

# 1    **The insert sequence in SARS-CoV-2 enhances spike protein cleavage by**

## 2    **TMPRSS**

3    Tong Meng<sup>1,2,10,11</sup>, Hao Cao<sup>3,4,11</sup>, Hao Zhang<sup>5,10,11</sup>, Zijian Kang<sup>6,10</sup>, Da Xu<sup>7,10</sup>, Haiyi  
4    Gong<sup>5,10</sup>, Jing Wang<sup>8</sup>, Zifu Li<sup>8</sup>, Xingang Cui<sup>7</sup>, Huji Xu<sup>4,6</sup>, Haifeng Wei<sup>5</sup>, Xiuwu Pan<sup>7</sup>,  
5    Rongrong Zhu<sup>9</sup>, Jianru Xiao<sup>5\*</sup>, Wang Zhou<sup>4,10\*</sup>, Liming Cheng<sup>1\*</sup>, Jianmin Liu<sup>8\*</sup>.

6    1 Division of Spine, Department of Orthopedics, Tongji Hospital affiliated to Tongji  
7    University School of Medicine, 200065 Shanghai, China

8    2 Tongji University Cancer Center, School of Medicine, Tongji University, 200092  
9    Shanghai, China

10    3 School of Life Science and Biopharmaceutics, Shenyang Pharmaceutical University,  
11    103 Wenhua Road, 110016 Shenyang, China

12    4 Peking-Tsinghua Center for Life Sciences, TsinghuaUniversity, 100084 Beijing,  
13    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 Deparntment of Urology, The Third Affiliated Hospital of Second Military Medical  
19    University, 201805 Shanghai, China

20    8 Department of Neurosurgery, Changhai hospital, Second Military Medical  
21    University, 200003 Shanghai, China

1    9 Key Laboratory of Spine and Spinal Cord Injury Repair and Regeneration of  
2    Ministry of Education, Orthopaedic Department of Tongji Hospital, School of Life  
3    Science and Technology, Tongji University, 200092 Shanghai, China  
4    10 Qiu-Jiang Bioinformatics Institute, 200003 Shanghai, China  
5    11 These authors contributed equally to this work: Tong Meng, Hao Cao, Hao Zhang  
6    \*e-mail: chstroke@163.com; limingcheng@tongji.edu.cn; brilliant212@163.com;  
7    jianruxiao83@163.com

## 8    **Abstract**

9    At the end of 2019, the SARS-CoV-2 induces an ongoing outbreak of pneumonia in  
10    China<sup>1</sup>, even more spread than SARS-CoV infection<sup>2</sup>. The entry of SARS-CoV into  
11    host cells mainly depends on the cell receptor (ACE2) recognition and spike protein  
12    cleavage-induced cell membrane fusion<sup>3,4</sup>. The spike protein of SARS-CoV-2 also  
13    binds to ACE2 with a similar affinity, whereas its spike protein cleavage remains  
14    unclear<sup>5,6</sup>. Here we show that an insertion sequence in the spike protein of  
15    SARS-CoV-2 enhances the cleavage efficiency, and besides pulmonary alveoli,  
16    intestinal and esophagus epithelium were also the target tissues of SARS-CoV-2.  
17    Compared with SARS-CoV, we found a SPRR insertion in the S1/S2 protease  
18    cleavage sites of SARS-CoV-2 spike protein increasing the cleavage efficiency by the  
19    protein sequence alignment and furin score calculation. Additionally, the insertion  
20    sequence facilitates the formation of an extended loop which was more suitable for  
21    protease recognition by the homology modeling and molecular docking. Furthermore,

1 the single-cell transcriptomes identified that ACE2 and TMPRSSs are highly  
2 coexpressed in AT2 cells of lung, along with esophageal upper epithelial cells and  
3 absorptive enterocytes. Our results provide the bioinformatics evidence for the  
4 increased spike protein cleavage of SARS-CoV-2 and indicate its potential target  
5 cells.

## 6 **Introduction**

7 At the end of 2019, a rising number of pneumonia patients with unknown pathogen  
8 emerged from Wuhan to nearly the entire China<sup>7</sup>. A novel coronavirus was isolated  
9 and based on its phylogeny, taxonomy and established practice, the Coronavirus  
10 Study Group (CSG) recognized it as a sister to severe acute respiratory syndrome  
11 coronaviruses (SARS-CoVs) and labeled it as severe acute respiratory syndrome  
12 coronavirus 2 (SARS-CoV-2)<sup>1,8</sup>. Although SARS-CoV-2 is generally less pathogenic  
13 than SARS-CoV and Middle East respiratory syndrome coronavirus (MERS-CoV), it  
14 has a relatively high transmissibility<sup>9</sup>.

15 With regard to human coronavirus, the transmissibility and infectivity is largely  
16 controlled by the spike (S) surface envelope protein<sup>10</sup>. Its surface unit (S1) mediates  
17 the entry into host cells by binding to cell receptor and the transmembrane unit (S2)  
18 subunit regulates the fusion of viral and cellular membranes<sup>3</sup>. Prior to membrane  
19 fusion, the S protein should be cleaved and activated to allow for the fusion peptide  
20 releasing onto host cell membranes (Fig. 1a)<sup>11</sup>. SARS-CoV-2 uses the same cell  
21 receptor (angiotensin converting enzyme II, ACE2) as SARS-CoV, with a similar

1 binding affinity, whereas their transmissibility and infectivity are different<sup>5,6,12,13</sup>.

2 Thus, the different virus transmission and infectivity may be associated with the

3 differentiated protease-induced S protein cleavage between SARS-CoV-2 and

4 SARS-CoV.

5 The transmembrane serine proteases (TMPRSSs) were the main host cell proteases on

6 the cell membrane<sup>14</sup>. The substrate specificity of TMPRSSs are almost similar and

7 revealing a strong preference for arginine or lysine residues in the P1 position.

8 Nowadays, their hydrolytic effects of TMPRSSs have been widely reported in

9 SARS-CoV and MERS-CoV pneumonia<sup>15</sup>. In the SARS-CoV-infected alveolar cells,

10 TMPRSSs, especially the TMPRSS2 and TMPRSS11D, cleave the SARS-CoV S

11 protein (SARS-S) at residue R667 (the S1/S2 cleavage site) and residue R797 (the S2'

12 cleavage site) (Fig. 1a)<sup>15,16</sup>. Besides cleaving S protein, they can also promote viral

13 spread in the host by cleaving ACE2 (Fig. 1b)<sup>14,17</sup>. Although SARS-CoV-2 and

14 SARS-CoV share the same host cell receptor with a similar affinity, however, the

15 SARS-CoV-2 S protein cleavage induced by TMPRSS remains unclear which may be

16 associated with the viral infectivity<sup>4,5</sup>.

## 17 **Results**

### 18 **The comparison of the S1/S2 and S2' cleavage sites between SARS-CoV-2 and**

### 19 **SARS-CoV**

20 Generally, compared with SARS-CoV, the major differences in SARS-CoV-2 are the

21 three short insertions in the N-terminal domain and four out of five key residues

1 changes in the receptor-binding motif<sup>5</sup>. Here we used the alignment, furin score and  
 2 homology modeling to compare the sequence of the S1/S2 and S2' cleavage sites  
 3 between SARS-CoV-2 and SARS-CoV (Fig. 1c). The amino acid sequence of the  
 4 S1/S2 and S2' cleavage sites among ten beta-coronavirus were then analyzed and we  
 5 found that compared with SARS, there was an insertion sequence (SPRR) in the  
 6 S1/S2 cleavage sites of SARS-CoV-2 (Fig. 2a). The furin score was next used to  
 7 identify the cleavage efficiency of the insertion sequence in SARS-CoV-2. Its furin  
 8 score was 0.688, which was obviously higher than that of the corresponding sequence  
 9 in SARS-CoV (0.139), indicating that the insertion sequence may increase the  
 10 cleavage efficiency by proteases (Fig. 2b).

11 The structures of SARS-S and SARS-CoV-2 S protein were presented in Extended  
 12 Data Fig. 1a and 1b, along with their structural superimposition (Extended Data Fig.  
 13 1c). The structural comparison of homology modeling SARS-CoV-2 S protein with  
 14 SARS-S protein (PDB: 5x5b) showed that a exposed loop was formed by the insertion  
 15 which comprised R682 and R683 (S1/S2 site) on the surface of SARS-CoV-2 S  
 16 protein, and no significant difference of them in S2' site (Fig. 2c, d).

