

1 The transmembrane serine protease inhibitors are potential antiviral drugs for
2 2019-nCoV targeting the insertion sequence-induced viral infectivity
3 enhancement

4 Tong Meng^{1,2,10}†, Hao Cao^{3,4}†, Hao Zhang^{5,10}†, Zijian Kang^{6,10}, Da Xu^{7,10}, Haiyi
 5 Gong^{5,10}, Jing Wang⁸, Zifu Li⁸, Xingang Cui⁷, Huji Xu^{4,6}, Haifeng Wei⁵, Xiuwu Pan⁷,
 6 Rongrong Zhu⁹, Jianru Xiao^{5*}, Wang Zhou^{3,10*}, Liming Cheng^{1*}, Jianmin Liu^{8*}.

7 1 Division of Spine, Department of Orthopedics, Tongji Hospital affiliated to Tongji
 8 University School of Medicine, 200065 Shanghai, China

9 2 Tongji University Cancer Center, School of Medicine, Tongji University, 200092
 10 Shanghai, China

11 3 School of Life Science and Biopharmaceutics, Shenyang Pharmaceutical University,
 12 103 Wenhua Road, Shenyang, 110016, PR China

13 4 Peking-Tsinghua Center for Life Sciences, Tsinghua University, Beijing, P.R. China

14 5 Department of Orthopaedic Oncology, Changzheng Hospital, Second Military
 15 Medical University, 200003 Shanghai, China

16 6 Department of Rheumatology and Immunology, Changzheng Hospital, Second
 17 Military Medical University, 200003 Shanghai, China

18 7 Department of Urology, The Third Affiliated Hospital of Second Military Medical
 19 University, 201805 Shanghai, China

20 8 Department of Neurosurgery, Changzheng Hospital, Second Military Medical
 21 University, 200003 Shanghai, China

22 9 Key Laboratory of Spine and Spinal Cord Injury Repair and Regeneration of
 23 Ministry of Education, Orthopaedic Department of Tongji Hospital, School of Life
 24 Science and Technology, Tongji University, 200092 Shanghai, China

25 10 Qiu-Jiang Bioinformatics Institute, 200003 Shanghai, China

26 †These authors contributed equally to this work, and all should be considered first
 27 author.

28 *Correspondence to: chstroke@163.com (Jian-Min Liu)

1 limingcheng@tongji.edu.cn (Li-Ming Cheng)

2 brilliant212@163.com (Wang Zhou)

3 jianruxiao83@163.com (Jian-Ru Xiao)

4 **Abstract**

5 In December 2019, 2019 novel coronavirus (2019-nCoV) induced an ongoing
6 outbreak of pneumonia in Wuhan, Hubei, China. It enters into host cell via cellular
7 receptor recognition and membrane fusion. The former is based on angiotensin
8 converting enzyme II (ACE2). In the latter process, type II transmembrane serine
9 proteases (TTSPs) play important roles in spike protein cleavage and activation. In
10 this study, we used the single-cell transcriptomes of normal human lung and
11 gastroenteric system to identify the ACE2- and TTSP-coexpressing cell composition
12 and proportion. The results revealed that TMPRSS2 was highly co-expressed with
13 ACE2 in the absorptive enterocytes, upper epithelial cells of esophagus and lung AT2
14 cells, implying the important role of TMPRSS2 in 2019-nCoV infection. Additionally,
15 sequence and structural alignment showed that 675-QTQTNSPRRARSVAS-679 was
16 the key sequence mediating 2019-nCoV spike protein, and there was a inserted
17 sequence (680-SPRR-683). We speculated that this insertion sequence especially the
18 exposed structure at R682 and R683 may enhance the recognition and cleavage
19 activity of TMPRSS2 and then increase its viral infectivity. In conclusion, this study
20 provides the bioinformatics and structure evidence for the increased viral infectivity
21 of 2019-nCoV and indicates transmembrane serine protease inhibitors as the antiviral
22 treatment options for 2019-nCoV infection targeting TMPRSS2.

23 **Introduction**

24 At the end of 2019, a rising number of pneumonia patients with unknown pathogen
25 emerged from Wuhan to nearly the entire China[1]. A novel coronavirus was isolated
26 from the human airway epithelial cells and named 2019 novel coronavirus
27 (2019-nCoV)[2]. By analyzing complete genome sequences, 2019-nCoV has an 86.9%
28 nucleotide sequence identity to a severe acute respiratory syndrome (SARS)-like

1 coronavirus detected in bats (bat-SL-CoVZC45, MG772933.1) and is suggested to be
2 the species of SARS related coronaviruses (SARSr-CoV) by pairwise protein
3 sequence analysis[2].

4 As a human coronavirus (HCoV), its prerequisite of infection is entering host cells.
5 During this process, the spike (S) glycoprotein of HCoV plays an important role[3].
6 The surface unit (S1) of S protein mediates the entry into host cells by binding to
7 cellular receptor and the transmembrane unit (S2) subunit harbors the functional
8 elements for the fusion of viral and cellular membranes [4]. It should be noted that the
9 cleavage of S protein is priming to be activated which is essential for viral infectivity
10 [5]. Thus, this process not only needs virus-binding components (cell receptors), but
11 also requires virus protein-cleaving components (cell proteases)[6].

12 In 2019-nCoV infection, a metalloproteinase, angiotensin converting enzyme II (ACE2)
13 is proved to be the cellular receptor, same as SARS-CoV infection[7, 8]. As the host
14 cell protease of ACE2, the type II transmembrane serine proteases (TTSPs), such as
15 TMPRSS2 and TMPRSS11D, can cleave and activate the SARS-CoV S protein
16 (SARS-S) for membrane fusion and also cleave ACE2 to promote viral uptake [9, 10].

17 Thus, we suppose that TTSPs may also play a significant role in 2019-nCoV
18 infection.

19 In order to identify the ACE2- and TTSP-coexpressing cell composition and
20 proportion, we used the single-cell transcriptomes of normal human lung and
21 gastroenteric system based on the public databases. TMPRSS2 was highly
22 co-expressed with ACE2 in absorptive enterocytes, upper epithelial cells of esophagus
23 and lung AT2 cells. In addition, we also explore the cleavage effects of TMPRSS2 in
24 2019-nCoV based on its structure. A striking finding was that the insertion sequence
25 in 2019-nCoV S protein may increase the cleavage activity of TMPRSS2. Thus, as a
26 previously reported antiviral target, transmembrane serine protease inhibitors may
27 also be used for 2019-nCoV treatment targeting TMPRSS2.

28 **Materials and methods**

Data Sources

Single cell transcriptome data were obtained from Single Cell Portal (https://singlecell.broadinstitute.org/single_cell) , Human Cell Atlas Data Portal. (<https://data.humancellatlas.org>) and Gene Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/>). Esophageal data were obtained from the research published by E Madisson et al containing 6 esophageal and 5 lung tissue samples[11]. Three lung datasets were obtained from GSE130148 [12], GSE122960[13] and GSE128169[14], including four, five and eight lung tissues respectively. GSE134520 included 6 gastric mucosal samples from 3 non-atrophic gastritis and 2 chronic atrophic gastritis patients[15]. GSE134809 comprises and 11 noninflammatory ileal samples from Crohn's disease patients[16]. The data from Christopher S *et al* included 12 normal colon samples[17].

Quality Control

Cells would be identified as poor-quality ones if they met one of the following thresholds: (1) The number of expressed genes fewer than 200 or greater than 5000. (2) more than 20% of UMIs were mapped to mitochondrial or ribosomal genes.

Data Integration, Dimension Reduction and Cell Clustering

We perform Different methods to process data for according to the downloaded data.

Esophagus dataset: Rdata were obtained from the research published by E. Madisson *et al*[11]. Dimension Reduction and clustering had already been implemented by the authors.

Lung, stomach and ileum datasets: We utilized functions in the Seurat package to normalize and scale the single-cell gene expression data[18].

