# Journal Pre-proof

Serum proteomics identify biomarkers associated with the pathogenesis of idiopathic pulmonary fibrosis

Lan Wang, Minghui Zhu, Yan Li, Peishuo Yan, Zhongzheng Li, Xiuping Chen, Juntang Yang, Xin Pan, Huabin Zhao, Shenghui Wang, Hongmei Yuan, Mengxia Zhao, Xiaogang Sun, Ruyan Wan, Fei Li, Xiaobo Wang, Hongtao Yu, Ivan Rosas, Chen Ding, Guoying Yu
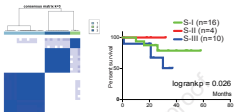
Please cite this article as: Wang L, Zhu M, Li Y, Yan P, Li Z, Chen X, Yang J, Pan X, Zhao H, Wang S, Yuan H, Zhao M, Sun X, Wan R, Li F, Wang X, Yu H, Rosas I, Ding C, Yu G, Serum proteomics identify biomarkers associated with the pathogenesis of idiopathic pulmonary fibrosis, *Molecular & Cellular Proteomics* (2023), doi: https://doi.org/10.1016/j.mcpro.2023.100524.
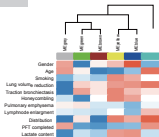
This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.
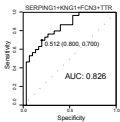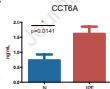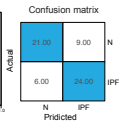
Molecular subtypes

Aging-associated module

Indicator correlated with lactic acid content

Biomarker combinations

1   **Serum proteomics identify biomarkers associated with the pathogenesis of idiopathic**

2   **pulmonary fibrosis**

3

4   Lan Wang[1#], Minghui Zhu[2#], Yan Li[3#], Peishuo Yan[1#], Zhongzheng Li[1], Xiuping Chen[3], Juntang

5   Yang[1], Xin Pan[1], Huabin Zhao[1], Shenghui Wang[1], Hongmei Yuan[1], Mengxia Zhao[1], Xiaogang

6   Sun[1], Ruyan Wan[1], Fei Li[2], Xiaobo Wang[2], Hongtao Yu[2], Ivan Rosas[4], Chen Ding[3*], Guoying

7   Yu[1*]

8   1.  State Key Laboratory of Cell Differentiation and Regulation, Henan International Joint

9       Laboratory of Pulmonary Fibrosis, Henan center for outstanding overseas scientists of

10      pulmonary fibrosis, College of Life Sciences, Institute of Biomedical Science, Henan

11      Normal University, Xinxiang, Henan 453007, China

12  2.  Henan Provincial Chest Hospital, Zhengzhou, Henan 450003, China

13  3.  State Key Laboratory of Genetic Engineering, Human Phenome Institute, Institutes of

14      Biomedical Sciences, and School of Life Sciences, Zhongshan Hospital, Fudan University,

15      Shanghai 200433, China.

16  4.  Division of Pulmonary, Critical Care and Sleep Medicine, Baylor College of Medicine,

17      Houston, TX 77030, USA

18  *Corresponding authors

19  #Contribute equally

20

21  **Running title**：Serum proteomics identify biomarkers of IPF

22

23

24

25  **Address for correspondence:**

26      Guoying Yu, PhD

27      Henan Normal University, 46 Jianshe Road, Xinxiang, Henan 453007, China

28      Email: guoyingyu@htu.edu.cn

29

30 **Abstract**

31 The heterogeneity of idiopathic pulmonary fibrosis (IPF) limits its diagnosis and treatment. The

32 association between the pathophysiological features and the serum protein signatures of IPF

33 currently remains unclear. The present study analyzed the specific proteins and patterns

34 associated with the clinical parameters of IPF based on a serum proteomic dataset by Data-

35 Independent Acquisition (DIA) using mass spectrometry. Differentiated proteins in sera

36 distinguished in IPF patients into three subgroups in signal pathways and overall survival.