# **17 The insertion sequence of SARS-CoV-2 facilitating the TMPRSS recognition and** **18 S protein cleavage**

19 Structurally, TMPRSSs include extracellular domain, transmembrane domain and  
 20 intracellular domain in which extracellular domain is the main catalytic domain. They  
 21 show similar substrate-specificity and catalytic mechanism. Take TMPRSS2 as an

1 example. The catalytic triad consisted of H296, D345 and S441 and the substrate  
2 binding residue D435, a conserved aspartate residue, was located in the bottom of  
3 pocket<sup>18,19</sup>. The substrate binding pocket is deeper than most of serine proteinase  
4 (Extended Data Fig. 2a, b). The bottom of the catalytic pocket has a negatively  
5 charged aspartic acid residue which can facilitate the binding and stabilization of  
6 arginine or lysine residues in the P1 position<sup>18,19</sup>.  
7 Polypeptide substrate analogue KQLR included arginine, glutamine, leucine and  
8 lysine (Extended Data Fig. 2c). The substrate analogue could bind to the catalytic  
9 pocket of TMPRSS2 (Extended Data Fig. 2d, e). The conformation of the insertion  
10 sequence in SARS-CoV-2 S protein and TMPRSS2 was next simulated by molecular  
11 docking. We found the insertion sequence formed a loop which was easily recognized  
12 by the catalytic pocket of TMPRSS2 (Extended Data Fig. 2f, g). Thus, both the furin  
13 score and molecular docking revealed that the insertion sequence of SARS-CoV-2  
14 facilitates the TMPRSS2 recognition and S protein cleavage.

## 15 **The potential target tissues of COVID-19**

16 The entry of SARS-CoV-2 into host cells depends on the cell receptor recognition and  
17 cell proteases cleaving. Thus, the target cells should coexpress both the cell receptor  
18 ACE2 and cell proteases TMPRSSs. In order to identify the coexpressing cell  
19 composition and proportion, we utilized 3 datasets including 32 samples and built the  
20 largest single-cell transcriptome atlas of normal lung, the commonest infected organ  
21 of SARS-CoV-2.

1 After initial quality controls, a total of 113,045 cells and 29 sub-clusters were  
2 identified in the lung (Fig. 3a). The marker genes and dataset proportions of each  
3 sub-cluster were presented in Extended Data Fig. 3-4.

4 We detected the expression of ACE2 and TMPRSSs in 29 cell groups, in which the  
5 expression of the whole 17 TMPRSS genes is in the form of total signature value.

6 Pseudocoloring analysis was performed and we found that ACE2 was mainly expressed  
7 in AT2 cells and marked with red (Fig. 3b, c). The total 17 TMPRSS genes was found  
8 in AT1, AT2, airway secretory and ciliated cells colored with blue (Fig. 3b, d,  
9 Extended Data Fig. 5a). Thus, we found an obvious coexpression between TMPRSSs  
10 and ACE2 in AT2. Among the whole TMPRSS genes, TMPRSS1 and TMPRSS2  
11 were highly expressed in AT2 and AT1 cells, which were co-expressed with ACE2 in  
12 lung (Fig. 3b, Extended Data Fig. 5b). Due to the entry of virus into host cell is  
13 related to endocytosis, we also detected the endocytosis-related genes among different  
14 cells. We found that these genes had consistent distribution and highly expressed in  
15 AT1, AT2, airway secretory, ciliated cells and M2 macrophage (Extended Data Fig.  
16 5c).

17 Due to the RNA of SARS-CoV-2 was also found in the stool specimen of the  
18 SARS-CoV-2-infected patient<sup>20</sup>, the digestive system may also be the potential route  
19 of COVID-19. Thus, in addition to lung, 4 datasets with the single-cell transcriptomes  
20 of the esophagus, gastric, small intestine and colon were analyzed to identify the  
21 expression of ACE2 and TMPRSSs in the digestive system. The co-expression of

1 ACE2 and TMPRSS was analyzed in esophagus, stomach, small intestine and colon  
2 by 87947, 29678, 11218 and 47442 high-quality single cells, respectively (Extended  
3 Data Fig. 6a). The coexpression of ACE2 and total TMPRSS genes were found in the  
4 upper epithelial cells of esophagus, the absorptive enterocytes of ileum epithelia and  
5 the enterocytes of colon epithelia (Extended Data Fig. 6b-e, 7a-d).

6 As both ACE2 and TMPRSSs are expressed in the lung and digestive system, we next  
7 compared their relative expression values in the ACE2-expressing cells. A similar  
8 distribution was found between ACE2 and TMPRSSs in all the 9 clusters with high  
9 expressions in the esophageal upper epithelial cells, the ileal absorptive enterocytes  
10 and the colonic enterocytes (Fig. 4a). In addition, their expression of AT2 was  
11 relatively lower than that of epithelial cells in the digestive system. Among all the  
12 TMPRSSs, TMPRSS1 and TMPRSS2 were relatively highly expressed in AT2, and  
13 most TMPRSSs were highly found in the esophageal upper epithelial cells (Extended  
14 Data Fig. 8a). The endocytosis- and exocytosis-associated genes which are related to  
15 the entry of virus into host cells and virus infection were also detected in all the 9  
16 clusters. The endocytosis signature was more expressed in AT1 and AT2 cells,  
17 whereas the exocytosis signature was highly gathered in esophageal upper epithelial  
18 cells. It can explain that the commonest infected tissue in COVID-19 is pulmonary  
19 alveoli and SARS-CoV-2 can also be detected in the esophageal erosion (Fig. 4b)<sup>21</sup>.

20 The RNA-seq data of lung, esophagus, stomach, small intestine, colon-transverse and  
21 colon-sigmoid were obtained from GTEx database. The expressions of ACE2 and



1 TMPRSS2 also had a similar tendency and were highly expressed in small intestine  
2 and colon, while the TMPRSS11D was mainly found in the esophagus (Extended  
3 Data Fig. 8b).

#### 4 **Discussion**

5 The coronaviruses is the common infection source of respiratory, enteric and central  
6 nervous system in humans and other mammals<sup>22</sup>. At the beginning of the twenty-first  
7 century, two betacoronaviruses, SARS-CoV and MERS-CoV, result in persistent  
8 public panics and became the most significant public health events<sup>23</sup>. In December  
9 2019, a novel identified coronavirus (SARS-CoV-2) induced an ongoing outbreak of  
10 pneumonia in Wuhan, Hubei, China <sup>7</sup>. The rapidly increasing number of  
11 SARS-CoV-2-infected cases suggests that SARS-CoV-2 may be transmitted  
12 effectively among humans and give rise to a high pandemic potential <sup>7,8,24</sup>.  
13 Previous studies identified that SARS-CoV mutated between 2002 and 2004 to better  
14 bind to its cell receptor, replicate in human cells and enhance the virulence <sup>9</sup>. Thus, it  
15 is important to explore whether SARS-CoV-2 behaves like SARS-CoV to adapt to the  
16 host cell. Notably, SARS-CoV and SARS-CoV-2 share the same receptor protein  
17 ACE2<sup>5,25</sup>. Besides, the receptor-binding domain (RBD) in S protein of SARS-CoV-2  
18 binds to ACE2 with the similar affinity as SARS-CoV RBD does<sup>6</sup>. Thus, during the  
19 process of viral and host cellular membrane fusion, whether the specific structure of  
20 SARS-CoV-2 S protein seems better suited to be activated by host cell proteases may

1 be related to the different virus infectivities and transmissibilities between  
2 SARS-CoV and SARS-CoV-2<sup>6</sup>.  
3 In this study, we found the furin score of the S1/S2 cleavage sites in SARS-CoV-2  
4 was higher than that of SARS, implying a more degree of hydrolysis. Through the  
5 comparison of the two structures, R682, R683 and relative S680, P681 extended the  
6 original exposed loop combined with R685 of SARS-CoV-2, which was more  
7 suitable for hydrolysis by TMPRSSs. The substrate specificity of TMPRSSs are  
8 almost similar, revealing a strong preference for arginine or lysine residues in the P1  
9 position represented by R. More R (R682, R683 and R685) in the S1/S2 cleavage  
10 sites of SARS-CoV-2 can enhance the cleavage of S1 with S2, which means that the  
11 structural constraints of S1 on S2 is removed, and the fusion peptides in S2 are  
12 exposed and insert into the target host cell membrane, finally it increases the  
13 efficiency of fuse membranes<sup>18,19</sup>.  
14 By the way, some researchers previously supposed the SARS-CoV-2 was artificial  
15 due to four inserts in the S protein of SARS-CoV-2 from HIV sequence. However, the  
16 results of protein sequence alignment revealed that the similar sequence of the  
17 reported fourth insertion site (680-SPRR-683) in SARS-CoV-2 was commonly found  
18 in many beta-coronavirus. Therefore, we supposed that based on the current evidence,  
19 it is not scientific to consider the insertion sequence in SARS-CoV-2 S protein being  
20 artificial.