Unique molecular identifier (UMI) counts were normalized by the total number of UMIs per cell, multiplied by 10000 for normalization and log-transformed using the "NormalizeData" function. Then, multiple sample data within each dataset were merged using the "FindIntegrationAnchors" and "IntegrateData" functions. After identifying highly variable genes (HVGs) using the "FindVariableGenes" function a

1 principal component analysis (PCA) was performed on the single-cell expression
2 matrix using the “RunPCA” function.
3 We next utilized the “FindClusters” function in the Seurat package to conduct the cell
4 clustering analysis into a graph structure in PCA space after constructing a
5 K-nearest-neighbor graph based on the Euclidean distance in PCA space. Uniform
6 Manifold Approximation and Projection (UMAP) visualization was performed for
7 obtaining the clusters of cells.

8 *Colon Dataset* The single cell data was processed with the R packages LIGER[19]
9 and Seurat[18]. The gene matrix was first normalized to remove differences in
10 sequencing depth and capture efficiency among cells. variable genes in each dataset
11 were identified using the “selectGenes” function. Then we used the “optimizeALS”
12 function in LIGER to perform the integrative nonnegative matrix factorization and
13 select a k of 15 and lambda of 5.0 to obtain a plot of expected alignment. The
14 “quantileAlignSNF” function was then performed to build a shared factor
15 neighborhood graph to jointly cluster cells, then quantile normalizes corresponding
16 clusters. Next nonlinear dimensionality reduction was calculated using the
17 “RunUMAP” function and the results were visualized with UMAP.

18 **Identification of cell types and Gene expression analysis**

19 Clusters were annotated on the expression of known cell markers and the clustering
20 information provided in the articles. Then, we utilized the “RunALRA” function to
21 impute lost values in the gene expression matrix. The imputed gene expression was
22 shown in Feature plots and violin plots. we used “Quantile normalization” in the R
23 package preprocessCore (R package version 1.46.0.

24 <https://github.com/bmbolstad/preprocessCore>) to remove unwanted technical
25 variability across different datasets were further denoised to compare gene expression.

26 **Endocytosis and exocytosis-associated genes signature**

1 All pathways related to endocytosis or exocytosis were obtained from Harmonizome
2 dataset [20]. To detect the expression levels of functional genesets, mean expression
3 of an inflammation signature was calculated.

4 **External validation**

5 To minimize bias, external databases of Genotype-Tissue Expression (GTEx)[21], and
6 The Human Protein Atlas[22] were used to detect gene and protein expression of
7 ACE2 at the tissue levels including normal lung and digestive system, such as
8 esophagus, stomach, small intestine and colon.

9 **Modelling**

10 The structures of 2019-nCoV Spike protein and TMPRSS2 were generated by using
11 SWISS-MODEL online server[23]. The structures were marked, superimposed and
12 visualized by using Chimera[24].

13 **Results**

14 **Annotation of cell types**

15 In this study, 5 datasets with single-cell transcriptomes of esophagus, gastric, small
16 intestine and colon were analyzed, along with lung. Based on Cell Ranger output, the
17 gene expression count matrices were used to present sequential clustering of cells
18 according to different organs or particular clusters. The cell type identity in each
19 cluster was annotated by the expression of the known cell type markers.

20 **Cell type-specific expression of ACE2, TMPRSS2, TMPRSS11D and ADAM17**

21 After initial quality controls, 57,020 cells and 15 clusters were identified in the lung
22 (Fig. 1A). The detected cell types included ciliated, alveolar type 1 (AT1) and
23 alveolar type 2 (AT2) cells, along with fibroblast, muscle, and endothelial cells. The
24 identified immune cell types were T, B and NK cells, along with macrophages,
25 monocytes and dendritic cells (DC). ACE2 was mainly expressed in AT2 cells along
26 with AT1 and fibroblast cells, while TMPRSS2 was found in AT1 and AT2 cells; and
27 TMPRSS11D was expressed in AT1 cells, fibroblast and macrophage (Fig. 1B). The
28 Violin plots showed that ACE2 and TMPRSS2 were co-expressed in AT1 and AT2

1 cells. TMPRSS11D was not found in any clusters, whereas ADAM17 was found in all
2 clusters (Fig. 1C). The Immunohistochemistry (IHC) images of ACE2, TMPRSS2 and
3 TMPRSS11D in normal lung showed a similar result (Fig. 1D).

4 In the esophagus, 87,947 cells passed quality control and 14 cell types were identified.
5 Over 90% cells belong to four major epithelial cell types: upper, stratified, suprabasal,
6 and dividing cells of the suprabasal layer. ACE2 was highly expressed in upper and
7 stratified epithelial cells, while TMPRSS2 and TMPRSS11D were found in upper
8 epithelial cells. ADAM17 was found in almost all clusters (Fig. 2A).

9 A total of 29,678 cells and 10 cell types were identified in the stomach after quality
10 control with a high proportion of gastric epithelial cells, including antral basal gland
11 mucous cells (GMCs), pit mucous cells (PMCs), chief cells and enteroendocrine cells.
12 The expression of ACE2 and TMPRSS11D are relatively low in all the clusters, while
13 TMPRSS2 was found in GMCs and PMCs. ADAM17 was also found in all clusters
14 (Fig. 2B).

15 After quality controls, 11,218 cells and 5 cell types were identified in the ileum
16 epithelia (Fig. 2C). ACE2 and TMPRSS2 were highly expressed in absorptive
17 enterocytes. In the meantime, TMPRSS2 was also lowly expressed in goblet and
18 Paneth cells. TMPRSS11D was not found in any cluster, whereas ADAM17 was
19 found in all clusters except undifferentiated cells (Fig. 2C).

20 All the 47,442 cells from the colon were annotated after quality controls. Absorptive
21 and secretory clusters were identified in epithelial cells. The absorptive epithelial cells
22 included transit amplifying (TA) cells (TA 1, TA 2), immature enterocytes, and
23 enterocytes. The secretory epithelial cells comprised progenitor cells (secretory TA,
24 immature goblet) and for mature cells (goblet, and enteroendocrine). ACE2 was
25 mainly found in enterocytes and less expressed in immature enterocytes, so was
26 TMPRSS2. Other clusters also had a lower expression of TMPRSS2. TMPRSS11D
27 was not found in any cluster and ADAM17 was expressed in enteroendocrine cells
28 (Fig. 2D).

1 The expressions of ACE2, classic TTSPs (TMPRSS2, TMPRSS3, TMPRSS4,
2 TMPRSS6, TMPRSS11D, TMPRSS14 and ADAM17 were detected in lung and
3 digestive tract clusters. An almost consistent expression and distribution was only
4 found between ACE2 and TMPRSS2 in all the 9 clusters, with high expression in
5 absorptive enterocytes, upper epithelial cells of esophagus and lung AT2 cells (Fig.
6 3A). The endocytosis and exocytosis-associated genes which are related to the entry
7 of virus into host cells were also detected in all the 9 clusters. The endocytosis
8 signature were more expressed in colon and exocytosis signature were more
9 expressed in upper epithelial cells of esophagus (Fig. 3B). The RNA-seq data of lung,
10 esophagus, stomach, small intestine, colon-transverse and colon-sigmoid were
11 obtained from GTEx database. The expressions of ACE2 and TMPRSS2 also had a
12 similar tendency and were highly expressed in small intestine and colon, while the
13 TMPRSS11D was mainly found in the esophagus (Fig. 3C).

