

# 1        **Computational Deconvolution of Cell Type-Specific Gene** 2        **Expression in COPD and IPF Lungs Reveals Disease** 3        **Severity Associations**

4        Min Hyung Ryu<sup>1,3</sup>, Jeong H. Yun<sup>1,2,3</sup>, Kangjin Kim<sup>1,3</sup>, Michele Gentili<sup>1,3</sup>, Auyon Ghosh<sup>4</sup>, Frank  
5        Sciorba<sup>5</sup>, Lucas Barwick<sup>6</sup>, Andrew Limper<sup>7</sup>, Gerard Criner<sup>8</sup>, Kevin K. Brown<sup>9</sup>, Robert Wise<sup>10</sup>,  
6        Fernando J. Martinez<sup>11</sup>, Kevin R. Flaherty<sup>12</sup>, Michael H. Cho<sup>1,2,3</sup>, Peter J. Castaldi<sup>1,3,13</sup>, Dawn L.  
7        DeMeo<sup>1,2,3</sup>, Edwin K. Silverman<sup>1,2,3</sup>, Craig P. Hersh<sup>1,2,3,#</sup>, Jarrett D. Morrow<sup>1,3,#</sup>

8        <sup>1</sup> Channing Division of Network Medicine, and <sup>2</sup> Division of Pulmonary and Critical Care  
9        Medicine, <sup>13</sup> Division of General Internal Medicine and Primary Care, Department of Medicine,  
10        Brigham and Women's Hospital, Boston, Massachusetts

11        <sup>3</sup> Harvard Medical School, Boston, Massachusetts.

12        <sup>4</sup> Department of Medicine, Division of Pulmonary, Critical Care, and Sleep Medicine, SUNY  
13        Upstate Medical University, 750 East Adams Street, Syracuse, New York

14        <sup>5</sup> Division of Pulmonary, Allergy and Critical Care Medicine, University of Pittsburgh,  
15        Pittsburgh, Pennsylvania

16        <sup>6</sup> Emmes, Frederick, Maryland

17        <sup>7</sup> Division of Pulmonary and Critical Care Medicine, Department of Internal Medicine, Mayo  
18        Clinic, Rochester, Minnesota

19        <sup>8</sup> Thoracic Medicine and Surgery, Lewis Katz School of Medicine at Temple University,  
20        Philadelphia, Pennsylvania

21 <sup>9</sup> Department of Medicine, National Jewish Health, Denver, Colorado

22 <sup>10</sup> Department of Medicine, Johns Hopkins Medicine, Baltimore, Maryland

23 <sup>11</sup> Department of Medicine, Weill Cornell Medical College, NYPresbyterian Hospital, New  
24 York, New York

25 <sup>12</sup> Division of Pulmonary and Critical Care Medicine, University of Michigan Health System,  
26 Ann Arbor, Michigan

27 # co-senior authors.

28 **Corresponding Author:**

29 Craig P Hersh, MD, MPH

30 Channing Division of Network Medicine,

31 Brigham and Women's Hospital

32 181 Longwood Ave. Boston, MA 02115

33 [craig.hersh@channing.harvard.edu](mailto:craig.hersh@channing.harvard.edu)

34 Phone 617-525-0729

35 Fax 617-525-0958

36

37 **ABSTRACT**

38 Total word count: 349

39 **Rationale:** Chronic obstructive pulmonary disease (COPD) and idiopathic pulmonary fibrosis  
40 (IPF) are debilitating diseases associated with divergent histopathological changes in the lungs.  
41 At present, due to cost and technical limitations, profiling cell types is not practical in large  
42 epidemiology cohorts ( $n > 1000$ ). Here, we used computational deconvolution to identify cell  
43 types in COPD and IPF lungs whose abundances and cell type-specific gene expression are  
44 associated with disease diagnosis and severity.

45 **Methods:** We analyzed lung tissue RNA-seq data from 1026 subjects (COPD,  $n=465$ ; IPF,  
46  $n=213$ ; control,  $n=348$ ) from the Lung Tissue Research Consortium. We performed RNA-seq  
47 deconvolution, querying thirty-eight discrete cell-type varieties in the lungs. We tested whether  
48 deconvoluted cell-type abundance and cell type-specific gene expression were associated with  
49 disease severity.

50 **Results:** The abundance score of twenty cell types significantly differed between IPF and  
51 control lungs. In IPF subjects, eleven and nine cell types were significantly associated with  
52 forced vital capacity (FVC) and diffusing capacity for carbon monoxide ( $D_{LCO}$ ), respectively.  
53 Aberrant basaloid cells, a rare cells found in fibrotic lungs, were associated with worse FVC and  
54  $D_{LCO}$  in IPF subjects, indicating that this aberrant epithelial population increased with disease  
55 severity. Alveolar type 1 and vascular endothelial (VE) capillary A were decreased in COPD  
56 lungs compared to controls. An increase in macrophages and classical monocytes was associated  
57 with lower  $D_{LCO}$  in IPF and COPD subjects. In both diseases, lower non-classical monocytes  
58 and VE capillary A cells were associated with increased disease severity. Alveolar type 2 cells

59 and alveolar macrophages had the highest number of genes with cell type-specific differential  
60 expression by disease severity in COPD and IPF. In IPF, genes implicated in the pathogenesis of  
61 IPF, such as matrix metalloproteinase 7, growth differentiation factor 15, and eph receptor B2,  
62 were associated with disease severity in a cell type-specific manner.

63 **Conclusion:** Utilization of RNA-seq deconvolution enabled us to pinpoint cell types present in  
64 the lungs that are associated with the severity of COPD and IPF. This knowledge offers valuable  
65 insight into the alterations within tissues in more advanced illness, ultimately providing a better  
66 understanding of the underlying pathological processes that drive disease progression.

#### 67 **Keywords**

68 Chronic obstructive pulmonary disease, idiopathic pulmonary fibrosis, RNA sequencing,  
69 computational deconvolution, lung function tests, cell type-specific gene expression.

#### 70 **List of Abbreviations**

71 ATI, alveolar type 1 pneumocytes; ATII, alveolar type 2 pneumocytes; COPD, chronic  
72 obstructive pulmonary disease; DLCO, diffusing capacity for carbon monoxide; FACS,  
73 fluorescence-activated cell sorting; FDR, false discovery rate; FEV1, forced expiratory volume  
74 in one second; FVC, forced vital capacity; IIPs, idiopathic interstitial pneumonias; ILC A, type A  
75 innate lymphoid cells; ILD, interstitial lung disease; IPF, idiopathic pulmonary fibrosis; IQR,  
76 interquartile range; LTRC, Lung Tissue Research Consortium; PPI, protein-protein interaction;  
77 RNA-seq, RNA sequencing; SMC, smooth muscle cells; SMC, smooth muscle cells; TGF- $\beta$ ,  
78 transforming growth factor beta; VE Capillary A, vascular endothelial - aerocyte capillary; VE  
79 Capillary B, vascular endothelial - general capillary; VE Venous, vascular endothelial venous

80 cells; cMonocyte, classical monocytes; ncMonocyte, non-classical monocytes; pDC,

81 plasmacytoid dendritic cells, and; scRNA-seq, single-cell RNA sequencing

82

## 83 INTRODUCTION

84 Chronic obstructive pulmonary disease (COPD) and idiopathic pulmonary fibrosis (IPF) are  
85 debilitating chronic diseases of the lungs with progressive and complex pathobiology [1,2].  
86 COPD is characterized by airflow limitation, chronic airway inflammation, and lung  
87 parenchymal destruction [1]. IPF is characterized by cellular proliferation, interstitial  
88 inflammation, and fibrosis [2]. COPD and IPF are both related to long-term inhalation of noxious  
89 agents (e.g. tobacco smoking) and manifest in older adults as accelerated lung aging [3]. As  
90 such, both diseases are associated with significant morbidity, mortality, and a high economic  
91 burden to our society [4,5]. Therefore, there is an urgent need for disease prevention and  
92 improved treatments.

93 Genetics plays a role in predisposition to both diseases; eighty-two and nineteen loci have been  
94 associated with the risk of developing COPD or IPF, respectively [6,7]. COPD and IPF risk loci  
95 are enriched for pathways important in regulating cellular functions. For example, COPD risk  
96 loci are enriched for pathways regulating extracellular-matrix, cell-matrix adhesion, histone  
97 deacetylase binding, the Wnt-receptor signaling pathway, SMAD binding, and the MAPK  
98 cascade [6]. Similarly, IPF risk loci are enriched for pathways related to host defense, cell-cell  
99 adhesion, spindle assembly, transforming growth factor beta (TGF- $\beta$ ) signaling regulation, and  
100 telomere maintenance [8]. Furthermore, genetic factors are postulated to impact disease  
101 susceptibility in a cell type-specific and context specific manner. Therefore, improved molecular  
102 characterization of cells in the diseased lungs may provide insight into understanding disease  
103 pathobiology, paving the path to new therapeutics.