1 With the help of single cell sequencing, we found a strong co-expression between  
2 ACE2 and TMPRSSs, in especial TMPRSS1 and TMPRSS2, in lung AT2 cells, which  
3 was also the main infected cell type in SARS-CoV pneumonia<sup>26</sup>. Moreover, we also  
4 found the endocytosis-associated genes was highly expressed in AT2 cells, implying  
5 that endocytosis may also facilitate the entry of SARS-CoV-2 into host cells. As the  
6 alveolar stem-like cells, AT2 cells are in charge of surfactant biosynthesis,  
7 self-renewal and immunoregulation<sup>27</sup>. Thus, SARS-CoV-2 not only damages the AT2  
8 cells leading to the direct injury to alveoli, but also raises alveolar surface tension to  
9 induce dyspnea<sup>28</sup>. Additionally, the injured AT2 also damages the immunologic  
10 balance in alveoli and results in inflammatory cascade<sup>29</sup>. In addition, they are also  
11 highly co-expressed in absorptive enterocytes and upper epithelial cells of esophagus,  
12 implying that intestinal epithelium and esophagus epithelium may also be the  
13 potential target tissues. This can explain the cases whose SARS-CoV-2 was detected  
14 in the esophageal erosion or stool specimen, implying that the digestive system is a  
15 potential route of COVID-19<sup>7,20,21</sup>.  
16 Due to the critical role of TMPRSSs in influenza virus and coronavirus infections,  
17 serine protease inhibitors, such as camostat, nafamostat and leupeptin, have been used  
18 in the antiviral treatment targeting TMPRSSs with high antiviral activities<sup>14,30,31</sup>.  
19 Nowadays, Remdesivir (GS-5734) has been used in the treatment of SARS-CoV-2,  
20 however, the therapeutic effects are still unclear. Based on our results, we also  
21 supposed that TMPRSSs may also serve as candidate antiviral targets for

1 SARS-CoV-2 infection and the clinical trials of serine protease inhibitors should also  
2 be performed for COVID-19.

### 3 **Methods**

#### 4 **Structure modelling**

5 The structures of SARS-CoV-2 S protein and TMPRSS2 were generated by  
6 SWISS-MODEL online server<sup>32</sup>. The structures were marked, superimposed and  
7 visualized by Chimera<sup>33</sup>. To further explore the possible catalytic mechanism of the  
8 SARS-CoV-2 S protein cleaved by TMPRSS2, ZDOCK program was used to predict  
9 their interaction<sup>34</sup>. A total of 5000 models were generated and were set to 50 clusters,  
10 then the best scoring models from the 5 largest clusters were selected for further  
11 analysis.

#### 12 **Furin score**

13 The fragmentation maps, scoring and residue coverage analysis were conducted using  
14 arginine and lysine propeptide cleavage sites prediction algorithms ProP 1.0 server<sup>35</sup>.

#### 15 **Single cell transcriptome data sources**

16 Single cell transcriptome data were obtained from Single Cell Portal  
17 ([https://singlecell.broadinstitute.org/single\\_cell](https://singlecell.broadinstitute.org/single_cell)), Human Cell Atlas Data Portal  
18 (<https://data.humancellatlas.org>) and Gene Expression Omnibus (GEO;  
19 <https://www.ncbi.nlm.nih.gov/>). Esophageal and lung data were obtained from the  
20 research of E Madissoon *et al* containing 21 esophageal and 19 lung tissue samples<sup>36</sup>.  
21 Two lung datasets were further obtained from GSE122960<sup>38</sup> and GSE128169<sup>39</sup>,

1 including eight and five lung tissues respectively. GSE134520 included 6 gastric  
2 mucosal samples from 3 non-atrophic gastritis and 2 chronic atrophic gastritis  
3 patients<sup>40</sup>. GSE134809 comprises 11 noninflammatory ileal samples from Crohn's  
4 disease patients<sup>41</sup>. The data from Christopher S *et al* consisted of 12 normal colon  
5 samples<sup>42</sup>.

## 6 **Quality control**

7 Cells would be identified as poor-quality once (1) the number of expressed genes  
8 fewer than 200 or greater than 5000, or (2) more than 20% of UMIs being mapped to  
9 mitochondrial or ribosomal genes.

## 10 **Data Integration, Dimension Reduction and Cell Clustering**

11 Different methods were performed to process the downloaded data:

- 12 1. Esophagus dataset. Rdata were obtained and dimension reduction and clustering  
13 had already been implemented by the authors<sup>36</sup>.
- 14 2. Lung, stomach and ileum datasets. We utilized functions in the Seurat package to  
15 normalize and scale the single-cell gene expression data<sup>43</sup>. Unique  
16 molecular identifier (UMI) counts were normalized by the total number of UMIs  
17 per cell, multiplied by 10000 for normalization and log-transformed using the  
18 "NormalizeData" function. Then, multiple sample data within each dataset were  
19 merged using the "FindIntegrationAnchors" and "IntegrateData" functions. After  
20 identifying highly variable genes (HVGs) using the "FindVariableGenes" function  
21 a principal component analysis (PCA) was performed on the single-cell

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

7 3. Colon Dataset. The single cell data was processed with the R packages LIGER<sup>44</sup>  
8 and Seurat<sup>43</sup>. The gene expression matrix was first normalized to remove  
9 differences in sequencing depth and capture efficiency among cells. Variable  
10 genes in each dataset were identified using the “selectGenes” function. Then we  
11 used the “optimizeALS” function in LIGER to perform the integrative  
12 nonnegative matrix factorization and selecte a k of 15 and lambda of 5.0 to obtain  
13 a plot of expected alignment. The “quantileAlignSNF” function was then  
14 performed to builds a shared factor neighborhood graph to jointly cluster cells,  
15 then quantile normalizes corresponding clusters. Next nonlinear dimensionality  
16 reduction was calculated using the “RunUMAP” function and the results were  
17 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

1 shown in Feature plots and violin plots. We used “Quantile normalization” in the R  
2 package preprocessCore (R package version 1.46.0.  
3 <https://github.com/bmbolstad/preprocessCore>) to remove unwanted technical  
4 variability across different datasets. The data were further denoised to compare the  
5 gene expression levels of gene signature.  
6 Endocytosis or exocytosis associated genes were obtained from Harmonizome dataset  
7 <sup>45</sup>. Mean expressions of the genesets were calculated to compare the ability of  
8 endocytosis or exocytosis among clusters.  
9 To minimize bias, external databases of Genotype-Tissue Expression (GTEx)<sup>46</sup> was  
10 used to detect gene expression of ACE2, TMPRSS1 and TMPRSS2 at the tissue  
11 levels including normal lung and digestive system, such as esophagus, stomach, small  
12 intestine and colon.

13

14 **Acknowledgements** This study was jointly supported by the National Natural  
15 Science Foundation of China (Grants 81702659 and 81572746) and National Key  
16 R&D Program of China (Grants 2016YFA0100800).

17 **Author contributions** J.L., L.C., W.Z. and J.X. conceived the idea and directed the  
18 team. T.M., H.C., H.Z. and W.Z. designed and coordinated the analysis and  
19 characterization. H.Z., Z.K., D.X., H.G. performed single-cell sequencing and  
20 characterization under the guidance of X.C., H.X., and H.W.. Data collection and  
21 generation were performed by J.W., Z.L., R.Z. and X.P.. Data interpretation was

1 performed by J.L., L.C., W.Z. and J.X.. The alignment and structure comparison was  
2 performed by H.C. under the guidance of W.Z. The manuscript was written by T.M.,  
3 H.C., Z.K. and W.Z. All authors contributed to the analysis and discussion of the  
4 results leading to the manuscript.

5 **Competing interests** The authors declare no competing interests.

6

7

8 1 Gorbalenya AE, B. S., Baric RS, de Groot RJ, Drosten C, Gulyaeva AA,  
9 Haagmans BL, Lauber C, Leontovich AM, Neuman BW, Penzar D,  
10 Perlman S, Poon LL, Samborskiy D, Sidorov IA, Sola I, Ziebuhr J.  
11 Severe acute respiratory syndrome-related coronavirus: The species  
12 and its viruses – a statement of the Coronavirus Study Group. *bioRxiv*  
13 doi:doi: <https://doi.org/10.1101/2020.02.07.937862> (2020).

14 2 Zhong, N. S. *et al.* Epidemiology and cause of severe acute respiratory  
15 syndrome (SARS) in Guangdong, People's Republic of China, in  
16 February, 2003. *Lancet* 362, 1353-1358,  
17 doi:10.1016/s0140-6736(03)14630-2 (2003).

18 3 Hofmann, H. & Pohlmann, S. Cellular entry of the SARS coronavirus.  
19 *Trends in microbiology* 12, 466-472, doi:10.1016/j.tim.2004.08.008  
20 (2004).

21 4 Chan, J. F. *et al.* A familial cluster of pneumonia associated with the



1           2019 novel coronavirus indicating person-to-person transmission: a  
2           study of a family cluster. *Lancet*, doi:10.1016/s0140-6736(20)30154-9  
3           (2020).

4    5       P Zhou, X. Y., XG Wang, B Hu, L Zhang, W Zhang, HR Si, Y Zhu, B Li,  
5           CL Huang, HD Chen, J Chen, Y Luo, H Guo, RD Jiang, MQ Liu, Y Chen,  
6           XR Shen, X Wang, XS Zheng, K Zhao, QJ Chen, F Deng, LL Liu, B Yan,  
7           FX Zhan, YY Wang, GF Xiao, ZL Shi. A pneumonia outbreak  
8           associated with a new coronavirus of probable bat origin. *Nature*,  
9           doi:[https:// doi.org/10.1038/s41586-020-2012-7](https://doi.org/10.1038/s41586-020-2012-7) (2020).