14 **The structure of the SARS-S and 2019-nCoV S protein homo-trimers**

15 The structure of the SARS-S and 2019-nCoV S protein were compared. The insert
16 aa675-690 to SARS-S aa661-672 with the structural missed residues are colored
17 green (Fig. 4A). Although the 2019-nCoV S protein has similar structure with
18 SARS-S, there was a very obvious insertion sequence (QTQTNSPRRARSVASQS) in
19 2019-nCoV S protein. The insert aa675-690 of 2019-nCoV S protein that corresponds
20 to the insert region of SARS-S is colored yellow (Fig. 4B). The structural
21 superimpose of SARS-S (wheat) and 2019-nCoV S protein (cyan) (Fig. 4B).
22 The obvious insertion sequence in 2019-nCoV S protein is at residue R685, especially
23 at R682 and R683 (the orange frame in Fig. 5A and the yellow structure in Fig. 5B).
24 In addition, the insertion sequence of R682 and R683 was protruded from the
25 molecular surface (Fig. 5B, C).

26 **Structure and catalytic mechanism of TMPRSS2**

27 The catalytic triad comprised of H296, D345 and S441 are colored blue, green and
28 cyan, respectively, the substrate binding residue D435 which located in the bottom of

pocket is marked in red, the substrate binding pocket is deeper than most of serine proteinase (Fig. 6A, B). The bottom of TMPRSS2 has a negatively charged aspartic acid residue which can facilitate the binding and stabilization of the long-chain positively charged amino acid residue of substrate. Polypeptide substrate analogue KQLR was presented in Fig. 6C, with arginine, glutamine, leucine and lysine. The Fig. 6D and 6E revealed the state of substrate analogue binding to the catalytic pocket (Fig. 6D, E).

Discussion

The coronaviruses is the common infection source of enteric, respiratory, and central nervous system in humans and other mammals[25]. At the beginning of the twenty-first century, two betacoronaviruses, SARS-CoV and MERS-CoV, result in persistent public panics and became the most significant public health events[26]. In December 2019, a novel identified coronavirus (2019-nCoV) induced an ongoing outbreak of pneumonia in Wuhan, Hubei, China [27]. Till now, the pathogenic mechanism of 2019-nCoV is still unclear. In this study, we found the high co-expression between ACE2 and TMPRSS2 in the absorptive enterocytes, upper epithelial cells of esophagus and lung AT2 cells. Additionally, the insertion sequence (680-SPRR-683) in 2019-nCoV S protein may increase the cleavage activity of TMPRSS2 and enhance the viral infectivity, indicating that TMPRSS2 may serve as a candidate antiviral target for 2019-nCoV infection.

Similar to SARS-CoV and MERS-CoV, 2019-nCoV utilizes ACE2 as a receptor for host cell entry [8]. In human HeLa cells, expressing ACE2 from human, civet, and Chinese horseshoe bat can help many kinds of SARS-CoV enter into the cells including 2019-nCoV [8, 28-30]. During this process, the catalytic domain of ACE2 interacts with a defined receptor-binding domain (RBD) of CTD1 in S protein[7, 31]. The “up” and “down” transition of CTD1 allows ACE2 binding by regulating the relationship among CTD1, CTD2, S1-ACE2 complex and S2 subunit[32].

1 In this study, we found a strong co-expression between ACE2 and TMPRSS2. The
2 latter can cleave SARS-S near or at the cell surface along with another TTSP human
3 airway trypsin-like protease (HAT), also known as TMPRSS11D, and render host cell
4 entry independent of the endosomal pathway using cathepsin B/L[33]. Different from
5 cathepsin B/L, TMPRSS2 can also promote viral spread in the host [34]. In addition,
6 TMPRSS2 and HAT can also cleave ACE2 to augment about 30-fold viral infectivity
7 and during this process, the metalloprotease ADAM17 can compete with TMPRSS2
8 [9]. The detailed process is presented in Fig. 7. Due to the critical role of TMPRSS2
9 in influenza virus and coronavirus infections, serine protease inhibitors have been
10 used in the antiviral therapeutic strategy targeting TMPRSS2 with high antiviral
11 activities, such as camostat, nafamostat and leupeptin[34-37].

12 The TTSPs showed similar substrate-specificity and catalytic mechanism. As a TTSP,
13 TMPRSS2, located on the cell surface, includes extracellular domain, transmembrane
14 domain and intracellular domain structurally, in which extracellular domain is the
15 main catalytic domain. TMPRSS2 can catalyze the hydrolysis of substrate protein
16 with a strong substrate specificity. Moreover, based on previous study, as for TTSP
17 family members, the P1 position of the cleaved substrates is arginine[38, 39]. This is
18 mainly determined by the special structure of its catalytic pocket. The catalytic pocket
19 of TMPRSS2 is relatively deep, and its bottom has a negatively charged aspartic acid
20 residue which can facilitate the binding and stabilization of the long-chain positively
21 charged amino acid residue of substrate.

22 Generally, TMPRSS2 is shown to cleave SARS-S at residue R667 which is
23 significantly associated with SARS-S activation for cell-cell fusion[40]. As the S1/S2
24 cleavage site, R667 is often embedded in different cleavage motifs and cleaved during
25 S protein biogenesis[41]. Besides R667, residue R797 is also required for S protein
26 activation by TMPRSS2. As the S2' cleavage site, it is frequently cleaved during viral
27 entry[41]. Both the hydrolysis sites of R667 and R797 are corresponding to the
28 SARS-S sequence tagged site.

1 Similar to SARS-S, the 2019-nCoV S protein has three arginine hydrolysis sites
2 located on its surface (R685, R682 and R683), with two from the insertion sequence
3 in 2019-nCoV S protein (680-SPRR-683). As we above-mentioned, the main
4 hydrolysis site of TMPRSS2 is the long-chain positively charged amino acid residue
5 represented by R and more R often increase the viral infectivity[38, 39]. The insertion
6 sequence of R682 and R683 was exposed from the molecular surface. Based on the
7 catalytic mechanism of TMPRSS2, we supposed that the additional arginine residues
8 linked structure exposed from protein surface may be more conducive to the
9 recognition and cleavage activity of TMPRSS2, thereby enhancing the viral
10 infectivity of 2019-nCoV.

11 Some researchers believed that there were four inserts in the S protein of 2019-nCoV,
12 which may be the result of artificial modification. The SARS-S has been widely
13 studied, including vaccine development and drug design[42, 43]. The key residues
14 663-VSLLRSTSQ-671 involved in SARS-S protein cleavage site was missed in all
15 structures that have been resolved to date due to its flexibility. The reasonable
16 modification of protein/enzyme function is to modify the key sites of protein/enzyme
17 based on the existing structural studies. We can observe that in the key sequence of S
18 protein splicing, 675-QTQTNSPRRARSVAS-679, only two bases of 685-RS-686 are
19 same as SARS-S, which is most likely not designed for protein function.

20 **Conclusion**

21 ACE2 and TMPRSS2 were highly co-expressed to facilitate the entry of 2019-nCoV
22 into host cells. In addition, the insertion sequence in 2019-nCoV S protein may
23 increase the cleavage activity of TMPRSS2 and enhance the viral infectivity,
24 indicating that transmembrane serine protease inhibitors may serve as the antiviral
25 treatment options for 2019-nCoV infection targeting TMPRSS2.