104 Investigating the molecular and cellular aspects of pathological lungs in the context of these  
105 diseases holds great promise for developing preventative and treatment strategies. In particular,

106 single-cell RNA sequencing (scRNA-seq) has been used in COPD and IPF patients to search for  
107 putative disease-causing cell types. For example, scRNA-seq analysis of IPF lungs has identified  
108 aberrant basaloid cells, a rare, disease-enriched cell type [9]. In COPD lungs, scRNA-seq has  
109 identified a high metallothionein-expressing macrophage subpopulation enriched in advanced  
110 COPD and altered bioenergetics and cellular stress tolerance in an alveolar type 2 pneumocyte  
111 (ATII) subpopulation [10]. A recent multi-omic single-cell analysis revealed a CD8<sup>+</sup> T cell  
112 subpopulation (KLRG1+TEMRA cells) to be enriched in COPD lung tissue [11]. However, the  
113 number of subjects included in these prior studies was modest, limiting the generalization to a  
114 larger patient population.

115 Due to the cost and technical limitations, performing scRNA-seq or tissue dissection experiments  
116 combined with fluorescence-activated cell sorting are yet to be practical in large epidemiology  
117 cohorts (n>1000). Moreover, the impact of tissue dissociation on gene expression in  
118 fluorescence-activated cell sorting (FACS) and scRNA-seq protocols remains poorly understood.  
119 Given that COPD and IPF are heterogeneous diseases, molecular studies encompassing a wide  
120 range of subjects with cell type-specific resolution are needed to unravel the complex interplay  
121 of cells in disease pathophysiology. To this end, large-scale clinical and genomic data in  
122 population cohorts may be leveraged to advance our search for cellular drivers of COPD and IPF  
123 pathogenesis.

124 In the present study, we performed computational deconvolution with bulk lung homogenate  
125 RNA-seq data from 1,026 subjects in the Lung Tissue Research Consortium (LTRC). By  
126 leveraging the large-scale omics data, we tested the hypothesis that there are specific cell types  
127 whose abundance and cell type-specific gene expression are associated with disease severity in  
128 COPD and IPF subjects.

## 129 **METHODS**

### 130 **Study participants**

131 Research subjects undergoing clinically indicated thoracic surgery were recruited to participate  
132 in the LTRC, as previously described [12]. The participating centers' Institutional Review  
133 Boards approved the study, and all subjects provided written informed consent.

134 COPD subjects included in this analysis had forced expiratory volume in one second (FEV<sub>1</sub>) to  
135 forced vital capacity (FVC) ratio <0.70 and FEV<sub>1</sub> % predicted <80%. Spirometric severity was  
136 characterized by Global Initiative for Chronic Obstructive Lung Disease spirometry grades 2-4.  
137 COPD subjects had either pathological emphysema and no alternative pathological diagnosis  
138 (interstitial lung disease (ILD), idiopathic interstitial pneumonias (IIPs), sarcoidosis, constrictive  
139 bronchiolitis, cellular hypersensitivity pneumonitis, diffuse alveolar damage, or eosinophilic  
140 granuloma). Any individual meeting the physiological diagnostic criteria for COPD but with a  
141 clinical diagnosis of IPF or sarcoidosis was excluded from the COPD group.

142 IPF subjects had a clinical diagnosis of IPF based on the site's multidisciplinary diagnostic  
143 process of all available data instituted at each participating institution. Control subjects had  
144 normal spirometry with no pathologic diagnosis of ILD/IIPs, sarcoidosis, constrictive  
145 bronchiolitis, cellular hypersensitivity pneumonitis, diffuse alveolar damage, or eosinophilic  
146 granuloma.

### 147 **Computational deconvolution**

148 Computational deconvolution was performed using CIBERSORTx (available at  
149 <https://cibersortx.stanford.edu/>) [13]. The docker image obtained from CIBERSORTx website



150 was used with the Podman container image management engine on the Channing Division of  
151 Network Medicine GPU computing cluster. This provided computational efficiency beyond what  
152 was available through the CIBERSORTx web interface.

153 We used LTRC TOPMed Harmonized phenotype data set dated November 30, 2022 and freeze 1  
154 LTRC gene expression data set. Data are available on the NCBI database of Genotypes and  
155 Phenotypes (dbGaP), accession phs001662 (LTRC). LTRC RNA-seq data from TOPMed  
156 (<https://topmed.nhlbi.nih.gov>) are available through dbGaP. For the count matrix generation,  
157 isoform-level expression quantification was generated with Salmon (v1.3.0) pseudoalignment to  
158 GENCODE release 37 transcriptome and summarized to gene-level counts using tximeta  
159 (v1.8.5). For salmon alignment, `seq_bias_correct` and `gc_bias_correct` were set to TRUE.  
160 Deconvolution was performed on the entire LTRC dataset that passed RNAseq QC (n=1,555),  
161 irrespective of whether the subject was included in our final analysis. Batch effects (library  
162 preparation batch) were removed using `Combat_seq` in the `sva` R package, and the matrix was  
163 cpm normalized after batch effect removal. Genes that had cpm >1 in at least 20% of the LTRC  
164 dataset and had assigned HUGO Gene Nomenclature Committee symbols were used in the  
165 deconvolution.

166 In total, 23,097 genes were included in the deconvolution after the batch effect removal and  
167 filtering steps. A custom signature matrix from a reference scRNA-seq was generated using  
168 CIBERSORTx. The signature matrix is a specialized expression matrix of cell type-specific  
169 “barcode” genes which provides a reference atlas of known cellular signatures for the  
170 deconvolution procedure. For this process, the CIBERSORTx algorithm used scRNA-seq data  
171 on 31,943 lung cells from 44 ever-smokers: six control, seventeen COPD, and twenty-one IPF  
172 subjects)[9,10]. Of the 23,097 genes in the LTRC dataset, 19,655 were also in the scRNA-seq

173 dataset (42,406 features across 161,067 cells in the qc'ed dataset); these genes were used to train  
174 the CIBERSORTx algorithm. The CIBERSORTx signature matrix we generated is attached as a  
175 supplementary file. Thirty-eight discrete cell varieties were queried in the deconvolution; cells  
176 were labeled as per Adams et al (Supplemental Table E1) [9]. We chose to use this dataset for  
177 two main reasons: 1.) the dataset included a wide range of control, COPD, and IPF subjects. 2.)  
178 the dataset included disease-specific cell types such as aberrant basaloid cells. Moreover, the cell  
179 annotations for the scRNA-seq were shown to be consistent with automated annotation drawn  
180 from multiple cell type definition databases such as the Human Primary Cell Atlas and Blue  
181 ENCODE databases, as previously reported [9].

182 For the imputation of cell fraction, we used CIBERSORTx in fraction mode with single-cell  
183 mode set to TRUE and rmbatchSmode set to FALSE; i.e., batch correction and quantile  
184 normalization by the CIBERSORTx algorithm were disabled. Proportions were calculated for  
185 each sample with all the cell types proportions added up to 1. For deriving abundance scores for  
186 cell types, the computation was performed on the CIBERSORTx web interface as this specific  
187 function is disabled by the algorithm provided by the authors inside the docker image.

188 CIBERSORTx estimates the relative fraction of each cell type included in the signature matrix,  
189 such that the sum of all fractions is equal to 1 for a given bulk RNA-seq sample. Therefore, the  
190 number of cell types included in the signature matrix may impact the relative fraction of each  
191 cell type. To overcome this issue, we used CIBERSORTx absolute mode where the absolute  
192 abundance score was estimated by the median expression level of all genes in the signature  
193 matrix (matrix generated using the reference scRNA-seq matrix) divided by the median  
194 expression level of all genes in the sample mixture (LTRC gene expression) [14,15]. This  
195 approach allows relative abundance comparisons across samples and cell types.

196 Cell type-specific gene expression matrices were generated using CIBERSORTx high-resolution  
197 mode using the docker image and used in the subsequent analyses.

### 198 **Cell type-specific differential gene expression analysis**

199 We performed differential gene expression analysis in cell type-specific gene expression  
200 matrices to find out which genes, even after removing the cellular abundance effects, were  
201 differentially expressed between case and controls. Using cell type-specific gene expression  
202 matrices (gene-by-sample matrices for each cell type) generated from CIBERSORTx, we  
203 performed differential gene expression analysis using limma [16]. Cell type-specific differential  
204 gene expression was  $\log_2$ -transformed, and we included only the genes with varying levels in our  
205 analysis (a built-in function of CIBERSORTx). We tested the association between cell type-  
206 specific gene expression and disease severity separately in the COPD and IPF groups. In COPD  
207 subjects, disease severity was measured by lung function tests including forced expiratory  
208 volume in 1 second (FEV<sub>1</sub>) and diffusing capacity of the lungs for carbon monoxide as a percent  
209 predicted (D<sub>L</sub>CO %). In IPF subjects, disease severity was measured by forced vital capacity  
210 (FVC) and D<sub>L</sub>CO %. Linear models were adjusted for age, sex, height, ever smoking, and  
211 lifetime smoking intensity (in pack-years). Multiple testing correction was performed by the  
212 Benjamini-Hochberg procedure. Significance was determined at a false discovery rate (FDR) of  
213 5%.

### 214 **Functional enrichment analysis**

215 We performed functional enrichment analysis using the STRING database version 12.0  
216 (<https://string-db.org>) [17]. The reason for using STRING was to use a complementary method  
217 based on publicly available dataset to explore the functional consequences of differentially

218 expressed genes. Alongside the protein-protein interaction, we also report gene set enrichment  
219 results performed using cell type-specific gene expression data which is part of the STRING  
220 interactive online platform.