10   6       Tian XL, L. C., Huang A, Xia S, Lu SC, Shi ZL, Lu L, Jiang SB, Yang ZL,  
11           Wu YL, Ying TL. Potent binding of 2019 novel coronavirus spike protein  
12           by a SARS coronavirus-specific human monoclonal antibody. *bioRxiv*  
13           doi:doi: <https://doi.org/10.1101/2020.01.28.923011> (2020).

14   7       Huang, C. *et al.* Clinical features of patients infected with 2019 novel  
15           coronavirus           in           Wuhan,           China.           *Lancet*,  
16           doi:10.1016/s0140-6736(20)30183-5 (2020).

17   8       Zhu, N. *et al.* A Novel Coronavirus from Patients with Pneumonia in  
18           China, 2019. *The New England journal of medicine*,  
19           doi:10.1056/NEJMoa2001017 (2020).

20   9       Chen, J. Pathogenicity and Transmissibility of 2019-nCoV-A Quick  
21           Overview and Comparison with Other Emerging Viruses. *Microbes and*

1        *infection*, doi:10.1016/j.micinf.2020.01.004 (2020).

2    10    Walls, A. C. *et al.* Unexpected Receptor Functional Mimicry Elucidates  
3        Activation of Coronavirus Fusion. *Cell* 176, 1026-1039.e1015,  
4        doi:10.1016/j.cell.2018.12.028 (2019).

5    11    Gallagher, T. M. & Buchmeier, M. J. Coronavirus spike proteins in viral  
6        entry and pathogenesis. *Virology* 279, 371-374,  
7        doi:10.1006/viro.2000.0757 (2001).

8    12    Gui, M. *et al.* Cryo-electron microscopy structures of the SARS-CoV  
9        spike glycoprotein reveal a prerequisite conformational state for  
10       receptor binding. *Cell research* 27, 119-129, doi:10.1038/cr.2016.152  
11       (2017).

12   13    Zhao S, R. J., MUSA SS, Yang G, Lou Y, Gao D, Yang L, He D. .  
13       Preliminary estimation of the basic reproduction number of novel  
14       coronavirus (2019-nCoV) in China, from 2019 to 2020: A data-driven  
15       analysis in the early phase of the outbreak. . *bioRxiv*, 916395, doi:doi:  
16       <https://doi.org/10.1101/2020.01.23.916395> (2020).

17   14    Zhou, Y. *et al.* Protease inhibitors targeting coronavirus and filovirus  
18       entry. *Antiviral research* 116, 76-84, doi:10.1016/j.antiviral.2015.01.011  
19       (2015).

20   15    Millet, J. K. & Whittaker, G. R. Host cell proteases: Critical determinants  
21       of coronavirus tropism and pathogenesis. *Virus research* 202, 120-134,

- 1           doi:10.1016/j.virusres.2014.11.021 (2015).
- 2   16   Shirato, K., Kawase, M. & Matsuyama, S. Wild-type human
- 3           coronaviruses prefer cell-surface TMPRSS2 to endosomal cathepsins
- 4           for cell entry. *Virology* 517, 9-15, doi:10.1016/j.virol.2017.11.012
- 5           (2018).
- 6   17   Heurich, A. *et al.* TMPRSS2 and ADAM17 cleave ACE2 differentially
- 7           and only proteolysis by TMPRSS2 augments entry driven by the severe
- 8           acute respiratory syndrome coronavirus spike protein. *Journal of*
- 9           *virology* 88, 1293-1307, doi:10.1128/jvi.02202-13 (2014).
- 10   18   Herter, S. *et al.* Hepatocyte growth factor is a preferred in vitro substrate
- 11           for human hepsin, a membrane-anchored serine protease implicated in
- 12           prostate and ovarian cancers. *The Biochemical journal* 390, 125-136,
- 13           doi:10.1042/bj20041955 (2005).
- 14   19   Limburg, H. *et al.* TMPRSS2 Is the Major Activating Protease of
- 15           Influenza A Virus in Primary Human Airway Cells and Influenza B Virus
- 16           in Human Type II Pneumocytes. *Journal of virology* 93,
- 17           doi:10.1128/jvi.00649-19 (2019).
- 18   20   Holshue, M. L. *et al.* First Case of 2019 Novel Coronavirus in the United
- 19           States. *The New England journal of medicine*,
- 20           doi:10.1056/NEJMoa2001191 (2020).
- 21   21   Guan WJ, N. Z., Hu Y, Liang WH, Ou CQ, He JX, Liu L, Shan H, Lei CL,

1        Hui David S.C., Du B    Clinical characteristics of 2019 novel  
2        coronavirus infection in China. *medRxiv*, doi:doi:  
3        <http://dx.doi.org/10.1101/2020.02.06.20020974> (2020).

4    22    Perlman, S. & Netland, J. Coronaviruses post-SARS: update on  
5        replication and pathogenesis. *Nature reviews. Microbiology* 7, 439-450,  
6        doi:10.1038/nrmicro2147 (2009).

7    23    de Wit, E., van Doremalen, N., Falzarano, D. & Munster, V. J. SARS  
8        and MERS: recent insights into emerging coronaviruses. *Nature*  
9        *reviews. Microbiology* 14, 523-534, doi:10.1038/nrmicro.2016.81  
10        (2016).

11   24    Lee, P. I. & Hsueh, P. R. Emerging threats from zoonotic  
12        coronaviruses-from SARS and MERS to 2019-nCoV. *Journal of*  
13        *microbiology, immunology, and infection = Wei mian yu gan ran za zhi*,  
14        doi:10.1016/j.jmii.2020.02.001 (2020).

15   25    Li, W. *et al.* Angiotensin-converting enzyme 2 is a functional receptor for  
16        the SARS coronavirus. *Nature* 426, 450-454, doi:10.1038/nature02145  
17        (2003).

18   26    Kuiken, T. *et al.* Newly discovered coronavirus as the primary cause of  
19        severe acute respiratory syndrome. *Lancet* 362, 263-270,  
20        doi:10.1016/s0140-6736(03)13967-0 (2003).

21   27    Nabhan, A. N., Brownfield, D. G., Harbury, P. B., Krasnow, M. A. &

1        Desai, T. J. Single-cell Wnt signaling niches maintain stemness of  
2        alveolar type 2 cells. *Science (New York, N.Y.)* 359, 1118-1123,  
3        doi:10.1126/science.aam6603 (2018).

4    28   Barkauskas, C. E. *et al.* Type 2 alveolar cells are stem cells in adult lung.  
5        *The Journal of clinical investigation* 123, 3025-3036,  
6        doi:10.1172/jci68782 (2013).

7    29   Kroetz, D. N. *et al.* Type I Interferon Induced Epigenetic Regulation of  
8        Macrophages Suppresses Innate and Adaptive Immunity in Acute  
9        Respiratory Viral Infection. *PLoS pathogens* 11, e1005338,  
10        doi:10.1371/journal.ppat.1005338 (2015).

11   30   Shen, L. W., Mao, H. J., Wu, Y. L., Tanaka, Y. & Zhang, W. TMPRSS2:  
12        A potential target for treatment of influenza virus and coronavirus  
13        infections. *Biochimie* 142, 1-10, doi:10.1016/j.biochi.2017.07.016  
14        (2017).

15   31   Shin, W. J. & Seong, B. L. Type II transmembrane serine proteases as  
16        potential target for anti-influenza drug discovery. *Expert opinion on drug*  
17        *discovery* 12, 1139-1152, doi:10.1080/17460441.2017.1372417 (2017).

18   32   Biasini, M. *et al.* SWISS-MODEL: modelling protein tertiary and  
19        quaternary structure using evolutionary information. *Nucleic acids*  
20        *research* 42, W252-258, doi:10.1093/nar/gku340 (2014).

21   33   Pettersen, E. F. *et al.* UCSF Chimera--a visualization system for

1 exploratory research and analysis. *Journal of computational chemistry*  
2 25, 1605-1612, doi:10.1002/jcc.20084 (2004).

3 34 Wiehe, K. *et al.* ZDOCK and RDOCK performance in CAPRI rounds 3,  
4 4, and 5. *Proteins* 60, 207-213, doi:10.1002/prot.20559 (2005).

5 35 Duckert, P., Brunak, S. & Blom, N. Prediction of proprotein convertase  
6 cleavage sites. *Protein engineering, design & selection : PEDS* 17,  
7 107-112, doi:10.1093/protein/gzh013 (2004).

8 36 Madisson, E. *et al.* scRNA-seq assessment of the human lung, spleen,  
9 and esophagus tissue stability after cold preservation. *Genome biology*  
10 21, 1, doi:10.1186/s13059-019-1906-x (2019).

11 37 Vieira Braga, F. A. *et al.* A cellular census of human lungs identifies  
12 novel cell states in health and in asthma. *Nature medicine* 25,  
13 1153-1163, doi:10.1038/s41591-019-0468-5 (2019).