26 **Reference**

- 27 1. The L. Emerging understandings of 2019-nCoV. Lancet. 2020.
- 28 2. Zhu N, Zhang D, Wang W, Li X, Yang B, Song J, Zhao X, Huang B, Shi W,

- 1 Lu R, Niu P, Zhan F, Ma X, et al. A Novel Coronavirus from Patients with
- 2 Pneumonia in China, 2019. N Engl J Med. 2020.
- 3 3. Walls AC, Xiong X, Park YJ, Tortorici MA, Snijder J, Quispe J, Cameroni E,
- 4 Gopal R, Dai M, Lanzavecchia A, Zambon M, Rey FA, Corti D, et al.
- 5 Unexpected Receptor Functional Mimicry Elucidates Activation of Coronavirus
- 6 Fusion. Cell. 2019; 176: 1026-39.e15.
- 7 4. Hofmann H, Pohlmann S. Cellular entry of the SARS coronavirus. Trends
- 8 Microbiol. 2004; 12: 466-72.
- 9 5. Gallagher TM, Buchmeier MJ. Coronavirus spike proteins in viral entry and
- 10 pathogenesis. Virology. 2001; 279: 371-4.
- 11 6. Shulla A, Heald-Sargent T, Subramanya G, Zhao J, Perlman S, Gallagher
- 12 T. A transmembrane serine protease is linked to the severe acute respiratory
- 13 syndrome coronavirus receptor and activates virus entry. J Virol. 2011; 85:
- 14 873-82.
- 15 7. Gui M, Song W, Zhou H, Xu J, Chen S, Xiang Y, Wang X. Cryo-electron
- 16 microscopy structures of the SARS-CoV spike glycoprotein reveal a
- 17 prerequisite conformational state for receptor binding. Cell Res. 2017; 27:
- 18 119-29.
- 19 8. P Zhou XY, XG Wang, B Hu, L Zhang, W Zhang, HR Si, Y Zhu, B Li, CL
- 20 Huang, HD Chen, J Chen, Y Luo, H Guo, RD Jiang, MQ Liu, Y Chen, XR Shen,
- 21 X Wang, XS Zheng, K Zhao, QJ Chen, F Deng, LL Liu, B Yan, FX Zhan, YY

- 1 Wang, GF Xiao, ZL Shi. A pneumonia outbreak associated with a new
- 2 coronavirus of probable bat origin. *nature*. 2020.
- 3 9. Heurich A, Hofmann-Winkler H, Gierer S, Liepold T, Jahn O, Pohlmann S.
- 4 TMPRSS2 and ADAM17 cleave ACE2 differentially and only proteolysis by
- 5 TMPRSS2 augments entry driven by the severe acute respiratory syndrome
- 6 coronavirus spike protein. *J Virol*. 2014; 88: 1293-307.
- 7 10. Li F. Structure, Function, and Evolution of Coronavirus Spike Proteins.
- 8 *Annu Rev Virol*. 2016; 3: 237-61.
- 9 11. Madisson E, Wilbrey-Clark A, Miragaia RJ, Saeb-Parsy K, Mahbubani KT,
- 10 Georgakopoulos N, Harding P, Polanski K, Huang N, Nowicki-Osuch K,
- 11 Fitzgerald RC, Loudon KW, Ferdinand JR, et al. scRNA-seq assessment of
- 12 the human lung, spleen, and esophagus tissue stability after cold preservation.
- 13 *Genome Biol*. 2019; 21: 1.
- 14 12. Vieira Braga FA, Kar G, Berg M, Carpaij OA, Polanski K, Simon LM,
- 15 Brouwer S, Gomes T, Hesse L, Jiang J, Fasouli ES, Efremova M,
- 16 Vento-Tormo R, et al. A cellular census of human lungs identifies novel cell
- 17 states in health and in asthma. *Nat Med*. 2019; 25: 1153-63.
- 18 13. Reyfman PA, Walter JM, Joshi N, Anekalla KR, McQuattie-Pimentel AC,
- 19 Chiu S, Fernandez R, Akbarpour M, Chen CI, Ren Z, Verma R,
- 20 Abdala-Valencia H, Nam K, et al. Single-Cell Transcriptomic Analysis of
- 21 Human Lung Provides Insights into the Pathobiology of Pulmonary Fibrosis.

- 1 Am J Respir Crit Care Med. 2019; 199: 1517-36.
- 2 14. Valenzi E, Bulik M, Tabib T, Morse C, Sembrat J, Trejo Bittar H, Rojas M,
- 3 Lafyatis R. Single-cell analysis reveals fibroblast heterogeneity and
- 4 myofibroblasts in systemic sclerosis-associated interstitial lung disease. Ann
- 5 Rheum Dis. 2019; 78: 1379-87.
- 6 15. Zhang P, Yang M, Zhang Y, Xiao S, Lai X, Tan A, Du S, Li S. Dissecting
- 7 the Single-Cell Transcriptome Network Underlying Gastric Premalignant
- 8 Lesions and Early Gastric Cancer. Cell Rep. 2019; 27: 1934-47.e5.
- 9 16. Martin JC, Chang C, Boschetti G, Ungaro R, Giri M, Grout JA, Gettler K,
- 10 Chuang LS, Nayar S, Greenstein AJ, Dubinsky M, Walker L, Leader A, et al.
- 11 Single-Cell Analysis of Crohn's Disease Lesions Identifies a Pathogenic
- 12 Cellular Module Associated with Resistance to Anti-TNF Therapy. Cell. 2019;
- 13 178: 1493-508.e20.
- 14 17. Smillie CS, Biton M, Ordovas-Montanes J, Sullivan KM, Burgin G, Graham
- 15 DB, Herbst RH, Rogel N, Slyper M, Waldman J, Sud M, Andrews E, Velonias
- 16 G, et al. Intra- and Inter-cellular Rewiring of the Human Colon during
- 17 Ulcerative Colitis. Cell. 2019; 178: 714-30.e22.
- 18 18. Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, 3rd,
- 19 Hao Y, Stoeckius M, Smibert P, Satija R. Comprehensive Integration of
- 20 Single-Cell Data. Cell. 2019; 177: 1888-902.e21.
- 21 19. Welch JD, Kozareva V, Ferreira A, Vanderburg C, Martin C, Macosko EZ.

- 1 Single-Cell Multi-omic Integration Compares and Contrasts Features of Brain
- 2 Cell Identity. *Cell*. 2019; 177: 1873-87.e17.
- 3 20. Rouillard AD, Gundersen GW, Fernandez NF, Wang Z, Monteiro CD,
- 4 McDermott MG, Ma'ayan A. The harmonizome: a collection of processed
- 5 datasets gathered to serve and mine knowledge about genes and proteins.
- 6 Database (Oxford). 2016; 2016.
- 7 21. Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis:
- 8 multitissue gene regulation in humans. *Science*. 2015; 348: 648-60.
- 9 22. Uhlen M, Fagerberg L, Hallstrom BM, Lindskog C, Oksvold P, Mardinoglu
- 10 A, Sivertsson A, Kampf C, Sjostedt E, Asplund A, Olsson I, Edlund K,
- 11 Lundberg E, et al. Proteomics. Tissue-based map of the human proteome.
- 12 *Science*. 2015; 347: 1260419.
- 13 23. Biasini M, Bienert S, Waterhouse A, Arnold K, Studer G, Schmidt T, Kiefer
- 14 F, Gallo Cassarino T, Bertoni M, Bordoli L, Schwede T. SWISS-MODEL:
- 15 modelling protein tertiary and quaternary structure using evolutionary
- 16 information. *Nucleic Acids Res*. 2014; 42: W252-8.
- 17 24. Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng
- 18 EC, Ferrin TE. UCSF Chimera--a visualization system for exploratory research
- 19 and analysis. *J Comput Chem*. 2004; 25: 1605-12.
- 20 25. Perlman S, Netland J. Coronaviruses post-SARS: update on replication
- 21 and pathogenesis. *Nat Rev Microbiol*. 2009; 7: 439-50.

- 1 26. de Wit E, van Doremalen N, Falzarano D, Munster VJ. SARS and MERS:
2 recent insights into emerging coronaviruses. Nat Rev Microbiol. 2016; 14:
3 523-34.
- 4 27. Huang C, Wang Y, Li X, Ren L, Zhao J, Hu Y, Zhang L, Fan G, Xu J, Gu X,
5 Cheng Z, Yu T, Xia J, et al. Clinical features of patients infected with 2019
6 novel coronavirus in Wuhan, China. Lancet. 2020.
- 7 28. Yang XL, Hu B, Wang B, Wang MN, Zhang Q, Zhang W, Wu LJ, Ge XY,
8 Zhang YZ, Daszak P, Wang LF, Shi ZL. Isolation and Characterization of a
9 Novel Bat Coronavirus Closely Related to the Direct Progenitor of Severe
10 Acute Respiratory Syndrome Coronavirus. J Virol. 2015; 90: 3253-6.
- 11 29. Ge XY, Li JL, Yang XL, Chmura AA, Zhu G, Epstein JH, Mazet JK, Hu B,
12 Zhang W, Peng C, Zhang YJ, Luo CM, Tan B, et al. Isolation and
13 characterization of a bat SARS-like coronavirus that uses the ACE2 receptor.
14 Nature. 2013; 503: 535-8.
- 15 30. Hu B, Zeng LP, Yang XL, Ge XY, Zhang W, Li B, Xie JZ, Shen XR, Zhang
16 YZ, Wang N, Luo DS, Zheng XS, Wang MN, et al. Discovery of a rich gene
17 pool of bat SARS-related coronaviruses provides new insights into the origin of
18 SARS coronavirus. PLoS Pathog. 2017; 13: e1006698.
- 19 31. Li F, Li W, Farzan M, Harrison SC. Structure of SARS coronavirus spike
20 receptor-binding domain complexed with receptor. Science. 2005; 309:
21 1864-8.