221 Using the STRING interactive online platform, we queried active interaction sources and  
222 obtained confidence value in functional protein-protein interactions for protein network  
223 construction. We excluded any protein-protein interaction source that was based on text mining  
224 to reduce false positive signals. Active interaction sources include experiments, databases, co-  
225 expression, neighborhood, gene fusion, and co-occurrence. The list of genes used in the  
226 functional enrichment analysis are included in the Supplemental Table E2 and E3.

## 227 **RESULTS**

### 228 **Subjects**

229 465 subjects met the case criteria for COPD, 213 subjects met the case criteria for IPF, and 348  
230 subjects met the control criteria. Demographic and clinical characteristics of the 1,026 subjects  
231 included in our analysis are shown in [Table 1](#). Notably, IPF subjects were predominantly male  
232 (70%). The cohort included 90% of self-identified white subjects. COPD subjects were  
233 predominantly smokers (95.2% have ever smoked) and IPF and control subjects were 65.3% and  
234 67.8% ever smokers, respectively.

### 235 **Cellular composition differences among COPD, IPF, and controls**

236 Of the thirty-eight cell types queried in the deconvolution, twenty-seven cell types were detected  
237 in at least 10% of samples. Of these, there were nineteen cell types whose median proportion was  
238 greater than 1% in any one of the groups, as shown in [Figure 1](#).

239 We compared the cell abundance score between COPD, IPF, and control subjects, adjusting for  
240 age, sex, height, ever smoking, and smoking pack-years. [Figure 1](#) summarizes cell types whose  
241 abundance scores were significantly different ( $FDR < 0.05$ ) between COPD and control subjects  
242 and between IPF and control subjects, respectively. VE Capillary A and ATI were lower in  
243 COPD tissue compared to controls. Nine cell types were decreased and eleven were increased in  
244 IPF compared to controls.

### 245 **Associations between cell-type abundance and disease severity in COPD and IPF lungs**

246 Next, we identified cell types whose abundance scores in COPD and IPF lungs were associated  
247 with disease severity measured by FEV<sub>1</sub> (COPD), FVC (IPF), and D<sub>L</sub>CO (COPD and IPF). In  
248 COPD subjects, there were two and six cell types that were significantly associated with FEV<sub>1</sub>  
249 and D<sub>L</sub>CO, respectively ([Table 2](#)). In IPF subjects, there were eleven and nine cell types that  
250 were significantly associated with FVC and D<sub>L</sub>CO, respectively ([Table 2](#)). Decreases in the  
251 abundances of type A capillary vascular endothelial cells and non-classical monocytes were  
252 associated with worse disease severity in both COPD and IPF subjects. In IPF, aberrant basaloid  
253 cells showed the strongest association with both FVC and D<sub>L</sub>CO. In fact, we performed  
254 additional analysis testing the association between cell abundance score and GAP index [18], a  
255 mortality predictive score based on gender (G), age (A), and physiological measures (P; FVC,  
256 and D<sub>L</sub>CO) in IPF, and found that aberrant basaloid cells had one of the strongest associations  
257 with the index (Supplemental Table E4).

258 **Associations between cell type-specific gene expression and disease severity in COPD and**  
259 **IPF lungs**

260 We estimated cell type-specific gene expression for cell types whose median proportion was  
261 greater than 1%: ATII, Alveolar Macrophage, SMC, Fibroblast, ATI, Myofibroblast, VE  
262 Capillary B, B Plasma, VE Capillary A, ILC A, VE Venous, Pericyte, and T Cytotoxic. [Table 3](#)  
263 summarizes the number of differentially expressed genes in COPD and IPF. Overall, there were  
264 more differentially expressed genes (FDR<0.05) in IPF lungs than in COPD lungs. ATII cells  
265 and alveolar macrophages were two cell types with the greatest number of genes with cell type-  
266 specific differential gene expression associated with disease severity in both diseases. Aberrant  
267 basaloid cells, despite being estimated to represent only 1.3 % (IQR: 0-3.5 %) of cell proportion  
268 in IPF subjects, had the second largest number of cell type-specific genes whose expression was  
269 positively associated with IPF severity.

270 Next, we tested the association between cell type-specific gene expression and disease severity in  
271 COPD and IPF subjects. We included all cell types whose median proportion was greater than  
272 1% in each disease group. In COPD subjects, cell types tested were Alveolar Macrophage, ATI,  
273 ATII, B Plasma, Fibroblast, ILC A, Myofibroblast, Pericyte, SMC, T Cytotoxic, VE Capillary A,  
274 VE Capillary B, and VE Venous. In IPF subjects, cell types tested included Aberrant Basaloid,  
275 Alveolar Macrophage, ATI, ATII, B Plasma, Fibroblast, ILC A, Myofibroblast, Pericyte, SMC,  
276 T Cytotoxic, VE Capillary B, and VE Venous. [Figure 2](#) (and [Figure E2](#)) shows the number of  
277 genes with cell type-specific expression associated with lung function measures in COPD and  
278 IPF subjects. We also provide a list of all cell type-specific gene expression associations with  
279 disease severity in IPF and COPD ([Supplemental Table E5](#) and [6](#)). [Figure 2](#) also shows the  
280 number of genes with cell that overlap between the two different measures of disease severity.

281 Supplemental Tables E7 and E8 summarize the number of significant cell type-specific gene  
282 expressions associated with disease severity in COPD and IPF, respectively. Of note, besides the  
283 ATII cells, which were the most abundant cell types in the samples estimated using RNA-seq  
284 deconvolution, alveolar macrophages in COPD and aberrant basaloid cells had the highest  
285 number of genes associated with disease severity in both COPD and IPF. Hence, we chose these  
286 two cell types to perform functional enrichment analyses and highlight their upregulated  
287 function.

### 288 **Functional enrichment analysis of genes associated with COPD severity in alveolar** 289 **macrophages**

290 We performed functional enrichment analysis using the list of genes whose expression levels in  
291 alveolar macrophages were positively associated with COPD severity as measured by FEV<sub>1</sub> and  
292 DLCO. We queried all matched proteins encoded by the 77 genes identified in this cell type-  
293 specific differential gene expression analysis. In the protein-protein interaction (PPI) network  
294 analysis in the STRING database, we found significant functional enrichment with 144 edges  
295 (expected number of edges 60; PPI enrichment p-value  $<1 \times 10^{-16}$ ). [Figure 3](#) shows the PPI  
296 network for proteins encoded by the alveolar macrophage gene expression that is positively  
297 associated with COPD severity. The result of the functional enrichment analysis is included in  
298 the online Supplement Table E8. The most significantly enriched term was from the Reactome  
299 database for Eukaryotic Translation Elongation (Reactome term HSA-156842: FDR =  $1.25 \times 10^{-$   
300 <sup>11</sup>).

301 **Functional enrichment analysis of genes associated with IPF severity in aberrant basaloid**  
302 **cells**

303 We performed functional enrichment analysis using the list of genes whose expression levels in  
304 aberrant basaloid cells were positively associated with IPF severity as measured by FVC and  
305 DLCO. We queried all matched proteins encoded by the 185 genes identified in the cell type-  
306 specific differential gene expression analysis. We found significant functional enrichment with  
307 123 edges (expected number of edges 53; PPI enrichment p-value =  $2.22 \times 10^{-16}$ ). Figure 4 shows  
308 the PPI network for proteins encoded by the aberrant basaloid genes positively associated with  
309 IPF severity. The result of the functional enrichment analysis is included in the Online  
310 Supplement Table E9. Formation of the cornified envelope (STRING Cluster ID CL34114; FDR  
311 =  $5.85 \times 10^{-14}$ ) was indicated as the top most significant functional enrichment term.

312 **DISCUSSION**

313 We report the results of a computational tissue profiling analysis of bulk lung RNA-seq data  
314 from 1,026 subjects in the LTRC. We report the cellular composition and cell type-specific gene  
315 expression in lung tissue associated with disease severity in COPD and IPF subjects, extending  
316 the single-cell experiment discoveries from a modest sample size (<100 subjects) to a large  
317 population cohort (>1000 subjects). We trained a well-established and widely implemented  
318 computational RNA-seq deconvolution algorithm, CIBERSORTx [13,19,20], using publicly  
319 available scRNA-seq data from control, COPD, and IPF subjects [9].

320 We found that IPF lung tissues showed the most divergence from control lungs in cellular  
321 composition, with eighteen cell types whose abundance score was different from the controls,  
322 adjusting for covariates. Our results showed in a large IPF sample the association of aberrant



323 basaloid cells and their expression with IPF and IPF severity; the association with IPF severity  
324 has not been previously reported. We also found that abundances of eight cell types—  
325 ncMonocyte, Aberrant Basaloid, Macrophage, cMonocyte, T Cytotoxic, ATII, Alveolar  
326 Macrophage, and VE Capillary A—were associated with disease severity in the IPF subjects.  
327 Structural cells such ATII, aberrant basaloid cells, myofibroblasts, and fibroblasts were among  
328 the cell types with the most number of genes associated with IPF severity. Notably, we found  
329 that aberrant basaloid cells were enriched in IPF lungs, and that the abundance of this disease-  
330 enriched cell type increased as the disease severity increased. It is notable that aberrant basaloid  
331 proportions remained below 1% in COPD.