14 38 Reyfman, P. A. *et al.* Single-Cell Transcriptomic Analysis of Human  
15 Lung Provides Insights into the Pathobiology of Pulmonary Fibrosis.  
16 *American journal of respiratory and critical care medicine* 199,  
17 1517-1536, doi:10.1164/rccm.201712-2410OC (2019).

18 39 Valenzi, E. *et al.* Single-cell analysis reveals fibroblast heterogeneity  
19 and myofibroblasts in systemic sclerosis-associated interstitial lung  
20 disease. *Annals of the rheumatic diseases* 78, 1379-1387,  
21 doi:10.1136/annrheumdis-2018-214865 (2019).

1    40    Zhang, P. *et al.* Dissecting the Single-Cell Transcriptome Network  
2            Underlying Gastric Premalignant Lesions and Early Gastric Cancer.  
3            *Cell reports* 27, 1934-1947.e1935, doi:10.1016/j.celrep.2019.04.052  
4            (2019).

5    41    Martin, J. C. *et al.* Single-Cell Analysis of Crohn's Disease Lesions  
6            Identifies a Pathogenic Cellular Module Associated with Resistance to  
7            Anti-TNF Therapy. *Cell* 178, 1493-1508.e1420,  
8            doi:10.1016/j.cell.2019.08.008 (2019).

9    42    Smillie, C. S. *et al.* Intra- and Inter-cellular Rewiring of the Human Colon  
10           during Ulcerative Colitis. *Cell* 178, 714-730.e722,  
11           doi:10.1016/j.cell.2019.06.029 (2019).

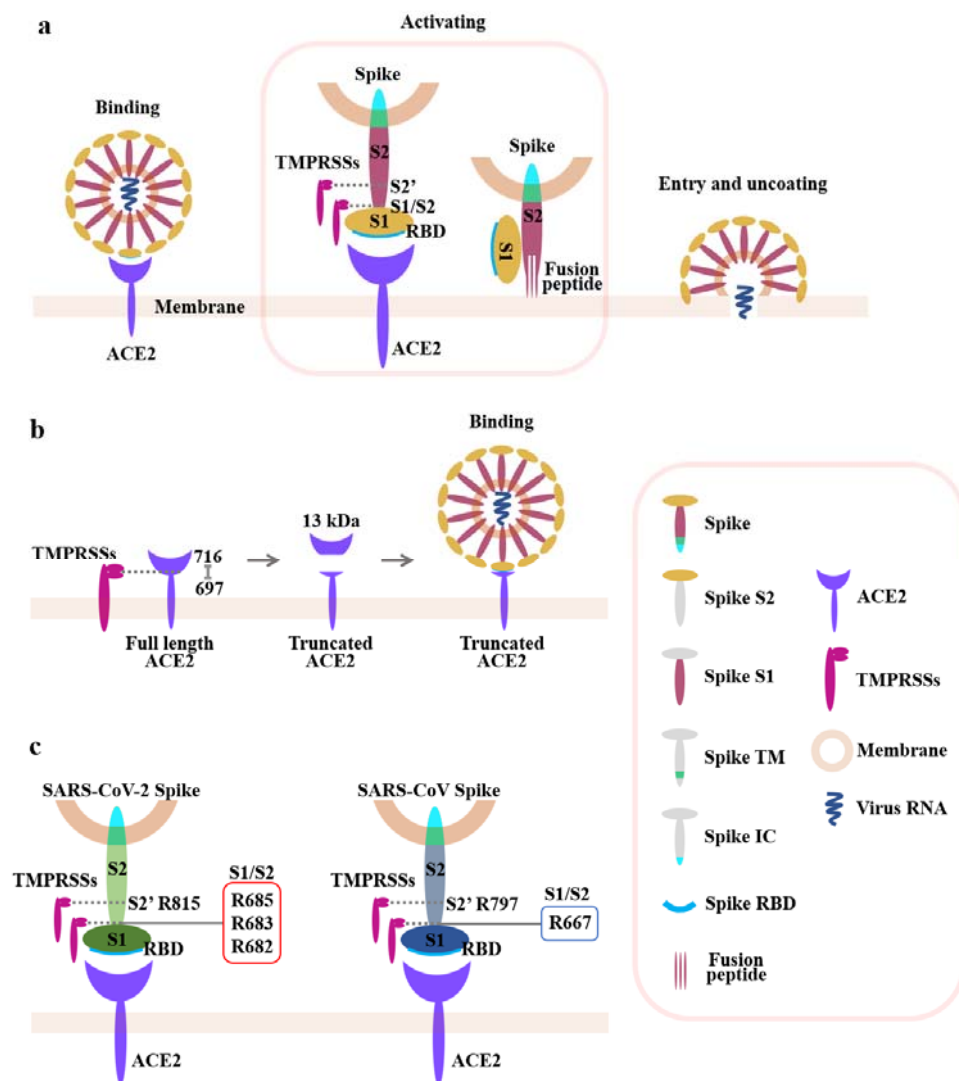
12   43    Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* 177,  
13           1888-1902.e1821, doi:10.1016/j.cell.2019.05.031 (2019).

14   44    Welch, J. D. *et al.* Single-Cell Multi-omic Integration Compares and  
15           Contrasts Features of Brain Cell Identity. *Cell* 177, 1873-1887.e1817,  
16           doi:10.1016/j.cell.2019.05.006 (2019).

17   45    Rouillard, A. D. *et al.* The harmonizome: a collection of processed  
18           datasets gathered to serve and mine knowledge about genes and  
19           proteins. *Database : the journal of biological databases and curation*  
20           2016, doi:10.1093/database/baw100 (2016).

21   46    Human genomics. The Genotype-Tissue Expression (GTEx) pilot

1 analysis: multitissue gene regulation in humans. *Science (New York,*  
2 *N. Y.)* **348**, 648-660, doi:10.1126/science.1262110 (2015).  
3 47 Uhlen, M. *et al.* Proteomics. Tissue-based map of the human proteome.  
4 *Science (New York, N. Y.)* **347**, 1260419, doi:10.1126/science.1260419  
5 (2015).  
6



7  
8 **Fig. 1 The schematic diagram of the project.**

9 a. The entry of SARS-CoV into host cells: The spike protein of SARS-CoV binds to

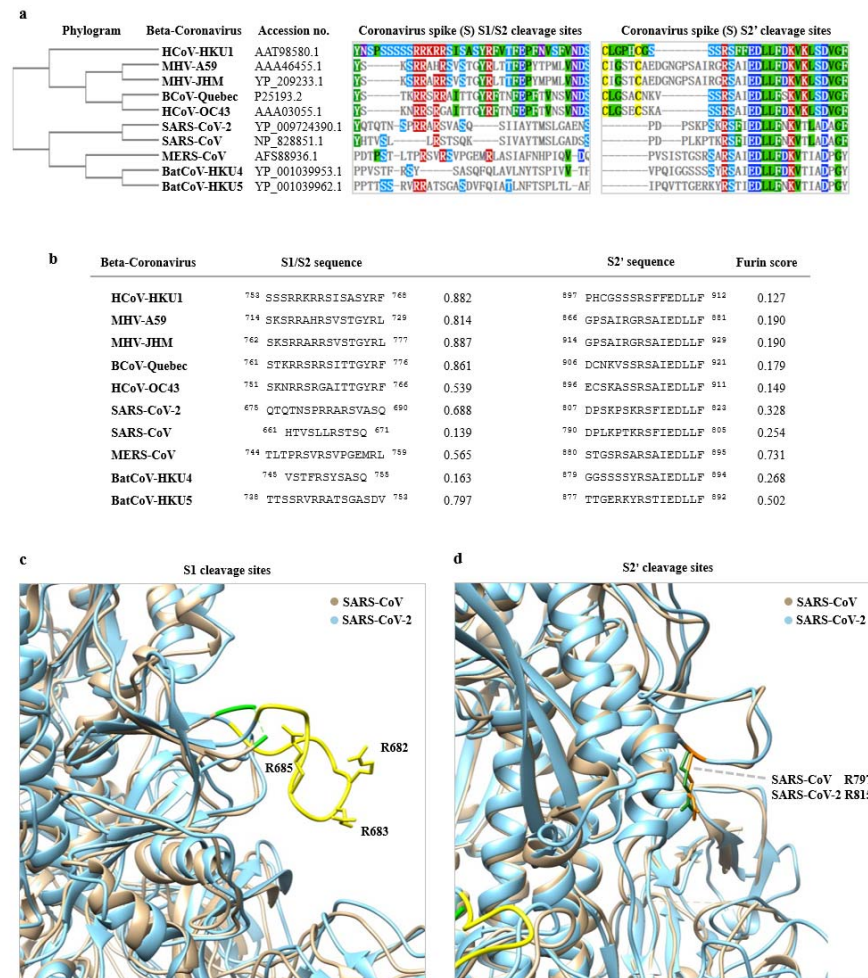


1 ACE2 through its S1 subunit for viral recognition. Then it is cleaved by  
2 TMPRSS2 at the S1/S2 boundary or within S2 subunit, which removes the  
3 structural constraint of S1 on S2, and releases the internal fusion peptide  
4 combined with the Spike TM domain for the fusion of viral and cellular  
5 membranes. Finally, the viral genomes enter into the host cells.