- 1 32. Song W, Gui M, Wang X, Xiang Y. Cryo-EM structure of the SARS
2 coronavirus spike glycoprotein in complex with its host cell receptor ACE2.
3 PLoS Pathog. 2018; 14: e1007236.
- 4 33. Shirato K, Kawase M, Matsuyama S. Wild-type human coronaviruses
5 prefer cell-surface TMPRSS2 to endosomal cathepsins for cell entry. Virology.
6 2018; 517: 9-15.
- 7 34. Zhou Y, Vedantham P, Lu K, Agudelo J, Carrion R, Jr., Nunneley JW,
8 Barnard D, Pohlmann S, McKerrow JH, Renslo AR, Simmons G. Protease
9 inhibitors targeting coronavirus and filovirus entry. Antiviral Res. 2015; 116:
10 76-84.
- 11 35. Shen LW, Mao HJ, Wu YL, Tanaka Y, Zhang W. TMPRSS2: A potential
12 target for treatment of influenza virus and coronavirus infections. Biochimie.
13 2017; 142: 1-10.
- 14 36. Shin WJ, Seong BL. Type II transmembrane serine proteases as potential
15 target for anti-influenza drug discovery. Expert Opin Drug Discov. 2017; 12:
16 1139-52.
- 17 37. Yamamoto M, Matsuyama S, Li X, Takeda M, Kawaguchi Y, Inoue JI,
18 Matsuda Z. Identification of Nafamostat as a Potent Inhibitor of Middle East
19 Respiratory Syndrome Coronavirus S Protein-Mediated Membrane Fusion
20 Using the Split-Protein-Based Cell-Cell Fusion Assay. Antimicrob Agents
21 Chemother. 2016; 60: 6532-9.

- 1 38. Herter S, Piper DE, Aaron W, Gabriele T, Cutler G, Cao P, Bhatt AS, Choe
2 Y, Craik CS, Walker N, Meininger D, Hoey T, Austin RJ. Hepatocyte growth
3 factor is a preferred in vitro substrate for human hepsin, a
4 membrane-anchored serine protease implicated in prostate and ovarian
5 cancers. *Biochem J.* 2005; 390: 125-36.
- 6 39. Limburg H, Harbig A, Bestle D, Stein DA, Moulton HM, Jaeger J, Janga H,
7 Harges K, Koepke J, Schulte L, Koczulla AR, Schmeck B, Klenk HD, et al.
8 TMPRSS2 Is the Major Activating Protease of Influenza A Virus in Primary
9 Human Airway Cells and Influenza B Virus in Human Type II Pneumocytes. *J*
10 *Virol.* 2019; 93.
- 11 40. Bertram S, Glowacka I, Muller MA, Lavender H, Gnirss K, Nehlmeier I,
12 Niemeyer D, He Y, Simmons G, Drosten C, Soilleux EJ, Jahn O, Steffen I, et al.
13 Cleavage and activation of the severe acute respiratory syndrome coronavirus
14 spike protein by human airway trypsin-like protease. *J Virol.* 2011; 85:
15 13363-72.
- 16 41. Millet JK, Whittaker GR. Host cell proteases: Critical determinants of
17 coronavirus tropism and pathogenesis. *Virus Res.* 2015; 202: 120-34.
- 18 42. McPherson C, Chubet R, Holtz K, Honda-Okubo Y, Barnard D, Cox M,
19 Petrovsky N. Development of a SARS Coronavirus Vaccine from Recombinant
20 Spike Protein Plus Delta Inulin Adjuvant. *Methods Mol Biol.* 2016; 1403:
21 269-84.

- 1 43. Du L, He Y, Zhou Y, Liu S, Zheng BJ, Jiang S. The spike protein of
- 2 SARS-CoV--a target for vaccine and therapeutic development. Nat Rev
- 3 Microbiol. 2009; 7: 226-36.

4

5

6

7

8

9

10

11

12

13

14

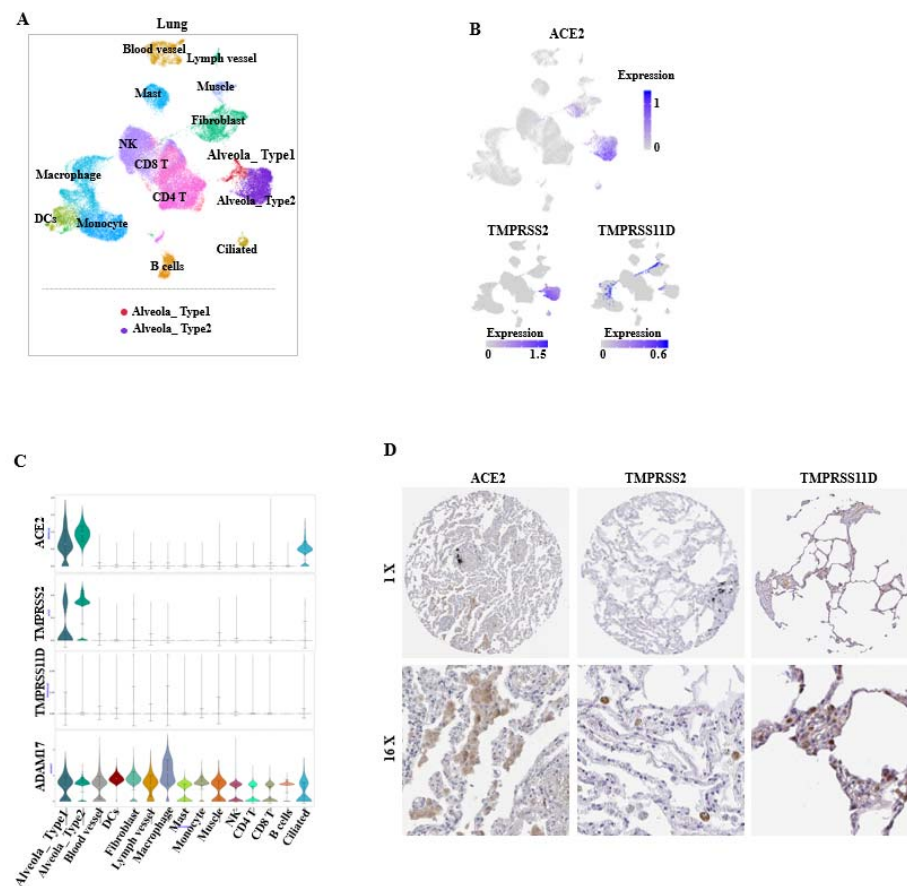


Figure 1. Single-cell analysis and immunohistochemistry of normal lung tissue.

(A). UMAP plots showing the landscape of cells from normal lung tissue. Fifteen clusters are colored, distinctly labeled.

(B). Feature plots demonstrating expression of ACE2, TMPRSS2 and TMPRSS11D across fifteen clusters.

(C). Violin plots showing the expression of ACE2, TMPRSS2, TMPRSS11D and ADAM17 across clusters. The expression is measured as the log2 (TP10K+1) value.