332 In aberrant basaloid cells, expression levels of matrix metalloproteinase 7 (*MMP7*), growth  
333 differentiation factor 15 (*GDF15*), and eph receptor B2 (*EPHB2*), were negatively associated  
334 with FVC or DLCO. In other words, the expression of these genes increased in more severe  
335 disease. These genes and the protein they encode have been implicated in the pathogenesis of  
336 IPF [21–24]. Our data supports the notion that GDF15 may be circulating biomarker reflective of  
337 aberrant basaloid cells in the airway epithelium [23]. We also found that *EPHB2* level in  
338 myofibroblasts was positively associated with IPF severity, extending the previous scRNA-seq  
339 finding that demonstrated increase level of *EPHB2* in IPF subjects compared to controls [9].

340 The functional enrichment analysis showed that the formation of the cornified envelope and  
341 keratinization were functionally enriched in aberrant basaloid cells with increasing severity of  
342 IPF. The cornified cell envelope is a highly insoluble and extremely tough structure that forms  
343 under the epithelium to help the epithelium defend against reactive oxygen species [25]. This  
344 may result from and/or be a contributing factor to the tissue fibrosis in IPF; however, alteration  
345 in this cellular function has not been implicated in IPF previously. Therefore, this result will

346 require further validation at the protein level. In addition to these functions, the protein  
347 interaction network analysis also highlighted the increased expression of matrix metalloproteases  
348 such as *MMP7*, *MMP10*, and *MMP1*, along with their functionally associated genes such as  
349 lipopolysaccharide binding protein (*LBP*), lipocalin 2 (*LCN2*), and transcobalamin1 (*TCNI*), in  
350 aberrant basaloid cells with increased disease severity. These results suggest that increased  
351 abundance of aberrant basaloid cells and their gene expression of cellular processes involved in  
352 aberrant barrier formation and extracellular matrix modification is associated with IPF severity.

353 We also showed that cellular composition is different between COPD and controls and that there  
354 were several cell types whose abundance was associated with COPD severity. There was a  
355 significant decrease in alveolar type 1 cells and capillary type A vascular endothelial cells in  
356 COPD lungs compared to controls. Capillary type A vascular endothelial cells were also  
357 negatively associated with increasing disease severity as measured by FEV<sub>1</sub> and D<sub>L</sub>CO. This  
358 observation provides additional evidence linking endothelial injury to COPD and extends earlier  
359 findings that identified injury to pulmonary vessels in lung tissue from COPD patients [26].

360 Beyond the pulmonary vasculature, the abundance of macrophage, ncMonocyte, and cMonocyte  
361 were associated with D<sub>L</sub>CO, but only ncMonocytes abundance was significantly associated with  
362 FEV<sub>1</sub>.

363 Monocytes and macrophages play an important role in pulmonary host defenses through their  
364 phagocytic activities and regulation of innate and adaptive immunity. The circulating monocyte  
365 pool and macrophages in tissue are composed of multiple subsets, each with a specialized  
366 function. Animal models and human *ex vivo* experiments have demonstrated the dysregulated  
367 functions of macrophage populations in COPD lungs [27]. Extensive molecular characterizations  
368 of immune cells in COPD, particularly the lung macrophage populations, have been conducted

369 using flow cytometry and other low-throughput molecular techniques [28–30]. However, due to  
370 the practicality of needing fresh samples and the experimental cost, tissue and immune profiling  
371 studies have been limited in terms of sample sizes (typically <100 subjects) and the small  
372 number of molecular targets. Recently, scRNA-seq studies with more molecular targets have  
373 been conducted and highlighted immunological dysregulation of monocytes and macrophages in  
374 COPD [9,10,31,32]. However, the number of COPD donors was small in these studies, and there  
375 was limited information on the disease phenotypes, which limited the ability to test for  
376 associations with disease severity, clinical outcomes, and pathological changes. Our  
377 computational tissue profiling in a large-scale cohort builds on this important body of work and  
378 extends the findings from scRNA-seq to an epidemiological cohort.

379 Given the important role alveolar macrophages play in COPD pathogenesis, we focused on this  
380 cell type for functional enrichment analysis, which highlighted that increased disease severity  
381 was associated with increased mRNA encoding for proteins involved in translation and energy  
382 metabolism. This finding agrees with previous studies that macrophage metabolic function is  
383 associated with COPD and supports the notion that metabolomic reprogramming of lung  
384 macrophages is important in the pathogenesis of COPD [33,34]. We provide the list of cell type-  
385 specific genes associated with COPD and IPF severity (Supplemental Table E5 and E6) for the  
386 community to explore using the cell type-specific functional enrichment using tools such as  
387 STRING database (<https://string-db.org/>) for other cell types.

388 There are some limitations of our study. First, RNA-seq based deconvolution methods are more  
389 suited for analysis of highly abundant cell types (cell types with frequency >1%) [13,19]. It is  
390 also influenced by the size of the cell type-specific transcriptome. This makes rare cell types with  
391 small transcriptomes challenging to study using the deconvolution approach. To overcome this

392 issue, future studies may combine RNA-seq deconvolution with results based on other omics  
393 (e.g., DNA methylation-based deconvolution). Also, careful enrichment of a cell type by FACS  
394 sorting may be required to study rarer cell populations. Second, bulk tissue analysis is limited in  
395 spatial resolution. This limits the understanding of the spatial distribution and interaction of cells  
396 in the diseased lungs. Nevertheless, our study informs which cell types may be the better  
397 candidates to be the focus of future spatial transcriptomic investigations. Finally, the study was  
398 limited to a population of predominantly white subjects with access to U.S. academic medical  
399 centers. This may limit the generalizability and calls for future efforts to include subjects from  
400 multi-ethnic and multi-national backgrounds.

## 401 **CONCLUSION**

402 In conclusion, we present here the cellular composition changes and cell type-specific gene  
403 expression associated with disease severity in COPD and IPF lungs. We document the cell types  
404 whose estimated abundance is associated with the severity of disease in COPD or IPF. We  
405 highlight two cell types—alveolar macrophages in COPD and aberrant basaloid cells in IPF—  
406 whose cell type-specific gene expressions were associated with clinical measures of disease  
407 severity. We also highlight the cell type-specific functional enrichment pointing to the altered  
408 cellular functions associated with disease severity. Using computational deconvolution, this  
409 study extends single-cell experimental discoveries from a modest sample size to a large  
410 population cohort and contributes to our understanding of tissue heterogeneity in COPD and IPF  
411 pathobiology. This knowledge offers insight into the alterations within lung tissue in advanced  
412 illness, providing a better understanding of the underlying pathological processes that drive  
413 disease progression.

414 **DECLARATIONS**

415 **Ethics approval and consent to participate**

416 The participating centers' Institutional Review Boards approved the study, and all subjects  
417 provided written informed consent.

418 **Consent for publication**

419 Not applicable.

420 **Availability of data and materials**

421 Data are available on the NCBI database of Genotypes and Phenotypes (dbGaP), accession  
422 phs001662 (LTRC). LTRC RNA-seq data from TOPMed (<https://topmed.nhlbi.nih.gov>) are  
423 available through dbGaP. The analysis results and code can be obtained by contacting the  
424 corresponding author with a reasonable request.

425 **Competing Interests**

426 Dr. Hersh reports grant support from Bayer, Boehringer-Ingelheim, and Vertex, and consulting  
427 fees from Chiesi, Sanofi, and Takeda, unrelated to this manuscript. Dr. Silverman reports grant  
428 support from Bayer and Northpond Laboratories. Dr. Cho reports grant support from Bayer.  
429 Dr. DeMeo reports grant support from Bayer and Alpha-1 Foundation. Dr. Castaldi reports grant  
430 support from Bayer, Sanofi and consulting fees from Verona Pharmaceuticals. Dr. Yun reports  
431 grant support from Bayer and consulting fees from Bridge Biotherapeutics, and travel  
432 reimbursement from the Korean Academy of Tuberculosis and Respiratory Disease unrelated to  
433 this manuscript. Dr. Flaherty reports grant funding from Boehringer Ingelheim unrelated to this  
434 manuscript. Dr. Martinez reports grant supports from NHLBI, AstraZeneca, Chiesi, Boehringer-

435 Ingelheim, GalaxoSmithCline, Novartis, Polarean, Sanofi/Regeneron, Sunovion, and TEVA  
436 Pharmaceuticals. Dr. Martinez reports receiving consulting fee from AstraZeneca, Boehringer-  
437 Ingelheim and Bristol Myers Squibb. Dr. Wise reports receiving consulting fees from  
438 Boehringer-Ingelheim, AstraZeneca, Abb-Vie, and Galderma.

#### 439 **Funding**

440 Present work was supported by grants from NHLBI (R01HL166231, P01HL114501,  
441 R01HL133135, and X01HL139404), K25 HL136846, K08 HL146972, Alpha-1 Foundation  
442 Research Grant, and TOPMed Fellowship. MHC was supported by R01HL162813,  
443 R01HL153248, R01HL14.

#### 444 **Author contributions**

445 Concept and design: MHR, CPH, and JDM; data collection: JHY, FS, LB, AL, GC, KB, RW,  
446 FM, KF, MHC, PJC, DLD, EKS, CPH, and JDM; statistical support: MHR, KJK, MG, CPH,  
447 JDM; data analysis: MHR, JHY, KJK, MG, AG, and JDM; manuscript writing - draft: MHR,  
448 CPH, and JDM; manuscript writing - edit: all authors; funding: PJC, EKS, and CPH. All authors  
449 read and approved the final manuscript.