6 b. ACE2 cleaving by TMPRSSs: TMPRSS2 can also cleave ACE2 amino acids 697  
7 to 716, resulting in the shedding of 13kD ACE2 fragment in culture supernatants  
8 and augmented viral infectivity.

9 c. The difference between SARS-CoV-2 and SARS-CoV in the Spike protein  
10 cleavage: The Spike protein of SARS involves two cleavage sites recognized by  
11 TMPRSSs, one at arginine 667 and the other at arginine 797 (right). Compared  
12 with SARS-CoV, the Spike protein of SARS-CoV-2 (left) has an insertion  
13 sequence 680-SPRR-683 at the S1/S2 cleavage site. We speculated that R682,  
14 R683 and R685 (red box) could be used as the most suitable substrates for  
15 TMPRSSs, which can increase the Spike protein cleavage efficiency of TMPRSSs,  
16 promote its activation and enhance SARS-CoV-2 infection.

17



**Fig. 2 The two potential Spike protein cleavage sites of SARS-CoV and SARS-CoV-2 by TMPRSS2.**

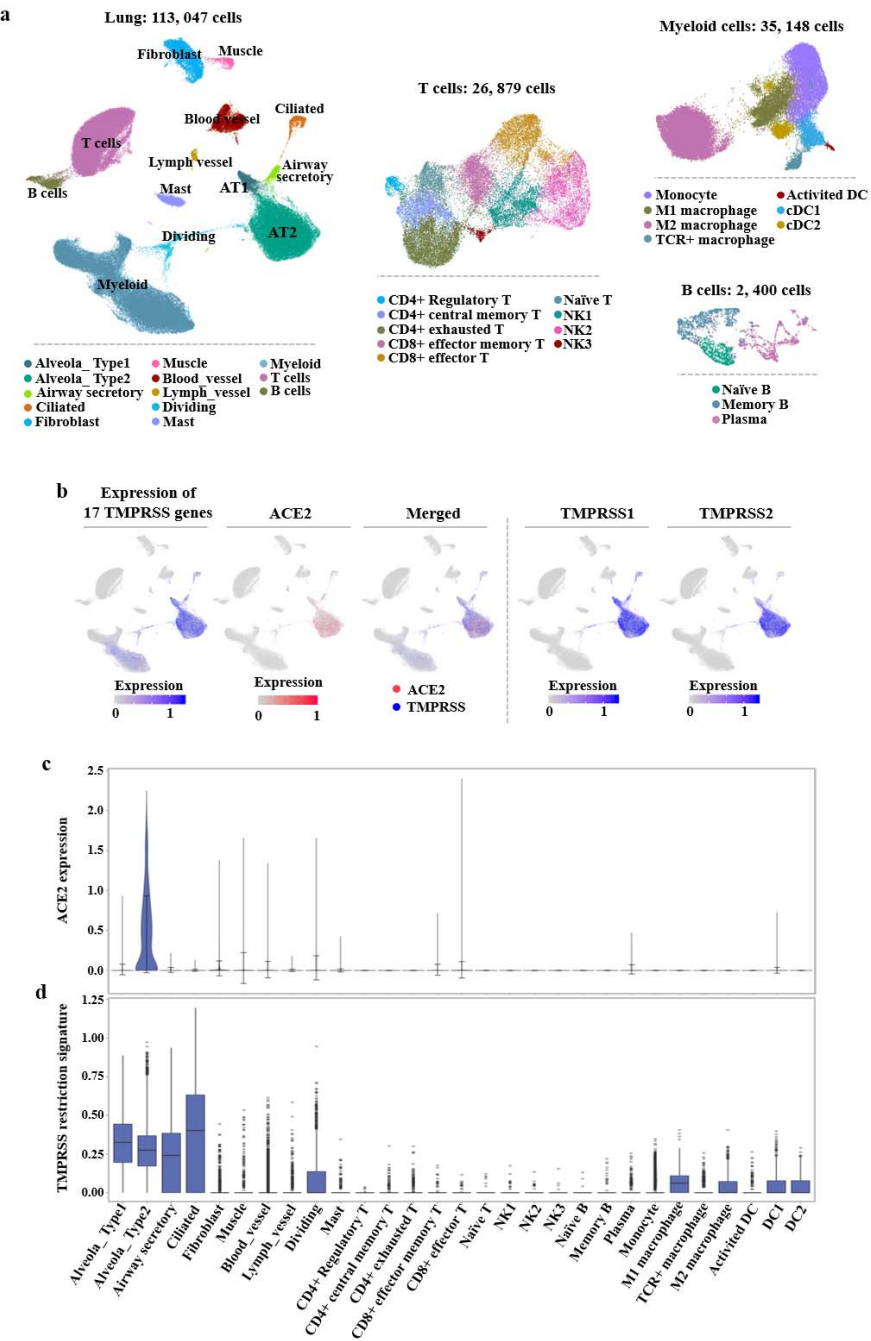
- Phylogenetic tree based on the protein sequences of Spike protein in SARS-CoV-2, SARS-CoV and other eight beta-coronaviruses are presented, along with the amino acid sequence alignment of two potential cleavage sites by TMPRSS2.
- The putative furin scores of the two potential cleavage sites of the ten coronaviruses.
- Structure comparison of the detailed Spike protein of the SARS-CoV and SARS-CoV-2. The insert 675-690 of SARS-CoV-2 Spike protein (yellow) and the corresponding loci to SARS-CoV Spike protein 661-672 (green). Three important

1 residues, R682, R683, R685, are specially marked.

2 d. The detail of c. The similarly SARS-CoV R797 with SARS-CoV-2 R815 are

3 marked with forest green and orange, respectively.