(D). Immunohistochemical images showing the expression of ACE2, TMPRSS2 and TMPRSS11D in lung tissues.

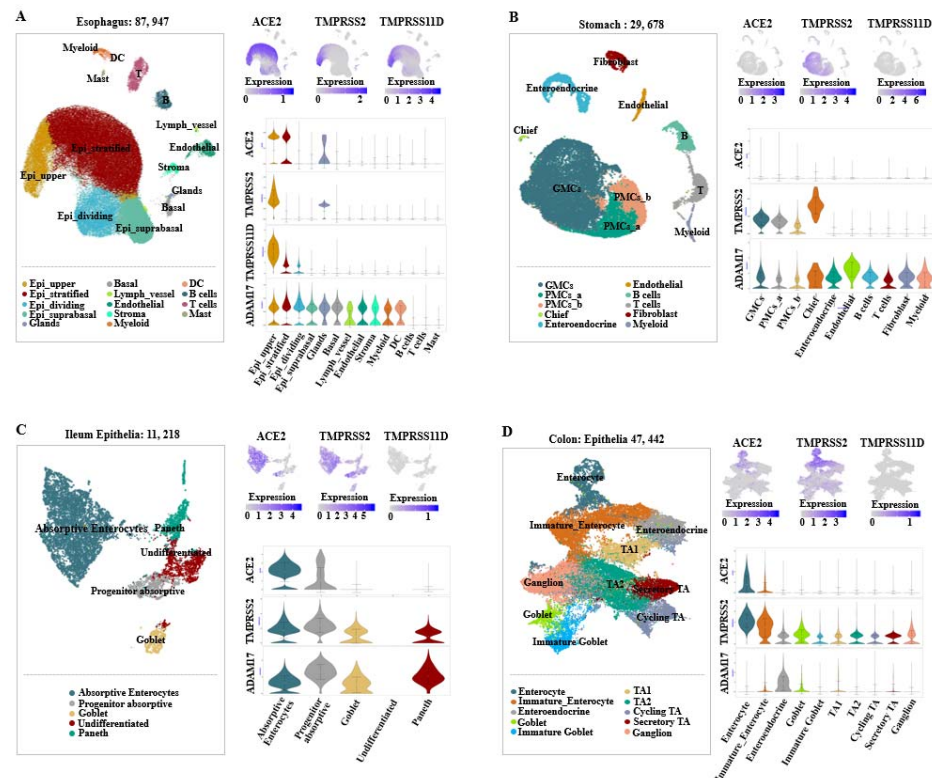


Figure 2. Single-cell analysis of esophageal cells, gastric mucosal cells, ileal epithelial cells and colonic epithelial cells.

(A). UMAP plots showing 87,947 esophageal cells. Fourteen clusters are colored, distinctively labeled. Feature plots demonstrating expression of ACE2, TMPRSS2 and TMPRSS11D across esophageal clusters. Violin plots showing the expression of ACE2, TMPRSS2, TMPRSS11D and ADAM17. The expression is measured as the $\log_2(\text{TP10K}+1)$ value.

(B). UMAP plots showing 29,678 gastric mucosal cells. Ten clusters are colored, distinctively labeled. Feature plots demonstrating expression of ACE2, TMPRSS2 and TMPRSS11D across gastric mucosal clusters. Violin plots showing the expression of ACE2, TMPRSS2, TMPRSS11D and ADAM17. The expression is measured as in (A).

(C). UMAP plots showing 11,218 ileal epithelial cells. Five clusters are colored, distinctively labeled. Feature plots demonstrating expression of ACE2, TMPRSS2 and TMPRSS11D across ileal epithelial clusters. Violin plots showing the expression of

1 ACE2, TMPRSS2, TMPRSS11D and ADAM17. The expression is measured as in
2 (A).
3 (D). UMAP plots showing 47,442 colonic epithelial cells. Ten clusters are colored,
4 distinctively labeled. Feature plots demonstrating expression of ACE2, TMPRSS2 and
5 TMPRSS11D across colonic epithelial clusters. Violin plots showing the expression
6 of ACE2, TMPRSS2, TMPRSS11D and ADAM17. The expression is measured as in
7 (A).

8

9

10

11

12

13

14

15

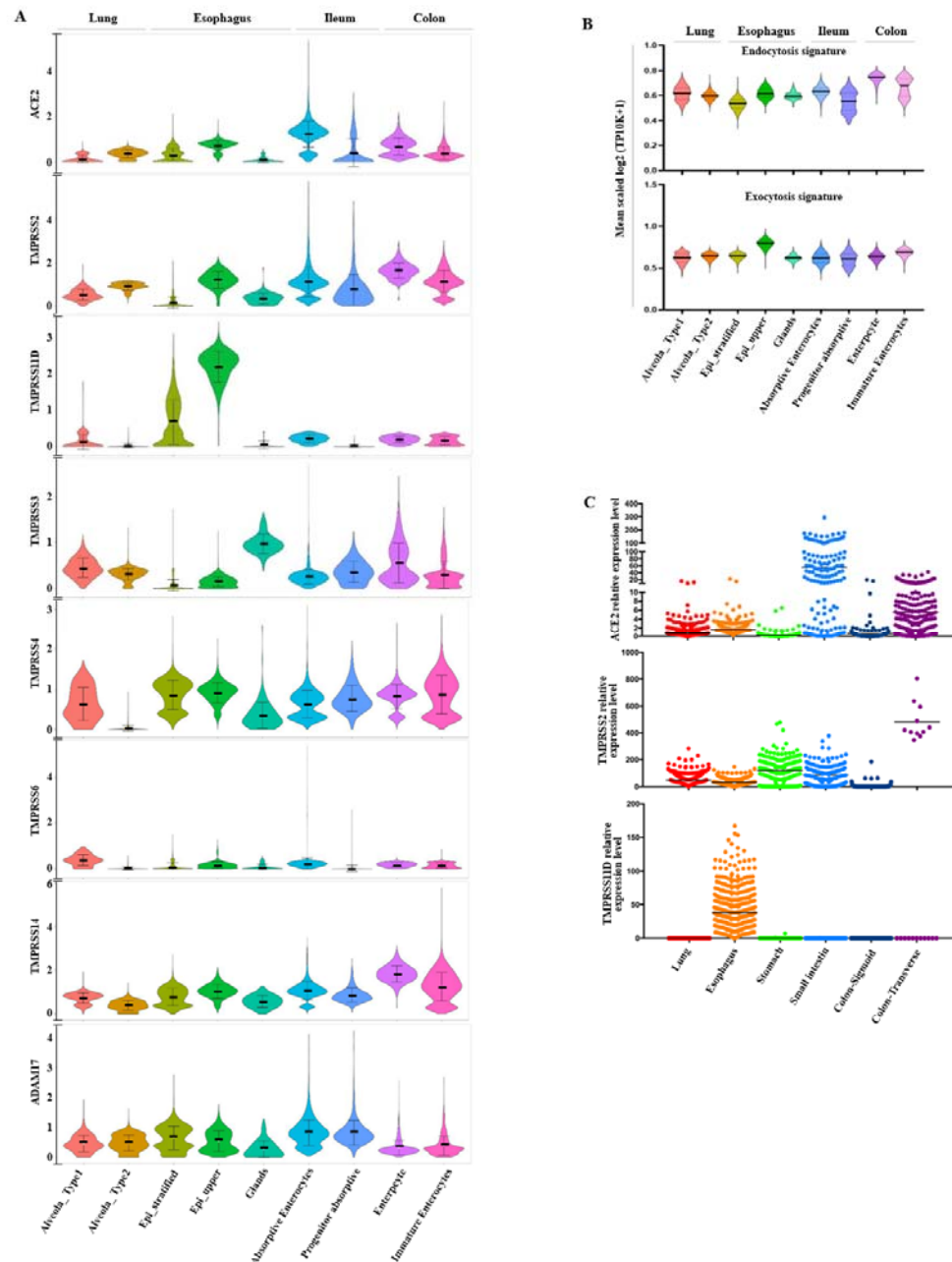


Figure 3. Expression levels of ACE2, TMPRSSs, ADAM17 and functional gene sets in lung and digestive tracts.