37 Aging-associated signatures by WGCNA coincidently provided clear and direct evidence that

38 aging is a critical risk factor for IPF rather than a single biomarker. LDHA and CCT6A

39 expression, which were associated with glucose metabolic reprogramming, were correlated

40 with high serum lactic acid content in the patients with IPF. Cross-model analysis and machine

41 learning showed that a combinatorial biomarker accurately distinguished IPF patients from

42 healthy subjects with an AUC of 0.848 (95% CI = 0.684–0.941) and validated from another

43 cohort and ELISA assay. This serum proteomic profile provides rigorous evidence that enables

44 understanding of the heterogeneity of IPF and protein alterations that could help in its diagnosis

45 and treatment decisions.

46 Keywords: Serum proteome, Molecular subtype, machine learning, indicator panel,

47 combinatorial biomarker

48

49

50

51

52

1

**Introduction**

IPF is a chronic and fatal progressive fibrotic lung disease with a reported median survival of 3–5 years (1) (2). The heterogeneity of IPF and the various pathophysiological mediators involved in its clinical progression limit its diagnosis and treatment. Aging is one of the critical risk factors for IPF, with increasing evidence highlighting the important role of senescence in IPF(3). Cellular senescence leads to DNA damage, cell cycle arrest, telomere shortening(4), mitochondrial dysfunction, metabolic reprogramming, resistance to apoptosis, and deficient autophagy. Mitochondrial dysfunction, including the leakage of high energy electrons from the electron transport chain (ETC), disrupted cristae, and a diminished capacity for oxidative phosphorylation, establish a close link between senescence and IPF(5). Metabolic dysfunction alters processes during lung tissue repair, as well as crucial metabolic pathways such as augmented glycolysis and increased fatty acid oxidation, which are important drivers of fibroblast activation(6). In particular, altered lactate metabolism may be an underlying feature of IPF and a novel clinical diagnostic marker(7, 8).

The use of machine learning tool did not reach a formal recommendation in American Thoracic Society（ATS）/European Respiratory Society（ERS）/Japanese Respiratory Society（JRS）/ Latin American Thoracic Society（ALAT) clinical practice guideline, but more of a consideration in specific circumstances at certain centers to identify diagnostic markers and to combine these molecular markers with current diagnostic modalities in the multidisciplinary diagnosis of IPF. Novel biomarkers integrated into clinical diagnosis can include circulating markers or molecular signatures obtained from less invasive sampling(9). To date, most biomarkers are the molecules abundant enriched and associated with pathophysiological

75    process in a specific disease. Proteomic strategies have allowed extensive assessment of larger

76    patient cohorts and the identification of novel biomarkers, while reducing the need for invasive

77    acquisition and analysis of blood and body fluids(10). Improvements in deep proteomes may

78    result in the identification of individual biomarkers or biomarker panels that may not be directly

79    involved in the disease pathophysiology and may only be associated with it. These biomarkers

80    may have the potential to better understand the pathophysiology of IPF, not only for diagnostic

81    but also for therapeutic purposes.

82    Previous studies found that aberrations in complement activation and oxidative damage,

83    haptoglobin-related protein were identified as candidate marker in IPF using the label-free

84    plasma proteomics(11). Here, we wished to gain further insights into the changed serum

85    proteomic of IPF patients, to obtain the proteins associated with the disease pathophysiology.

86    A global correlation network related to clinical traits was constructed, and machine learning

87    was used to identify a combinatorial biomarker.

88

89    **Materials and methods**

90    Experimental Design and Statistical Rationale

91    The purpose of this study to identify signatures associated with the pathogenesis of idiopathic

92    pulmonary fibrosis in serum from IPF patients. The workflow is depicted in Figure S1. Serum

93    samples were collected from 30 IPF patients as a cohort, IPF was diagnosed based on

94    ATS/ERS/JRS/ALAT Clinical Practice Guidelines (12). Subjects were obtained at diagnosis

95    and followed by physicians according to institutional practices, including by high-resolution

96    computed tomography (HRCT) and pulmonary function tests (PFTs). All patients with IPF

3

97     underwent routine blood tests, including measurements of lactic acid concentrations and some

98     antibodies. None of the included patients had evidence of autoimmune syndromes, malignancy,

99     infections, or drug or occupational exposures associated with lung fibrosis. Serum samples from