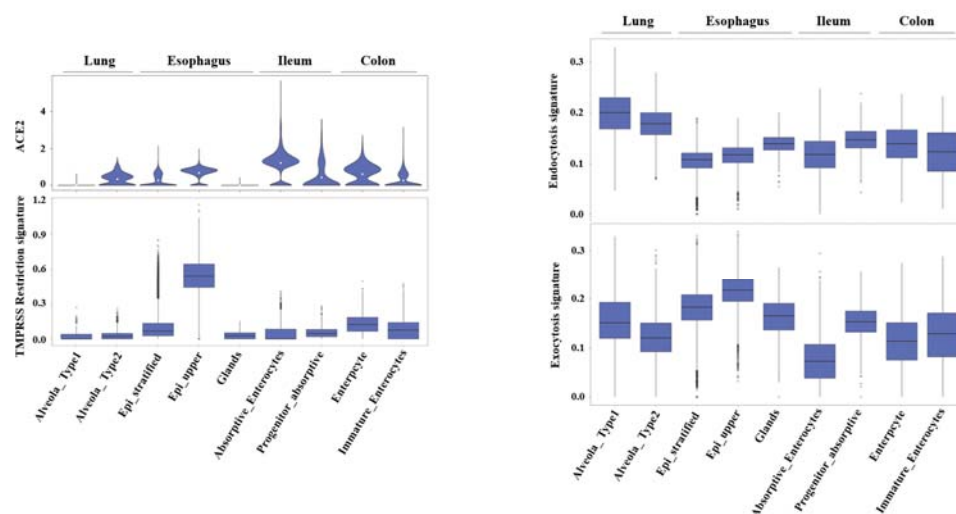
4



5

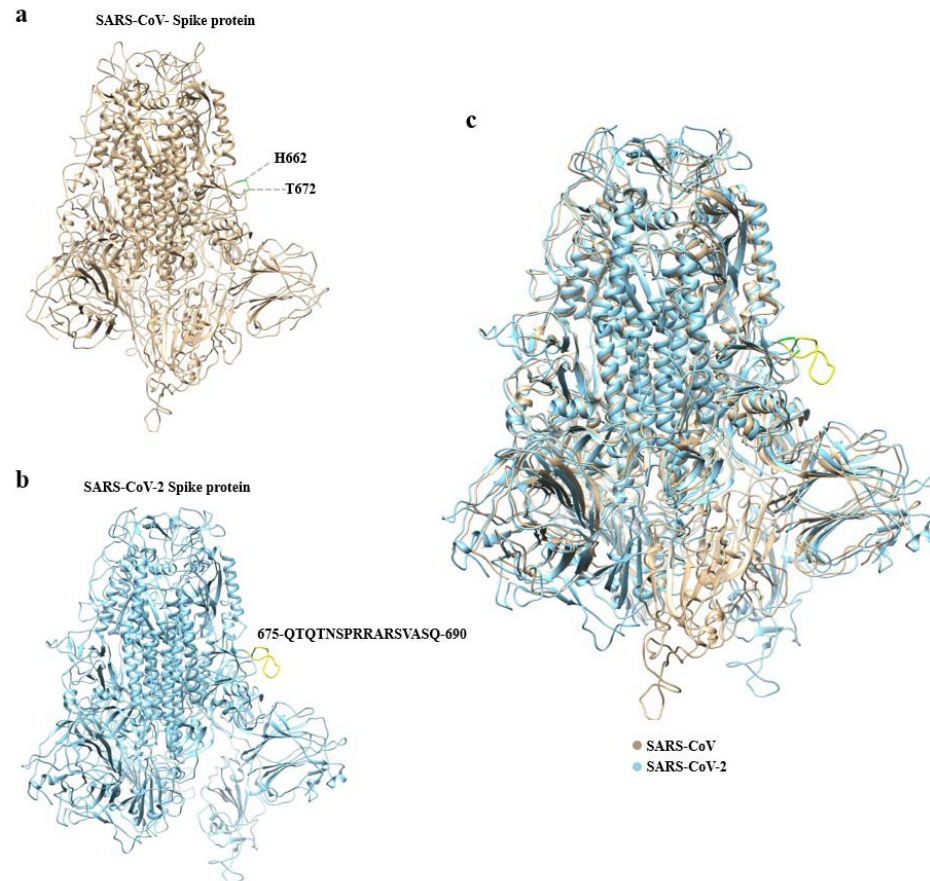
6 **Fig. 3 Single-cell analysis of the normal lung tissue.**

- 1 a. The UMAP plots of the landscape of lung cells. Thirteen clusters are colored,
- 2 distinctively labeled. T, B and myeloid cell subsets are further divided into finer
- 3 cell subsets according to the heterogeneity within the cell population.
- 4 b. The feature plots of the 17 TMPRSS genes, ACE2, TMPRSS1 and TMPRSS2.
- 5 c. The expression of ACE2 across clusters in the violin plot. The expression is
- 6 measured as the log2 (TP10K+1) value.
- 7 d. The mean expression of TMPRSS family genes across clusters in the boxplot. The
- 8 expression is measured as the mean log2 (TP10K+1) value.



9  
10 **Fig. 4 Expression levels of ACE2, TMPRSS restriction signature and functional**  
11 **gene sets in lung and digestive tracts.**

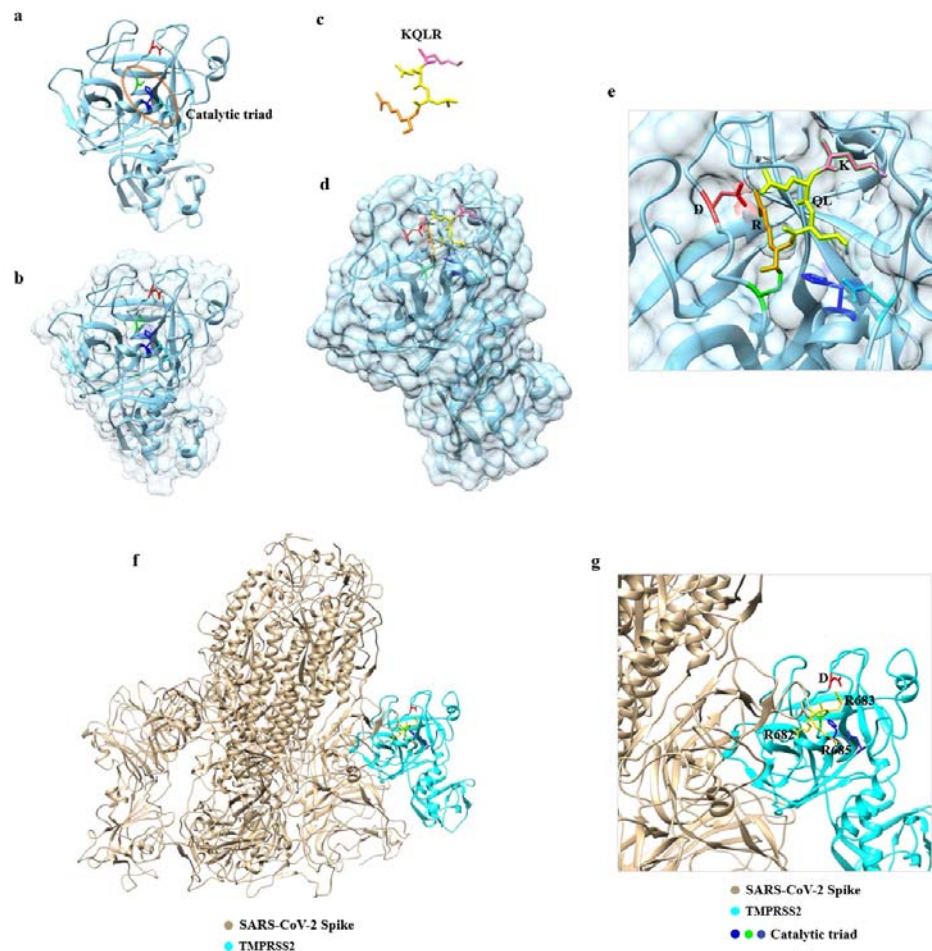
- 12 a. The expression levels of ACE2 and TMPRSS restriction signature in 2 lung
- 13 clusters and 7 digestive tract clusters. The expression is measured as the log2
- 14 (TP10K+1) value.
- 15 b. The expression levels of endocytosis and exocytosis-associated genes in 2 lung
- 16 clusters and 7 digestive tract clusters. The expression is measured as the log2
- 17 (TP10K+1) value.



**Extended Data Fig. 1 The overall structure of the Spike protein in SARS-CoV and SARS-CoV-2 homo-trimers**

- a. The 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 marked with green.
- b. The structure of the SARS-CoV-2 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 marked with yellow.
- c. The structural superimpose of Spike protein in the SARS-CoV (yellow) and SARS-CoV-2 (blue).





## Extended Data Fig. 2 The structure and catalytic mechanism of TMPRSS2

a-b. The overall structure and surface of TMPRSS2 (Modelled by SWISS-MODEL).

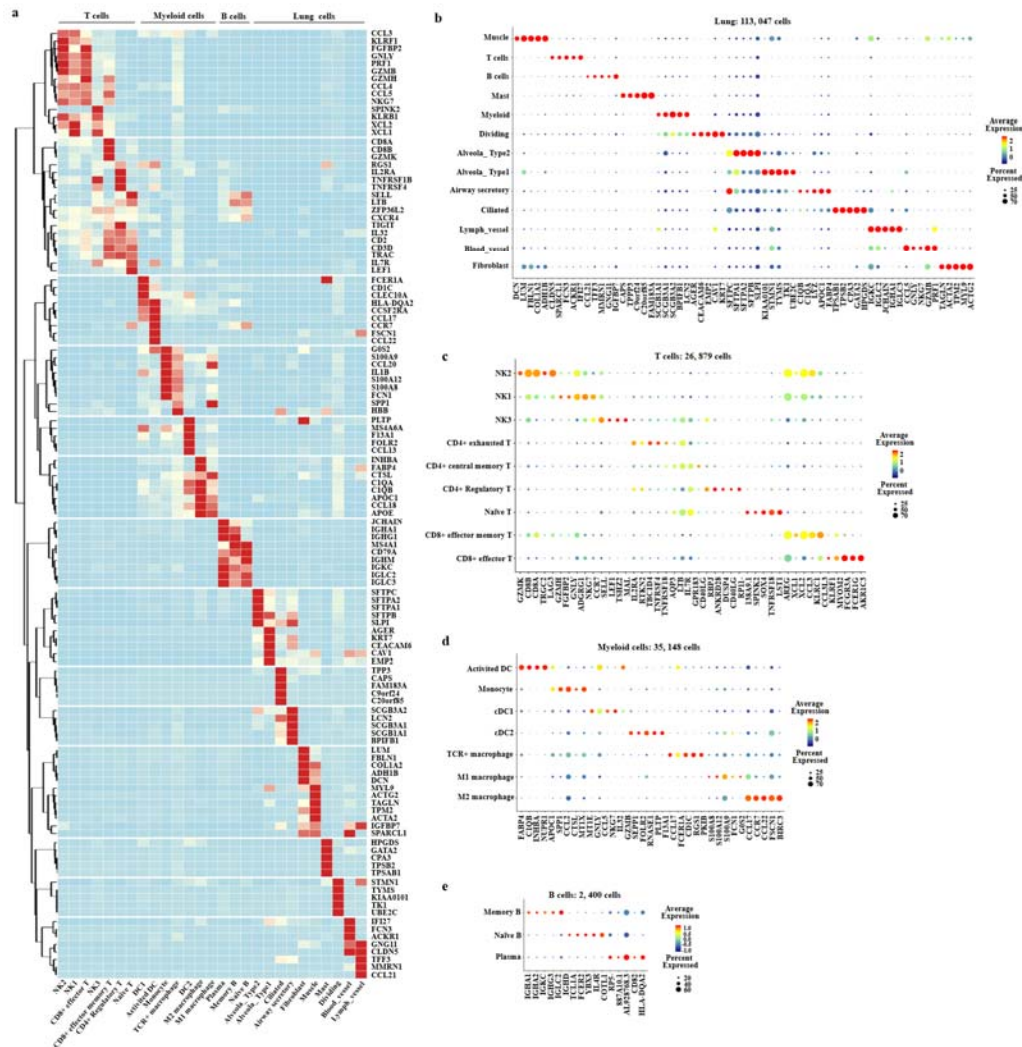
The TMPRSS2, catalytic triad comprised of H296, D345 and S441 are marked with cyan, blue, cyan and green, respectively. The substrate binding residue D435 located in the bottom of pocket is marked with red.

c. The polypeptide substrate analogue KQLR. The cleavage site Arg is marked with orange. Gln and Leu are marked with yellow. Lys is marked with pink.

d. The state of substrate analogue binding in the catalytic pocket. The state of substrate analogue binding in the catalytic pocket.