(A). Violin plots showing the expression levels of ACE2, TMPRSSs and ADAM17 in 2 lung clusters and 7 digestive tract clusters. The gene expression matrix was normalized and denoised to remove unwanted technical variability across the

1 datasets.

2 (B). Violin plots showing the expression levels of endocytosis and
3 exocytosis-associated genes. The expression is measured as the mean log2 (TP10K+1)
4 value.

5 (C). Expression levels of ACE2, TMPRSS2 and ADAM17 at RNA level in different
6 tissues. The expression is measured as the pTPM value in the RNA-seq data from the
7 GTEx database.

8

9

10

11

12

13

14

15

16

17

18

19

20

21

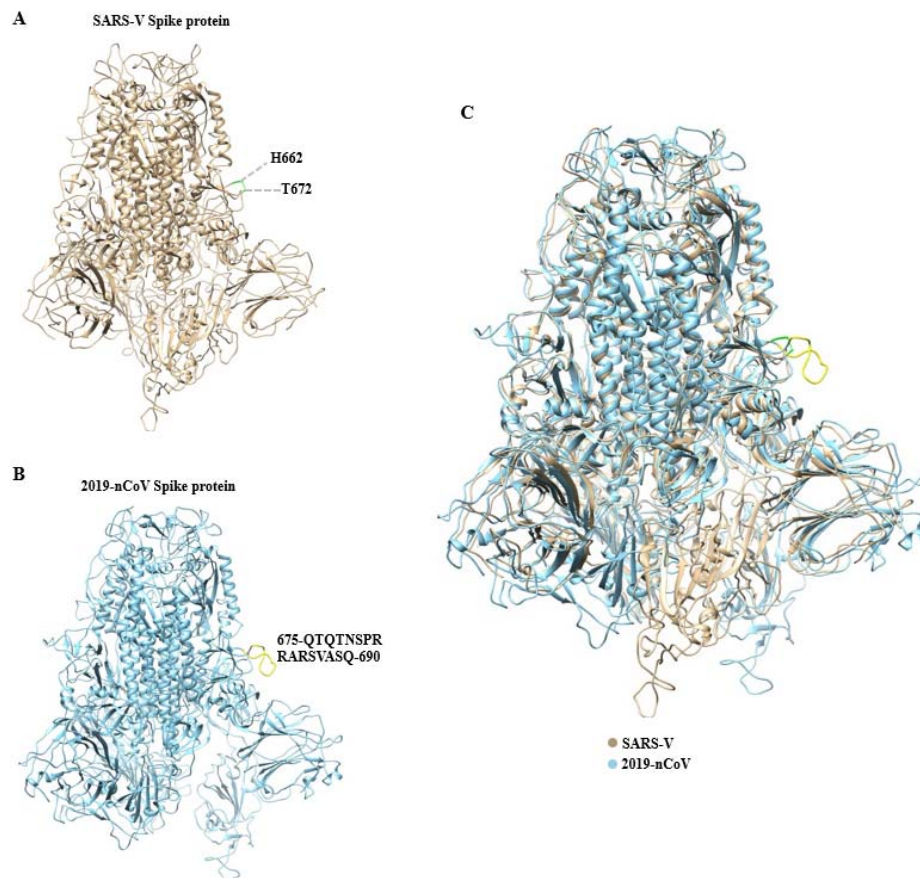


Fig 4. Overall structure of the SARS-CoV Spike protein and 2019-nCoV Spike protein homo-trimers

(A). Structure of the SARS-CoV Spike protein (from PDB: 5X5B). The insert aa675-690 to SARS-CoV Spike protein aa661-672 with the structural missed residues are colored green.

(B). Structure of the 2019-nCoV Spike protein (Modelled by SWISS-MODEL). The insert aa675-690 of 2019-nCoV Spike protein that corresponds to the insert region of SARS-V Spike protein is colored yellow

(C). The structural superimpose of SARS-CoV Spike protein (yellow) and 2019-nCoV Spike protein (blue)

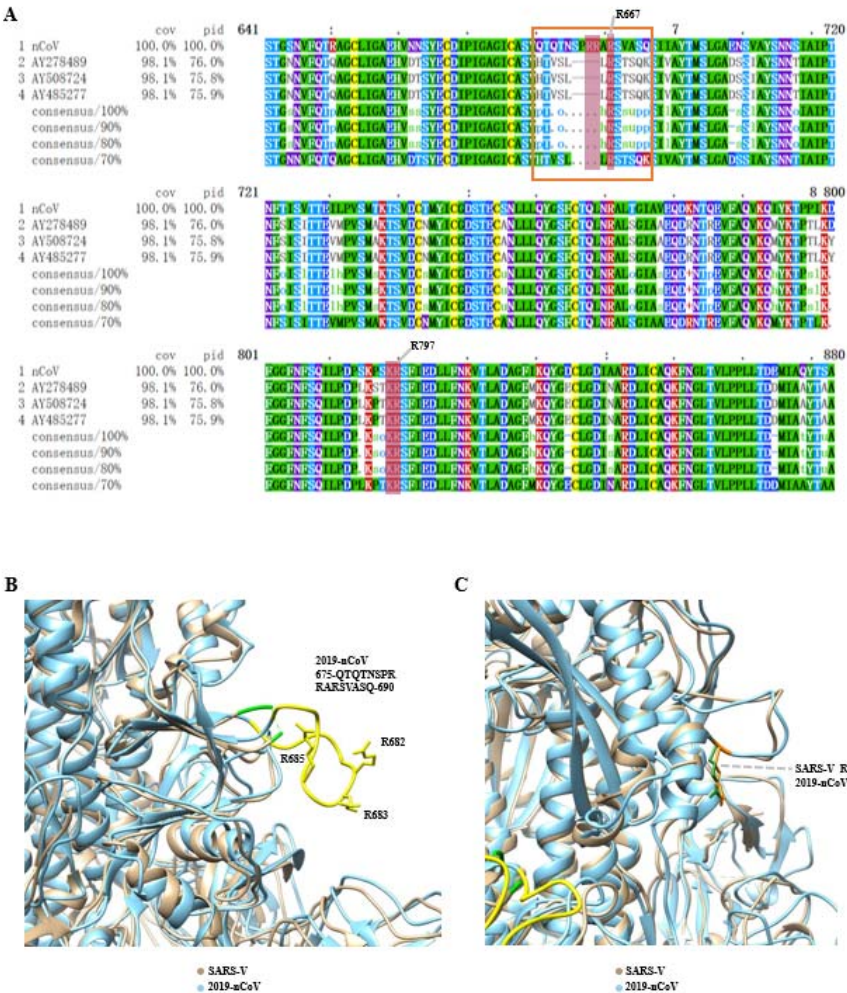
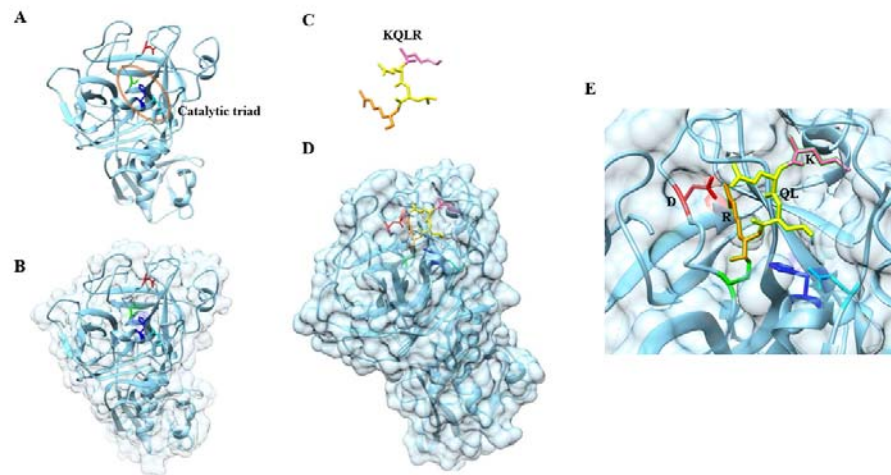


Fig 5. The two potential cleavage sites of SARS-CoV Spike protein and 2019-nCoV Spike protein by TMPRSS2.