100    30 healthy volunteers were collected as a control group, of which all participants underwent a

101    full medical examination prior to inclusion in the study. The validation cohort consisted of an

102    additional patient with IPF for the ELISA. The study was approved by the Henan Provincial

103    Chest Hospital Medical Research Ethics Committee (No. 2020-03-06). Oral and written

104    informed consent was obtained from all participants of this study. All samples used in this study

105    were collected at Henna Provincial Chest Hospital according to the guidelines in the

106    Declaration of Helsinki. The demographic and clinical characteristics of the cohorts are

107    provided, which including the summary data with statistics on age, sex, smoking status in Table

108    1 and other characteristics in the Table S1. A public IPF cohort (PRIDE project PXD010965)

109    that included 19 healthy individuals and 17 IPF patients was used to validate the accuracy of

110    the machine-learning-based classification of IPF. The animal handling procedures followed the

111    Henan Normal University Institutional Animal Care and Use Committee (IACUC, SMKX-

112    2019S002) guidelines, which coordinate with the Association of Animal Behavior and National

113    Regulations.

114    **Serum sample preparation**

115    Blood samples from IPF patients and healthy volunteers were taken from a vein in the cubital

116    fossa. The blood collection was done into commercial Monovette tubes containing tripotassium

117    ethylenediaminetetraacetic acid as the anticoagulant and whole blood glass tubes with

118    anticoagulant. The samples were centrifuged for serum separation (2000 rpm for 10 min, +4 ∘C)

4

119 immediately after collection. The supernatant was frozen at −80 ∘C before liquid

120 chromatography-mass spectrometry (LC-MS) analysis.

121 The 14 most abundant serum proteins were removed from each sample using commercial

122 depletion kits (High-Select™ Top14 Abundant Protein Depletion Mini Spin Columns),

123 according to the manufacturer's instructions. Following depletion, the proteins were denatured,

124 reduced, alkylated, digested into peptides, and desalted using a C-18 column for LC-MS/MS

125 analysis.

126

127 **High-pressure liquid chromatography and mass spectrometry**

128 Samples were subjected to LC-MS/MS, consisting of an EASY-nLC 1200 system coupled to a

129 nano-electrospray ion source and a Fusion Lumos Orbitrap (Thermo Fisher Scientific). Purified

130 peptides were separated on 150 μm I.D. × 15 cm columns (C18, 1.9 μm, 120Å, Dr. Maisch

131 GmbH). Each column was loaded with about 0.5 μg peptides in buffer A (0.1% formic acid),

132 followed by elution at a flow rate of 450 nL/min with a linear gradient of 3–30% of buffer B

133 (0.1% formic acid, 80% (v/v) acetonitrile) for 35 min, 75% buffer B for 7 min, 98% buffer B

134 for 1 min, and a wash with 98% buffer B for 2 min. The column temperature was maintained at

135 60°C using a Peltier element containing an oven developed in house.

136 MS spectra were acquired with a Data-Independent Acquisition (DIA) method. The DIA-MS

137 method consisted of an MS1 scan from 300 to 1,400 m/z range (AGC target of $4 \times 10^5$, maximum

138 injection time of 50 ms) at a resolution of 60,000 and 30 DIA segments (AGC target of $5 \times 10^4$,

139 maximum injection time of 22 ms) at a resolution of 15,000.

140

**Library-based DIA data analysis and quality control**

141

142 To build the spectral library, we acquired 128 DDA files on a Fusion Lumos Orbitrap mass

143 spectrometer in DDA mode, which was used as reference spectra libraries. A library was built

144 by Skyline-daily (22.2.1.278, University of Washington, USA) for DIA analysis, which were

145 composed of various body fluids and organ tissue samples from 64 individuals, covering blood,

146 hydrothorax, joint effusion, bile, ascites, cerebrospinal, urine, etc., with a deep fractionation

147 ranging from 7 to 31. For Skyline library building, carbamidomethyl (C) was set as the fixed

148 modification, and acetyl (protein N-term) and oxidation (M) were set as the variable

149 modifications. Two missed trypsin cleavages were allowed. Precursor ion score charges were