#### 450 **Acknowledgements**

451 NHLBI TOPMed: Lung Tissue Research Consortium Molecular data from the Trans-Omics in  
452 Precision Medicine (TOPMed) program was supported by the National Heart, Lung, and Blood  
453 Institute (NHLBI). RNASeq for “NHLBI TOPMed: Lung Tissue Research Consortium”  
454 (phs001662) was performed at the Northwest Genomics Center (HHSN268201600032I). Core  
455 support including centralized genomic read mapping and genotype calling, along with variant

456 quality metrics and filtering were provided by the TOPMed Informatics Research Center  
457 (3R01HL-117626-02S1; contract HHSN268201800002I). Core support including phenotype  
458 harmonization, data management, sample-identity QC, and general program coordination were  
459 provided by the TOPMed Data Coordinating Center (R01HL-120393; U01HL-120393; contract  
460 HHSN268201800001I). We gratefully acknowledge the studies and participants who provided  
461 biological samples and data for TOPMed.

462

463 **TABLES**

464 *Table 1: LTRC subject demographics and lung function tests*

	Control	COPD	IPF
n	348	465	213
Age (mean (SD))	61.51 (12.53)	63.35 (9.18)	63.55 (8.37)
Sex = Female sex (%)	211 (60.6)	210 (45.2)	64 (30.0)
Race (%)			
White	314 (90.2)	423 (91.0)	191 (89.7)
Asian	0 ( 0.0)	0 ( 0.0)	4 ( 1.9)
Black	22 ( 6.3)	29 ( 6.2)	8 ( 3.8)
Hispanic	10 ( 2.9)	9 ( 1.9)	4 ( 1.9)
Other race	2 ( 0.6)	4 ( 0.9)	6 ( 2.8)
BMI (mean (SD))	28.95 (5.97)	26.30 (5.22)	29.80 (5.44)
Ever Smoking (%)	215 (67.8)	415 (95.2)	128 (65.3)
Pack years of smoking (mean (SD))	20.12 (27.32)	47.16 (31.73)	18.84 (24.18)
FEV1/FVC (mean (SD))	0.77 (0.06)	0.45 (0.15)	0.83 (0.07)
FEV1 pp (mean (SD))	95.87 (12.61)	41.72 (20.27)	65.91 (19.06)
FVC pp (mean (SD))	96.09 (12.69)	68.36 (18.50)	60.28 (17.75)
DLCO % (mean (SD))	73.22 (15.44)	42.92 (18.87)	38.04 (19.59)

465 Abbreviations: COPD, chronic obstructive pulmonary disease; IPF, idiopathic pulmonary  
466 fibrosis; SD, standard deviation; BMI, body mass index; FEV<sub>1</sub>, forced expiratory volume in 1s;  
467 FVC, forced vital capacity; DLCO, diffusing capacity of the lungs for carbon monoxide as a  
468 percent predicted.

469



470 *Table 2: Cell-type transcriptome abundance score associated with disease severity in COPD and*  
 471 *IPF.*

Disease	Outcome	Cell type	Beta	95% CI	Adjusted p value
COPD	FEV <sub>1</sub>	VE Capillary A	0.11	0.06,0.16	0.001
COPD	FEV <sub>1</sub>	ncMonocyte	0.09	0.03,0.14	0.02
COPD	D <sub>L</sub> CO	ATI	4.79	2.99,6.59	<0.001
COPD	D <sub>L</sub> CO	Macrophage	-4.75	-7.02,-2.49	<0.001
COPD	D <sub>L</sub> CO	ncMonocyte	3.71	1.94,5.49	<0.001
COPD	D <sub>L</sub> CO	VE Capillary A	3.69	1.95,5.44	<0.001
COPD	D <sub>L</sub> CO	cMonocyte	-3.16	-5.47,-0.86	0.031
COPD	D <sub>L</sub> CO	ILC A	2.72	0.93,4.51	0.016
IPF	FVC	ncMonocyte	0.31	0.21,0.4	<0.001
IPF	FVC	Aberrant Basaloid	-0.24	-0.34,-0.14	<0.001
IPF	FVC	Macrophage	-0.20	-0.31,-0.1	0.001
IPF	FVC	cMonocyte	-0.19	-0.29,-0.09	0.001
IPF	FVC	T Cytotoxic	0.19	0.09,0.29	0.002
IPF	FVC	ATII	0.16	0.06,0.27	0.007
IPF	FVC	VE Venous	-0.16	-0.27,-0.06	0.012
IPF	FVC	Alveolar Macrophage	0.15	0.05,0.25	0.016
IPF	FVC	VE Capillary A	0.15	0.05,0.26	0.016
IPF	FVC	T	-0.14	-0.24,-0.03	0.029
IPF	FVC	pDC	-0.12	-0.23,-0.02	0.046
IPF	D <sub>L</sub> CO	Aberrant Basaloid	-6.42	-9.25,-3.58	<0.001
IPF	D <sub>L</sub> CO	Alveolar Macrophage	6.08	3.19,8.97	0.001
IPF	D <sub>L</sub> CO	ATII	5.65	2.8,8.51	0.001
IPF	D <sub>L</sub> CO	T Cytotoxic	5.07	2.04,8.1	0.006
IPF	D <sub>L</sub> CO	VE Capillary A	5.06	2.1,8.03	0.006
IPF	D <sub>L</sub> CO	ncMonocyte	4.79	1.89,7.69	0.006
IPF	D <sub>L</sub> CO	Macrophage	-4.64	-7.64,-1.65	0.01
IPF	D <sub>L</sub> CO	cMonocyte	-4.34	-7.53,-1.15	0.023
IPF	D <sub>L</sub> CO	T Regulatory	3.92	1.06,6.77	0.023

472 Statistical comparison was tested using linear regression adjusting for age, sex, height, ever  
 473 smoking and total pack-year. Beta was estimated using absolute value of outcome measures and  
 474 are estimated per one standard deviation change in CIBERSORTx absolute abundance score.  
 475 Pre-bronchodilator FEV<sub>1</sub> and DLCO percent predicted were used. Abbreviations: COPD, chronic

476 obstructive pulmonary disease; IPF, idiopathic pulmonary fibrosis; CI, confidence interval;  
477 FEV<sub>1</sub>, forced expiratory volume in 1s; FVC, forced vital capacity; DLCO, diffusing capacity of  
478 the lungs for carbon monoxide as a percent predicted; VE Capillary A, vascular endothelial -  
479 aerocyte capillary; ncMonocyte, non-classical monocytes; cMonocyte, classical monocytes; ATI,  
480 alveolar epithelial type 1 cells; ILC A, type A innate lymphoid cells; ATII, alveolar epithelial  
481 type 2 cells; VE Venous, vascular endothelial venous cells; pDC, plasmacytoid dendritic cells.  
482

483 *Table 3: Cell type-specific differential gene expression in COPD and IPF lungs compared to*  
 484 *control lungs*

Cell type	Total genes in analysis	Number of upregulated genes in COPD	Number of downregulated genes in COPD	Number of upregulated genes in IPF	Number of downregulated genes in IPF
ATII	10964	272	2088	3886	2967
Alveolar Macrophage	5614	228	772	1847	1275
SMC	4265	40	183	1922	604
Fibroblast	3500	89	254	1268	658
ATI	3362	252	108	819	578
Myofibroblast	3099	120	399	1225	612
VE Capillary B	2446	54	338	412	943
B Plasma	2325	116	50	1039	460
VE Capillary A	2063	73	284	404	669
ILC A	2049	8	1	332	240
VE Venous	1310	33	114	230	374
Pericyte	1210	26	56	326	275
T Cytotoxic	1139	7	2	192	81

485 Associations were tested using limma [16] on variable genes only. Significant association were  
 486 adjusted to FDR 5%. Abbreviations: COPD, chronic obstructive pulmonary disease; IPF,  
 487 idiopathic pulmonary fibrosis; ATII, alveolar epithelial type 2 cells; SMC, smooth muscle cells;  
 488 ATI, alveolar epithelial type 1 cells; VE Capillary B, vascular endothelial - general capillary; VE  
 489 Capillary A, vascular endothelial - aerocyte capillary; ILC A, type A innate lymphoid cells; VE  
 490 Venous, vascular endothelial venous cells.

491

492 **REFERENCES**

493 [1] Agustí A, Hogg JC. Update on the pathogenesis of chronic obstructive pulmonary  
494 disease. *New England Journal of Medicine* 2019;381:1248–56.  
495 <https://doi.org/10.1056/nejmra1900475>.

496 [2] Lederer DJ, Martinez FJ. Idiopathic pulmonary fibrosis. *New England Journal of*  
497 *Medicine* 2018;378:1811–23. <https://doi.org/10.1056/nejmra1705751>.

498 [3] Selman M, Martinez FJ, Pardo A. Why does an aging smoker’s lung develop idiopathic  
499 pulmonary fibrosis and not chronic obstructive pulmonary disease? *American Journal of*  
500 *Respiratory and Critical Care Medicine* 2019;199:279–85. [https://doi.org/10.1164/rccm.201806-](https://doi.org/10.1164/rccm.201806-1166pp)  
501 [1166pp](https://doi.org/10.1164/rccm.201806-1166pp).