(A). Amino acid sequence alignment of the SARS-CoV Spike protein with 2019-nCoV Spike protein. The insert 675-690 that corresponds to 2019-nCoV Spike protein is boxed in orange. Two potential cleavage sites by TMPRSS2, R667 and R797 are marked.

(B-C). Structural alignment in detail of the SARS-CoV Spike protein with 2019-nCoV Spike protein. Shown are the insert 675-690 of 2019-nCoV Spike protein (yellow) and the corresponding loci to SARS-CoV Spike protein 661-672 (green). Three important residues, R682, R683, R685, are specially marked in (B). The similarly SARS-CoV R797 with 2019-nCoV R815 are colored forest green and

orange, respectively (C).



2

3 **Fig 6. Structure and catalytic mechanism of TMPRSS2**

4 **(A-B).** Overall structure and surface of TMPRSS2 (Modelled by SWISS-MODEL).

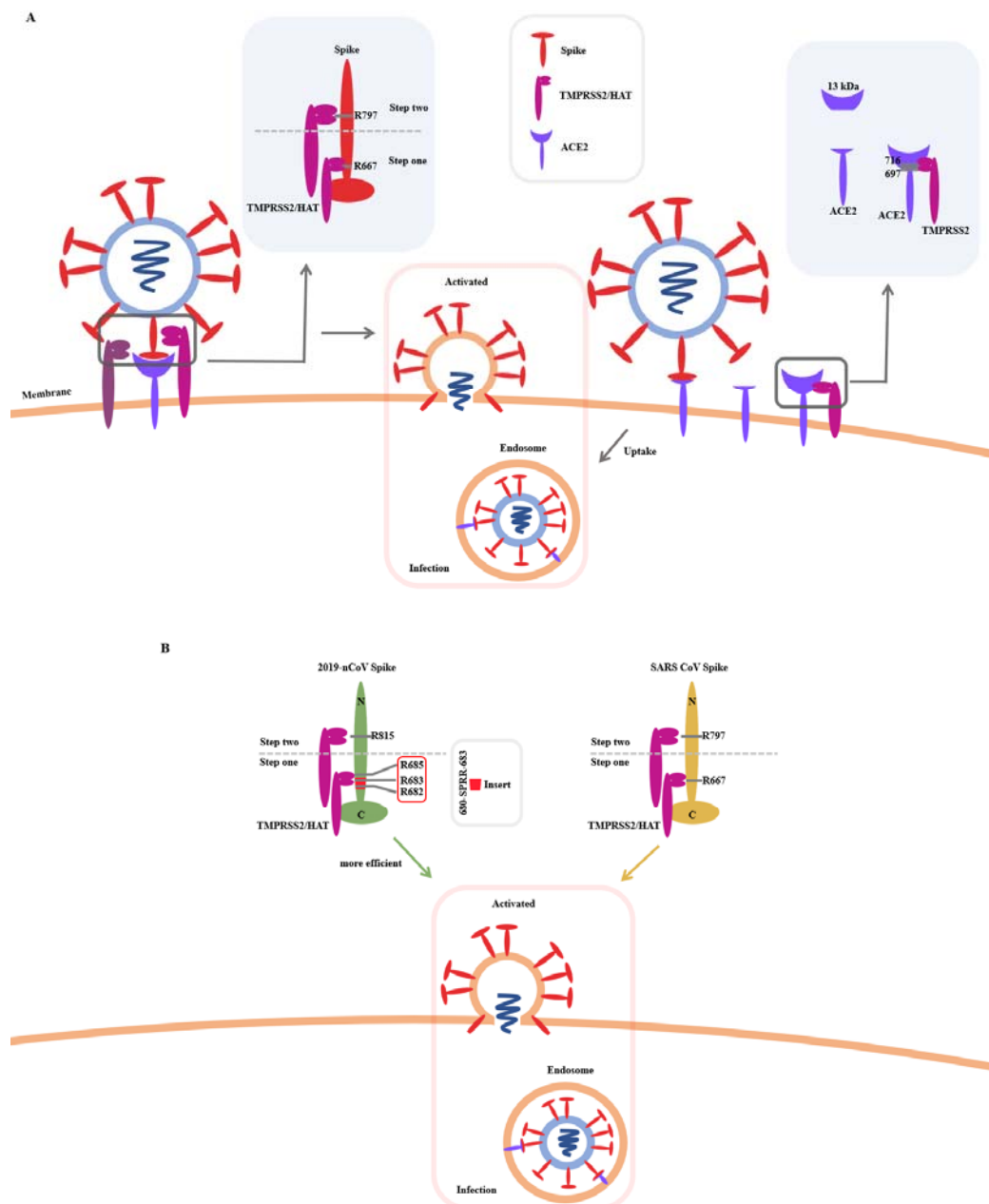
5 The catalytic triad comprised of H296, D345 and S441 are colored blue, green and
6 cyan, respectively. The substrate binding residue D435 which located in the bottom of
7 pocket is marked in red, the substrate binding pocket is deeper than most of serine
8 proteinase.

9 **(C).** Polypeptide substrate analogue KQLR. Cleaved site Arg is coloured orange. Gln,
10 Leu are colored yellow, and Lys is colored pink.

11 **(D-E).** The state of substrate analogue binding in the catalytic pocket, and the detail
12 shown in **(E)**. The state of substrate analogue binding in the catalytic pocket, and the
13 detail shown in E, Arg of substrate analogue is strongly interacted with D435, shown
14 in **(E)**, Arg of substrate analogue is strongly interacted with D435.

15

16



1

2 **Fig 7 : Role of TMPRSS2/HAT proteases in the cellular entry of 2019-nCoV.**

3 (A). Routes of coronavirus entering host cells. 2019-nCoV could enter host cells via
4 two distinct routes, depending on the availability of cellular proteases required for
5 activation of 2019-nCoV. The first route of activation can be achieved if the
6 2019-nCoV activating protease TMPRSS2 and ACE2 are co-expressed on the surface
7 of target cells. The spiked protein binds to ACE2 through its S1 subunit and is treated

1 by TMPRSS2 at the R667 or R797 in the S1/S2 site. This activates Spike protein and
2 allows 2019-nCoV fusion at the cell surface. 2019-nCoV was encapsulated into
3 cellular vesicles before transport of virions into host cell endosomes. Uptake can be
4 enhanced if TMPRSS2 also cleaves ACE2 amino acids 697 to 716, resulting in
5 shedding of 13kD ACE2 fragment in culture supernatants. The second route can be
6 pursued if there were no 2019-nCoV activating proteases expressed at the cell surface.
7 While binding of virion-associated Spike protein to ACE2, the virions are taken up
8 into endosomes, where 2019-nCoV could be cleaved and activated by the
9 pH-dependent protease.

10 (B). The difference between 2019-nCoV and SARS-CoV in activating the Spike
11 protein. The Spike protein of SARS involves two cleavage sites recognized by
12 TMPRSS2, one at arginine 667 and one at arginine 797. (right). Compared with
13 SARS-Cov, the S protein of 2019-nCoV (left)has an insertion sequence
14 680-SPRR-683 (grey box)at the TM cleavage site. We speculated that R682, R682
15 and R685 (red box) could be used as the most suitable substrates for TM, which can
16 increase the cleavage efficiency of TM to S protein, promote the activation of
17 2019-nCoV, and make 2019-nCoV more infectious.

18
19
20