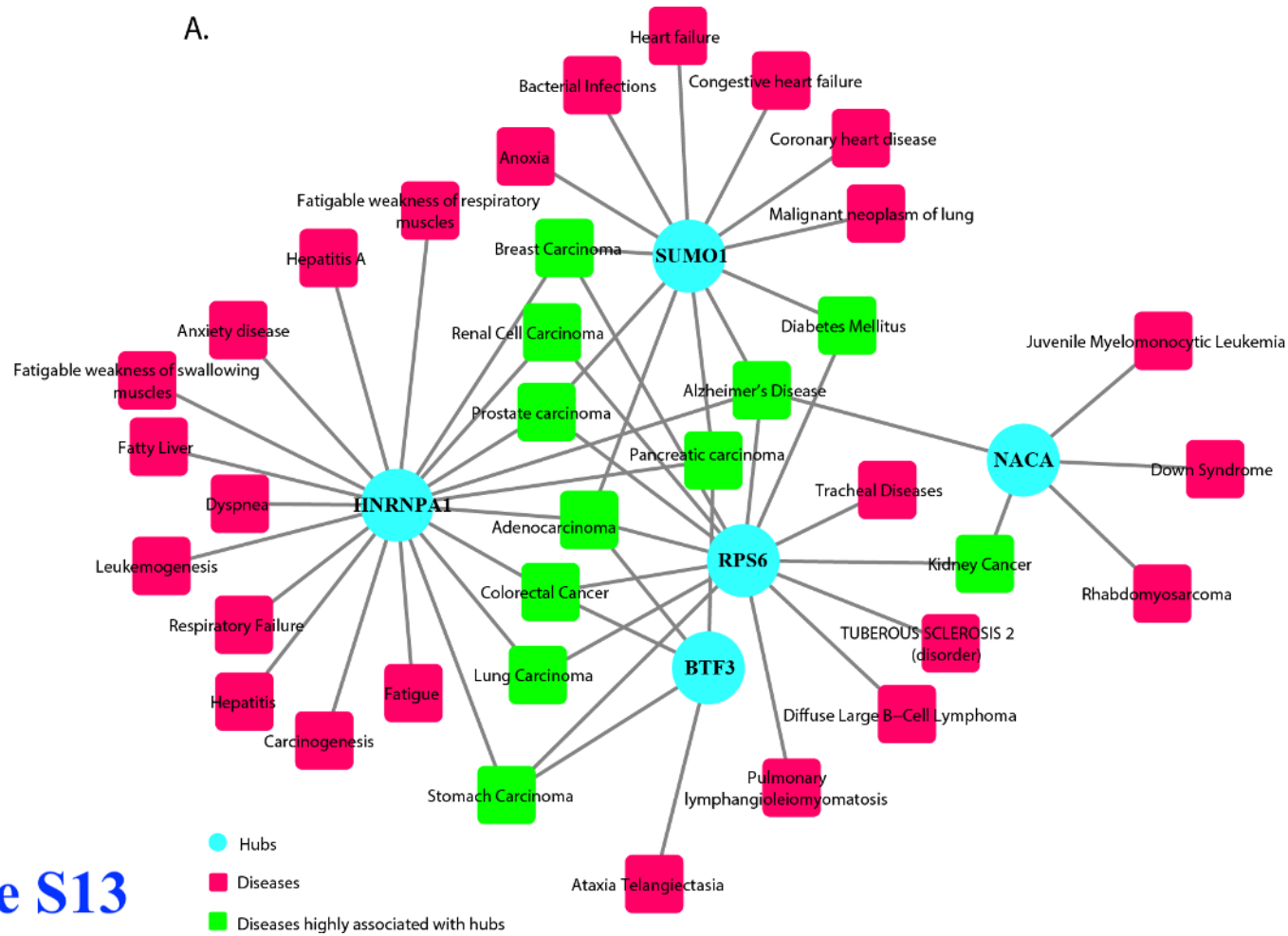
**Figure S10**



**Figure S11**



**Figure S12**



**Figure S13**