502 [4] Chen S, Kuhn M, Prettner K, Yu F, Yang T, Bärnighausen T, et al. The global economic  
503 burden of chronic obstructive pulmonary disease for 204 countries and territories in 202050:  
504 health-augmented macroeconomic modelling study. *The Lancet Global Health* 2023;11:e1183–  
505 93. [https://doi.org/10.1016/s2214-109x\(23\)00217-6](https://doi.org/10.1016/s2214-109x(23)00217-6).

506 [5] Wong AW, Koo J, Ryerson CJ, Sadatsafavi M, Chen W. A systematic review on the  
507 economic burden of interstitial lung disease and the cost-effectiveness of current therapies. *BMC*  
508 *Pulmonary Medicine* 2022;22. <https://doi.org/10.1186/s12890-022-01922-2>.

509 [6] Sakornsakolpat P, Prokopenko D, Lamontagne M, Reeve NF, Guyatt AL, Jackson VE, et  
510 al. Genetic landscape of chronic obstructive pulmonary disease identifies heterogeneous cell-  
511 type and phenotype associations. *Nature Genetics* 2019;51:494–505.  
512 <https://doi.org/10.1038/s41588-018-0342-2>.

- 513 [7] Allen RJ, Stockwell A, Oldham JM, Guillen-Guio B, Schwartz DA, Maher TM, et al.  
514 Genome-wide association study across five cohorts identifies five novel loci associated with  
515 idiopathic pulmonary fibrosis. *Thorax* 2022;77:829–33. [https://doi.org/10.1136/thoraxjnl-2021-](https://doi.org/10.1136/thoraxjnl-2021-218577)  
516 [218577](https://doi.org/10.1136/thoraxjnl-2021-218577).
- 517 [8] Allen RJ, Guillen-Guio B, Oldham JM, Ma S-F, Dressen A, Paynton ML, et al. Genome-  
518 wide association study of susceptibility to idiopathic pulmonary fibrosis. *American Journal of*  
519 *Respiratory and Critical Care Medicine* 2020;201:564–74. [https://doi.org/10.1164/rccm.201905-](https://doi.org/10.1164/rccm.201905-1017oc)  
520 [1017oc](https://doi.org/10.1164/rccm.201905-1017oc).
- 521 [9] Adams TS, Schupp JC, Poli S, Ayaub EA, Neumark N, Ahangari F, et al. Single-cell  
522 RNA-seq reveals ectopic and aberrant lung-resident cell populations in idiopathic pulmonary  
523 fibrosis. *Science Advances* 2020;6. <https://doi.org/10.1126/sciadv.aba1983>.
- 524 [10] Sauler M, McDonough JE, Adams TS, Kothapalli N, Barnthaler T, Werder RB, et al.  
525 Characterization of the COPD alveolar niche using single-cell RNA sequencing. *Nature*  
526 *Communications* 2022;13. <https://doi.org/10.1038/s41467-022-28062-9>.
- 527 [11] Villaseñor-Altamirano AB, Jain D, Jeong Y, Menon JA, Kamiya M, Haider H, et al.  
528 Activation of CD8<sup>+</sup> T cells in chronic obstructive pulmonary disease lung. *American Journal of*  
529 *Respiratory and Critical Care Medicine* 2023;208:1177–95.  
530 <https://doi.org/10.1164/rccm.202305-0924oc>.
- 531 [12] Yang IV, Pedersen BS, Rabinovich E, Hennessy CE, Davidson EJ, Murphy E, et al.  
532 Relationship of DNA methylation and gene expression in idiopathic pulmonary fibrosis.  
533 *American Journal of Respiratory and Critical Care Medicine* 2014;190:1263–72.  
534 <https://doi.org/10.1164/rccm.201408-1452oc>.

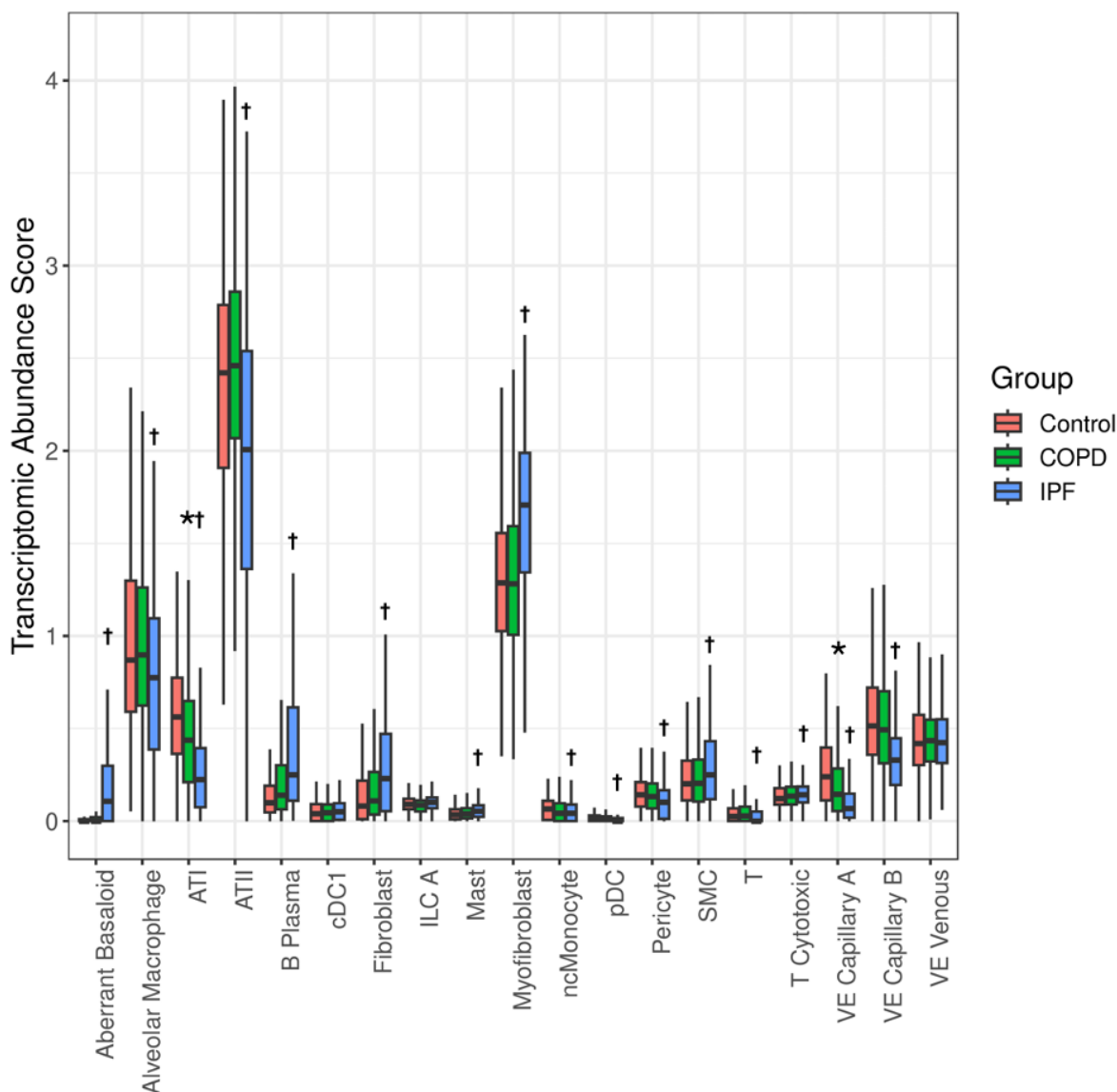
- 535 [13] Newman AM, Steen CB, Liu CL, Gentles AJ, Chaudhuri AA, Scherer F, et al.  
536 Determining cell type abundance and expression from bulk tissues with digital cytometry. *Nature*  
537 *Biotechnology* 2019;37:773–82. <https://doi.org/10.1038/s41587-019-0114-2>.
- 538 [14] Newman AM, Liu CL, Green MR, Gentles AJ, Feng W, Xu Y, et al. Robust enumeration  
539 of cell subsets from tissue expression profiles. *Nature Methods* 2015;12:453–7.  
540 <https://doi.org/10.1038/nmeth.3337>.
- 541 [15] Chen B, Khodadoust MS, Liu CL, Newman AM, Alizadeh AA. Profiling tumor  
542 infiltrating immune cells with CIBERSORT, Springer New York; 2018, p. 243–59.  
543 [https://doi.org/10.1007/978-1-4939-7493-1\\_12](https://doi.org/10.1007/978-1-4939-7493-1_12).
- 544 [16] Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, et al. limma powers differential  
545 expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Research*  
546 2015;43:e47–7. <https://doi.org/10.1093/nar/gkv007>.
- 547 [17] Szklarczyk D, Kirsch R, Koutrouli M, Nastou K, Mehryary F, Hachilif R, et al. The  
548 STRING database in 2023: proteinprotein association networks and functional enrichment  
549 analyses for any sequenced genome of interest. *Nucleic Acids Research* 2022;51:D638–46.  
550 <https://doi.org/10.1093/nar/gkac1000>.
- 551 [18] Ley B, Ryerson CJ, Vittinghoff E, Ryu JH, Tomassetti S, Lee JS, et al. A  
552 Multidimensional Index and Staging System for Idiopathic Pulmonary Fibrosis. *Annals of*  
553 *Internal Medicine* 2012;156:684. <https://doi.org/10.7326/0003-4819-156-10-201205150-00004>.
- 554 [19] Jin H, Liu Z. A benchmark for RNA-seq deconvolution analysis under dynamic testing  
555 environments. *Genome Biology* 2021;22. <https://doi.org/10.1186/s13059-021-02290-6>.