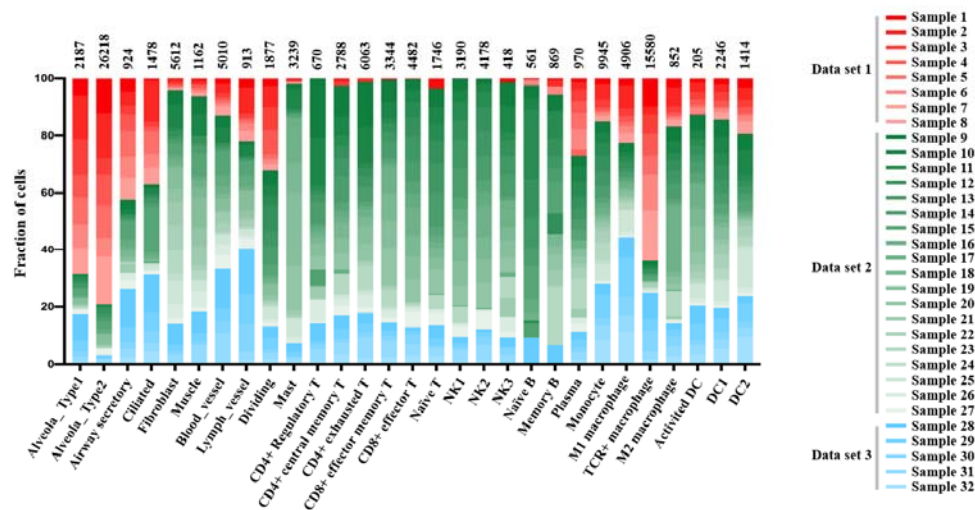
e. The detail of d. Arg of substrate analogue is strongly interacted with D435

- 1 f. The predicted state of SARS-CoV-2 Spike protein binding to the catalytic pocket of
- 2 TMPRSS2.
- 3 g. The detail of f. SARS-CoV-2 Spike protein and D345 of TMPRSS2 are marked
- 4 with wheat and medium blue, respectively.
- 5



6  
7 **Extended Data Fig. 3** Subset-specific markers.

- 8 a. The heatmap of marker genes (rows) across cell subsets (columns). The bubble
- 9 diagram of marker genes in thirteen clusters (b) and the sub-clusters of T cells (c),
- 10 B cells (d) and Myeloid cells (e).



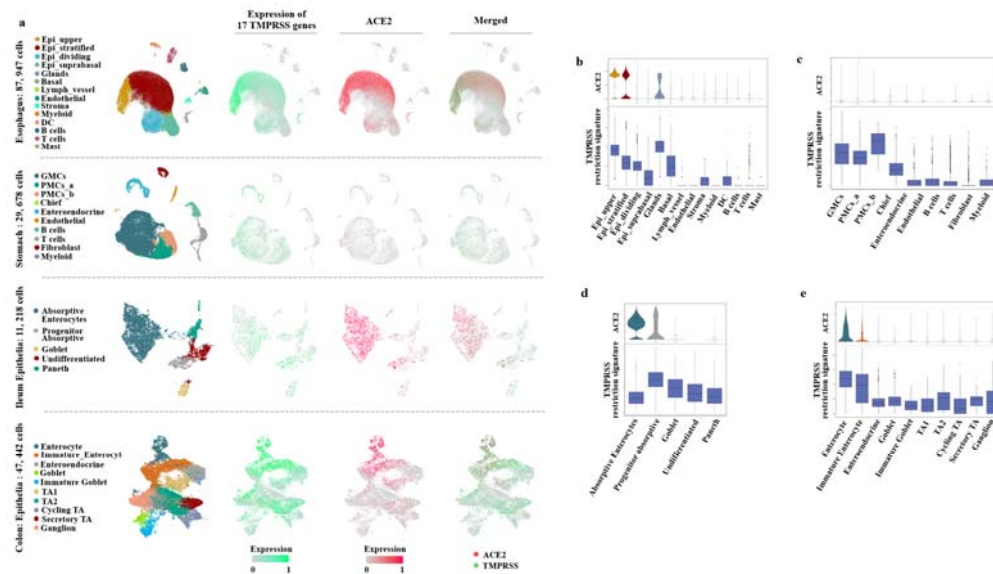
1

2 **Extended Data Fig. 4 All cell subset distributions across samples.**

3 The fractions of cells (y axis) in each cell subset (bars) that are derived from each  
4 sample in 3 databases (red, green and blue). The numbers of cells in each cluster are  
5 labeled above.



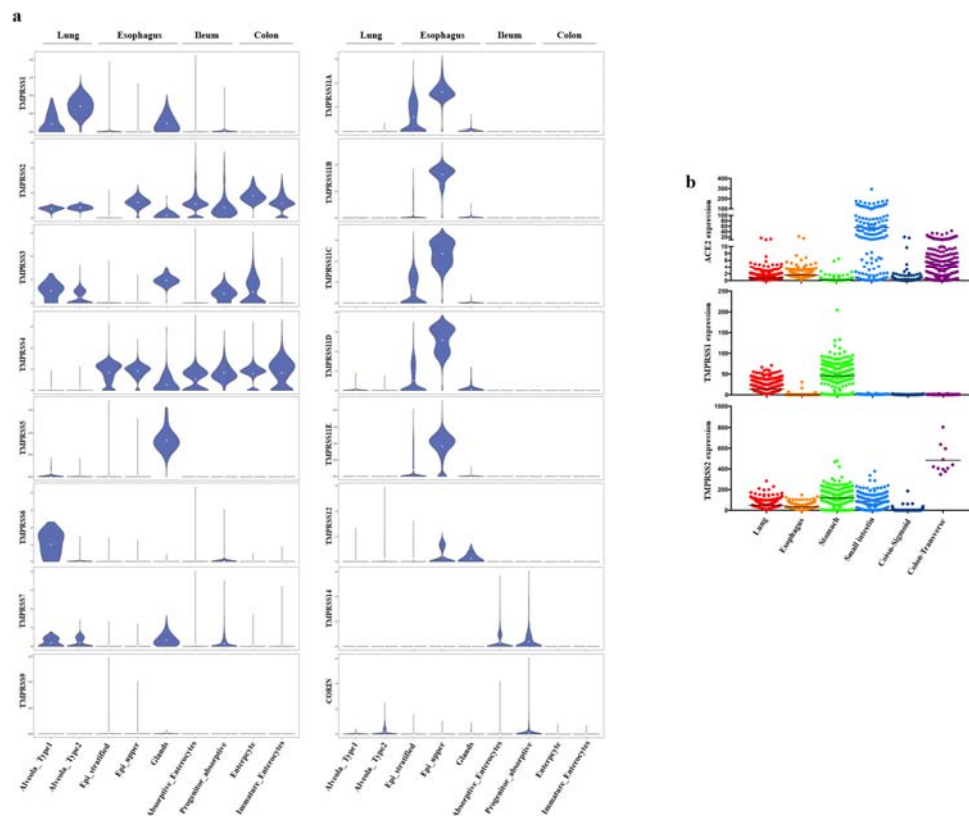




**Extended Data Fig. 6 The single-cell analysis of esophageal cells, gastric mucosal cells, ileal epithelial cells and colonic epithelial cells.**

- The UMAP plots of esophageal cells, gastric mucosal cells, ileal epithelial cells and colonic epithelial cells. The Feature plots show the expression of ACE2 (red) and TMPRSS family genes (green). The plots were merged to reveal the co-expression of these genes (brown).
- The expression levels ACE2 and TMPRSS restriction signature across clusters in esophagus (b), stomach(c), ileum(d) and colon(d). The expression is measured as the mean log<sub>2</sub> (TP10K+1) value.





## Extended Data Fig. 8 The expression levels of ASE2 and TMPRSS family genes in lung and digestive tracts

- The violin plots of TMPRSS family genes in lung and digestive tracts. The expression is measured as the mean log2 (TP10K+1) value.
- The expression levels of ACE2, TMPRSS1 and TMPRSS2 verified by RNA-seq data from the GTEx database.