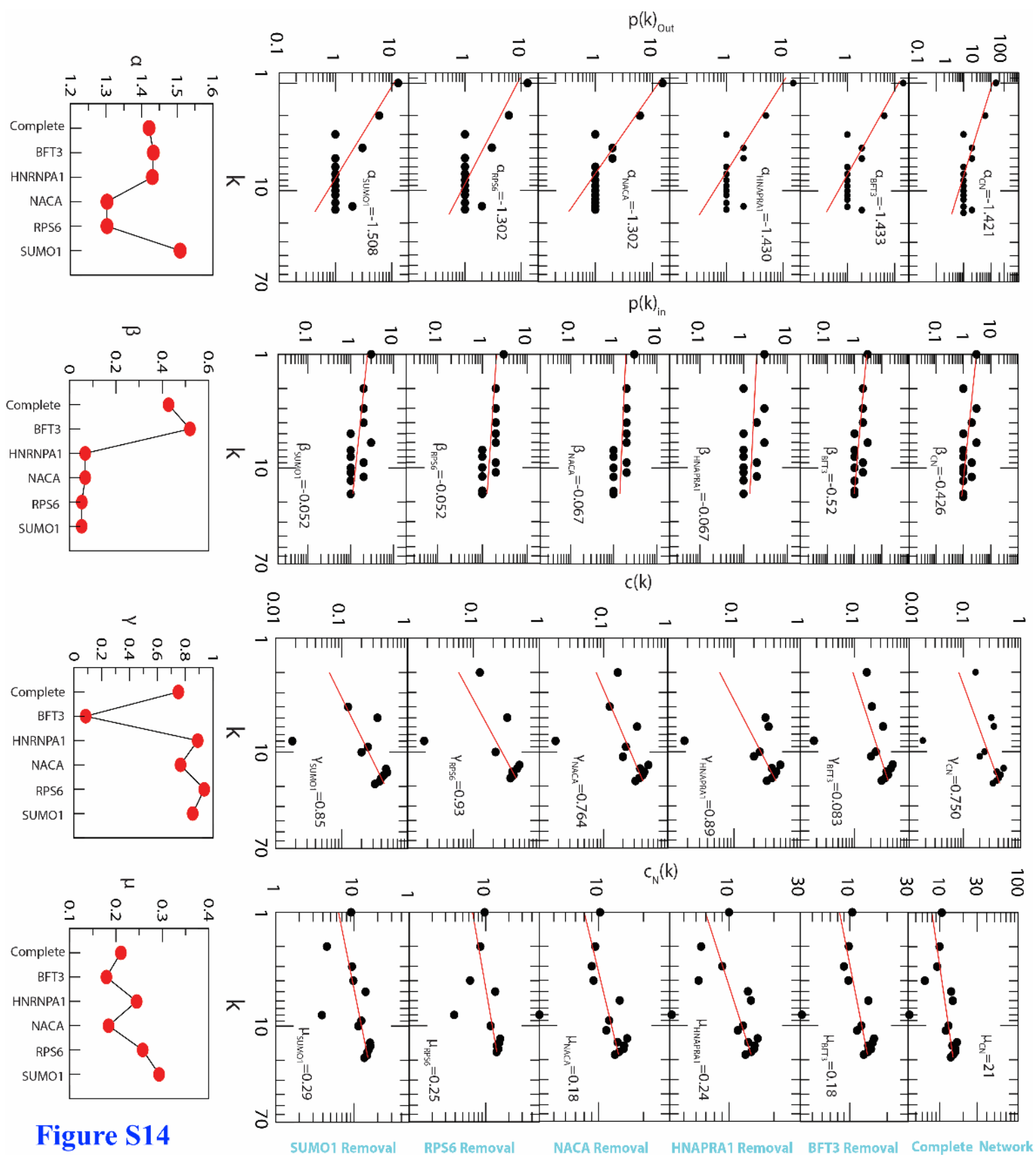


Figure S14

SUMO1 Removal RPS6 Removal NACA Removal HNRNPA1 Removal BFT3 Removal Complete Network