- 556 [20] Im Y, Kim Y. A comprehensive overview of RNA ceconvolution methods and their  
557 application. *Molecules and Cells* 2023;46:99–105. <https://doi.org/10.14348/molcells.2023.2178>.
- 558 [21] Bauer Y, White ES, Bernard S de, Cornelisse P, Leconte I, Morganti A, et al. MMP-7 is a  
559 predictive biomarker of disease progression in patients with idiopathic pulmonary fibrosis. *ERJ*  
560 *Open Research* 2017;3:00074–2016. <https://doi.org/10.1183/23120541.00074-2016>.
- 561 [22] Pardo A, Cabrera S, Maldonado M, Selman M. Role of matrix metalloproteinases in the  
562 pathogenesis of idiopathic pulmonary fibrosis. *Respiratory Research* 2016;17.  
563 <https://doi.org/10.1186/s12931-016-0343-6>.
- 564 [23] Zhang Y, Jiang M, Nouraie M, Roth MG, Tabib T, Winters S, et al. GDF15 is an  
565 epithelial-derived biomarker of idiopathic pulmonary fibrosis. *American Journal of Physiology-*  
566 *Lung Cellular and Molecular Physiology* 2019;317:L510–21.  
567 <https://doi.org/10.1152/ajplung.00062.2019>.
- 568 [24] Lagares D, Ghassemi-Kakroodi P, Tremblay C, Santos A, Probst CK, Franklin A, et al.  
569 ADAM10-mediated ephrin-B2 shedding promotes myofibroblast activation and organ fibrosis.  
570 *Nature Medicine* 2017;23:1405–15. <https://doi.org/10.1038/nm.4419>.
- 571 [25] Schäfer M, Werner S. The cornified envelope: a first line of defense against reactive  
572 oxygen species. *Journal of Investigative Dermatology* 2011;131:1409–11.  
573 <https://doi.org/10.1038/jid.2011.119>.
- 574 [26] Polverino F, Celli BR, Owen CA. COPD as an endothelial disorder: endothelial injury  
575 linking lesions in the lungs and other organs? (2017 Grover Conference Series). *Pulmonary*  
576 *Circulation* 2018;8:1–18. <https://doi.org/10.1177/2045894018758528>.

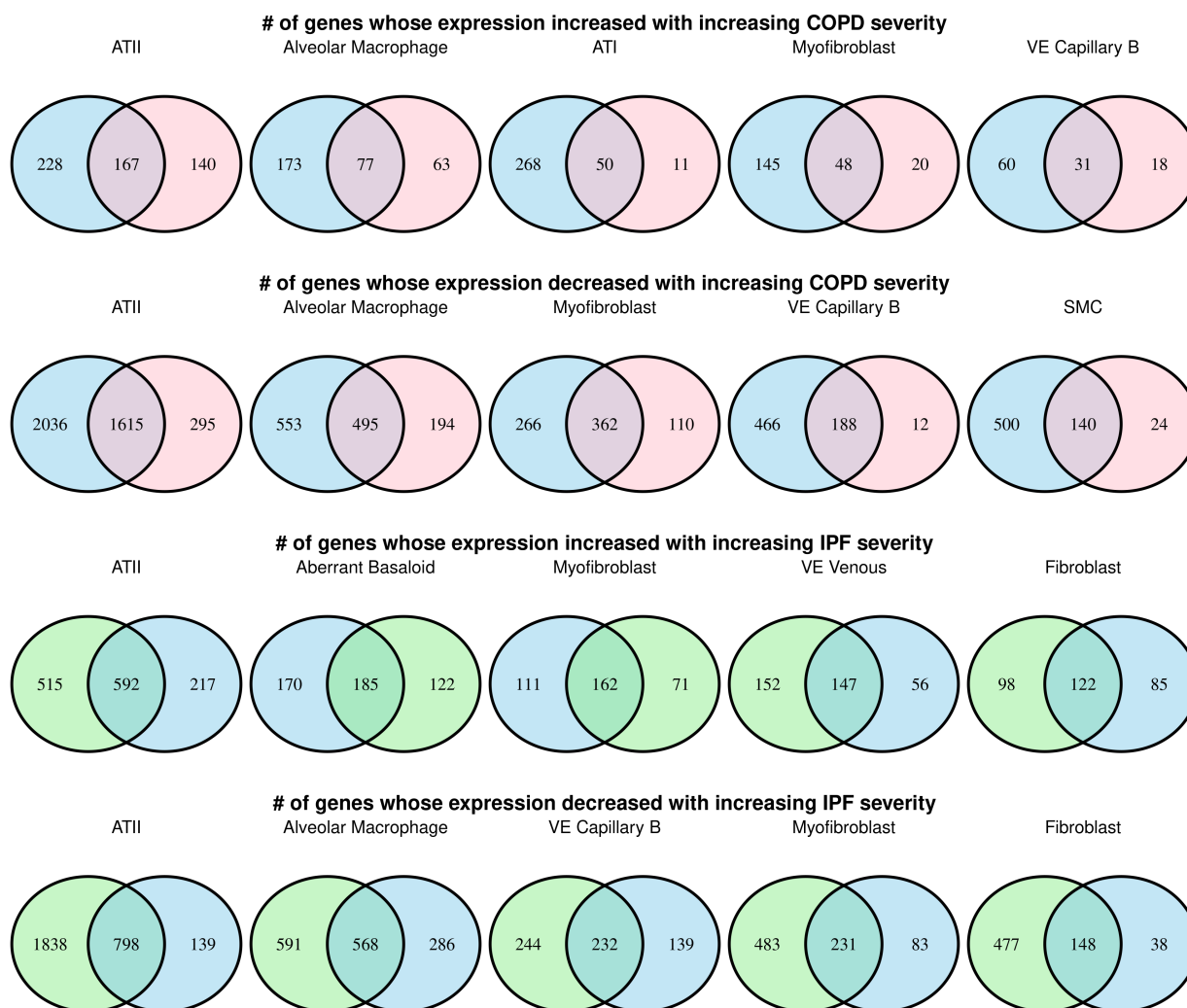
- 577 [27] Kapellos TS, Bassler K, Aschenbrenner AC, Fujii W, Schultze JL. Dysregulated  
578 functions of lung macrophage populations in COPD. *Journal of Immunology Research*  
579 2018;2018:1–19. <https://doi.org/10.1155/2018/2349045>.
- 580 [28] Tesfaigzi Y, Curtis JL, Petrache I, Polverino F, Kheradmand F, Adcock IM, et al. Does  
581 chronic obstructive pulmonary disease originate from different cell types? *American Journal of*  
582 *Respiratory Cell and Molecular Biology* 2023;69:500–7. [https://doi.org/10.1165/rcmb.2023-](https://doi.org/10.1165/rcmb.2023-0175ps)  
583 [0175ps](https://doi.org/10.1165/rcmb.2023-0175ps).
- 584 [29] Freeman CM, Curtis JL. Lung dendritic cells: shaping immune responses throughout  
585 chronic obstructive pulmonary disease progression. *American Journal of Respiratory Cell and*  
586 *Molecular Biology* 2017;56:152–9. <https://doi.org/10.1165/rcmb.2016-0272tr>.
- 587 [30] Dewhurst JA, Lea S, Hardaker E, Dungwa JV, Ravi AK, Singh D. Characterisation of  
588 lung macrophage subpopulations in COPD patients and controls. *Scientific Reports* 2017;7.  
589 <https://doi.org/10.1038/s41598-017-07101-2>.
- 590 [31] Morrow JD, Chase RP, Parker MM, Glass K, Seo M, Divo M, et al. RNA-sequencing  
591 across three matched tissues reveals shared and tissue-specific gene expression and pathway  
592 signatures of COPD. *Respiratory Research* 2019;20. <https://doi.org/10.1186/s12931-019-1032-z>.
- 593 [32] Huang Q, Wang Y, Zhang L, Qian W, Shen S, Wang J, et al. Single-cell transcriptomics  
594 highlights immunological dysregulations of monocytes in the pathobiology of COPD.  
595 *Respiratory Research* 2022;23. <https://doi.org/10.1186/s12931-022-02293-2>.
- 596 [33] Ogger PP, Byrne AJ. Macrophage metabolic reprogramming during chronic lung disease.  
597 *Mucosal Immunology* 2021;14:282–95. <https://doi.org/10.1038/s41385-020-00356-5>.



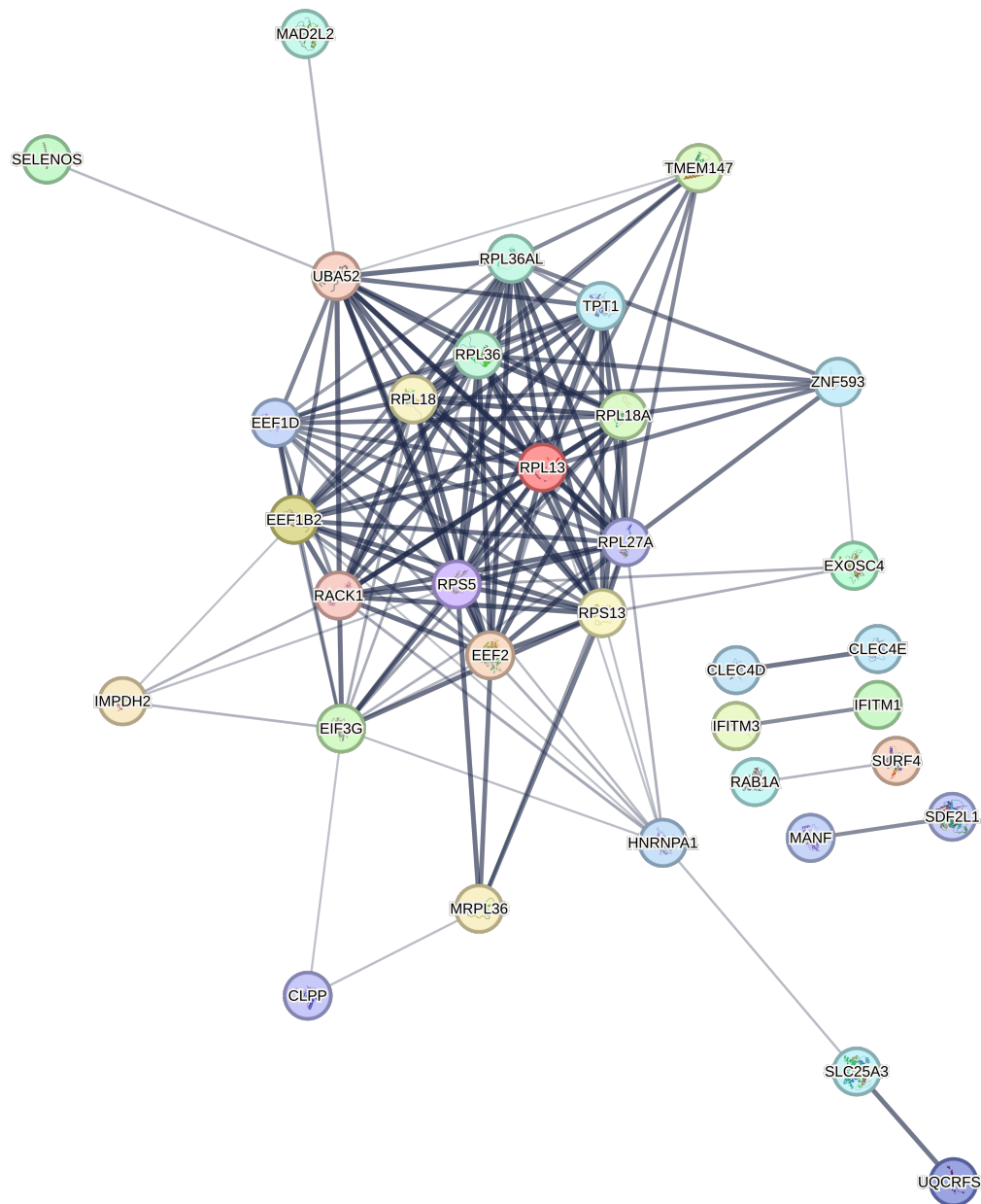
598 [34] Fujii W, Kapellos TS, Baßler K, Händler K, Holsten L, Knoll R, et al. Alveolar  
599 macrophage transcriptomic profiling in COPD shows major lipid metabolism changes. ERJ Open  
600 Research 2021;7:00915–2020. <https://doi.org/10.1183/23120541.00915-2020>.



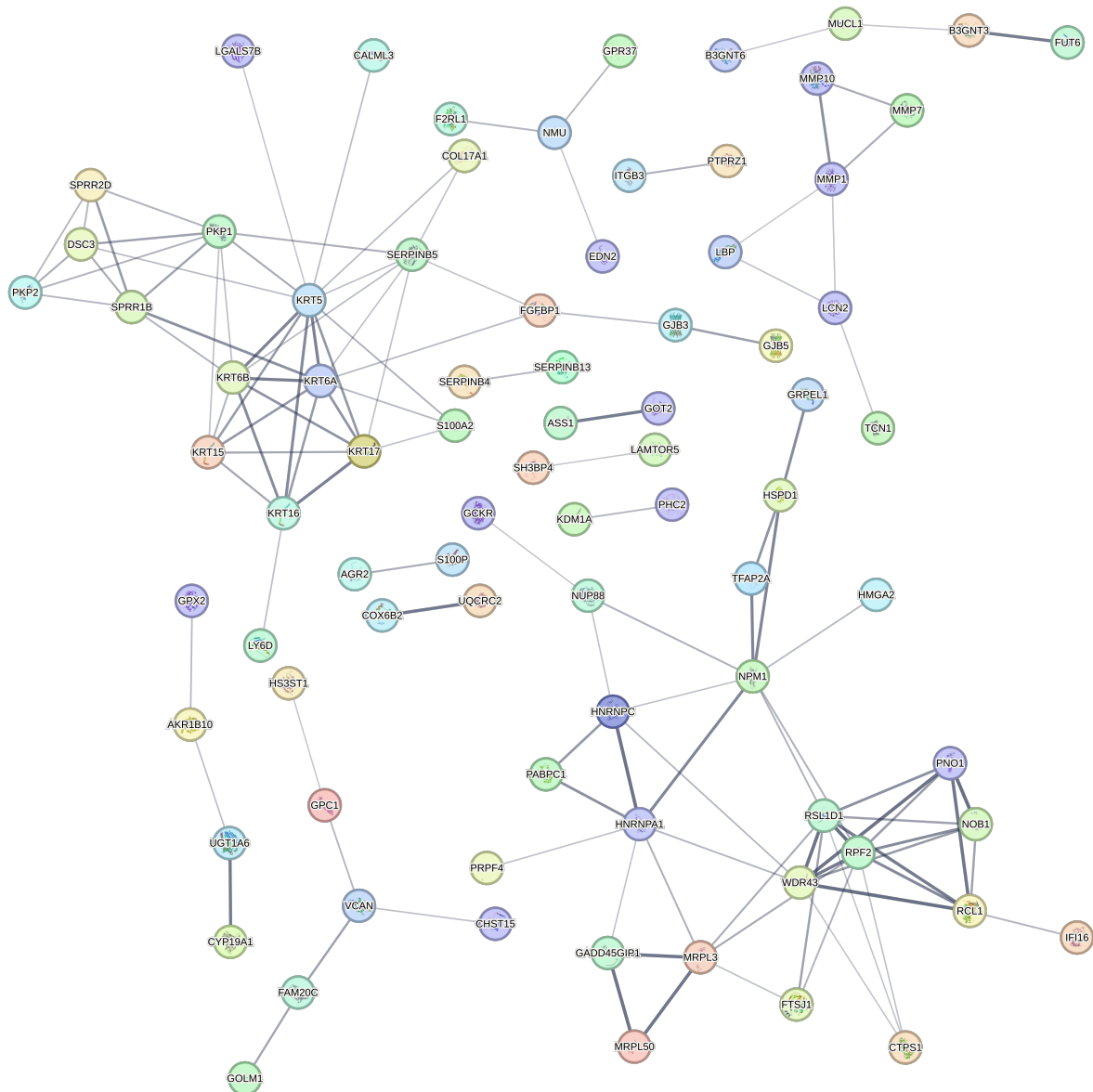
**Figure 1.** Boxplots showing cell type abundance score for each cell type split by disease status. Results are shown only for cell types detected in at least 10% of samples and whose median proportion was greater than 1%. Statistical comparison was tested using linear regression adjusting for age, sex, ever smoking and total pack-year. \* and † denote significant difference between COPD vs control and IPF vs control, respectively. Abbreviations: COPD, chronic obstructive pulmonary disease; IPF, idiopathic pulmonary fibrosis; ATI, alveolar epithelial type 1 cells; ATII, alveolar epithelial type 2 cells; cDC, classical dendritic cells; ILC, innate lymphoid cells; pDC, plasmacytoid dendritic cells; ncMonocyte, non-classical monocytes; SMC, smooth muscle cells; VE Capillary A, vascular endothelial - aerocyte capillary; VE Capillary B, vascular endothelial - general capillary; VE Venous, venous vascular endothelial.



**Figure 2.** Venn diagrams showing the five cell types with the most cell type-specific gene expression levels associated with disease severity in COPD and IPF lungs. Genes associated with DLCO, FEV1, and FVC are colored blue, pink, and green, respectively. Cell-types with a higher number of gene expressions associated with disease severity are ordered left to right.



**Figure 3.** Protein-protein interaction network for the proteins encoded by genes in alveolar macrophages positively associated with COPD severity. Edges represent protein-protein associations based on association confidence score calculated using STRING database (version 12.0). The edge line thickness indicates the strength of data support. Disconnected nodes in the network were hidden for illustrative purpose.



**Figure 4.** Protein-protein interaction network for proteins encoded by genes in aberrant basaloid cells that were positively associated with IPF severity. Edges represent protein-protein associations based on association confidence score calculated using STRING database (version 12.0). The edge line thickness indicates the strength of data support. Disconnected nodes in the network were hidden for illustrative purpose.