150 limited to +2, +3, and +4. The precursor and fragment tolerance were set as dynamic. Finally,

151 a library containing 68,781 peptides and 4,437 proteins was built. In our previous research, the

152 DIA library has been used for blood molecular markers for the pathophysiology and clinical

153 progress of COVID-19 (13). For Skyline analysis, the default setting was used for library-based

154 DIA analysis according to the standard workflow in Skyline

155 (https://skyline.ms/_webdav/home/software/Skyline/@files/tutorials). A total of 60 raw files'

156 reports were exported by Skyline DIA analysis, and were merged into an integrated expression

157 matrix including the expression of each single protein, of which all identified distinct peptides

158 were used for the corresponding protein quantification. The detection Q value was set to 5% at

159 the peptide and protein levels. Proteome qualification was performed as previously reported

160 with the iBAQ algorithm(14), followed by normalization to the fraction of the total (FOT),

161 defined as a protein's iBAQ divided by the total iBAQ of all identified proteins within one

162 sample, thus representing the normalized abundance of a particular protein across samples.

163 Finally, the FOT values were further multiplied by $10^5$ for ease of presentation, and missing

6

164    values were replaced by the minimal value.

165    The quality of proteomic data was ensured at multiple levels. Instrument performance was

166    evaluated using a whole cell extract of HEK293T cells. To avoid carryover, blank samples

167    (buffer A) were run after every five injections. The consistency of sample collection and

168    handling was validated by assessing the abundance of the quality markers FGA, FGB, and FGC.

169

170    **Differential protein analysis**

171    The differential expression of proteins in IPF patients and healthy controls was also analyzed

172    by Student's t-tests. Proteins differentially expressed with p-values < 0.05 and fold changes >

173    1.5 or < 2/3 were visualized using an R package heatmap. Between-group analysis of DEPs was

174    performed using paired two-class analysis of the same R package with an FDR threshold of

175    0.05.

176

177    **Pathway enrichment analysis and functional annotation**

178    The biological characteristics of the three IPF subtypes and the proteins differentially expressed

179    by IPF patients and healthy controls were determined by pathway enrichment analysis with

180    Reactome. The statistical significance of pathway enrichment was determined by Fisher's exact

181    test and pathways with an FDR threshold of 0.05 were regarded as being significantly regulated.

182

183    **Proteome molecular subtyping of IPF**

184    Prior to clustering analysis, proteins that were expressed in more than 25% of patient samples

185    were selected (n = 1190) (Table S4). The serum proteomic subtypes of IPF were identified by

7

186  consensus clustering (R package Consensus Cluster Plus v.1.48.0) (15). A total of 1190 proteins

187  were subjected to k-means clustering with up to six clusters. The consensus matrix of k = 3

188  showed clear among-cluster separation (Figure S3A), and the cumulative distribution function

189  of the consensus matrix for each k-value was measured. Clustering by k = 3 resulted in the

190  lowest proportion of ambiguous clustering. To determine the correlations between proteomic

191  subtypes and clinical features, categorical variables, including age, gender, smoke status, and

192  HRCT characteristics, were assessed by Fisher's exact tests.

193

194  **WGCNA analysis**

195  To identify differentially co-expressed gene modules, WGCNA was applied to the proteins that

196  were expressed in more than 67% of patient samples (n =687). WGCNA was performed in R

197  (R Core Team, 2019) using a WGCNA package (16). Module eigenproteins were calculated as

198  the first principal components of the co-expressed genes in the module (17, 18). The eigengenes

199  of each module were used to measure the association between a module and clinical information.

200  The eigengene-based connectivity (kME) was used to represent the strength of a gene's

201  correlation with other gene module members.

202

203  **Machine-learning-based selection of biomarker combinations of IPF**

204  Biomarker combinations were identified using the random forest method, a machine learning

205  method that can predict the value of a response variable. Data with coefficients of variation

206  (CV) less than 0.5 were selected as candidate reservoirs, with no more than four proteins

207  randomly selected to form the potential optimal biomarker combination (OBC), and 5,000

208  potential OBCs were prepared. Each candidate OBC was subjected to 5-fold cross-validation,

209  with the original dataset randomly divided 4:1 into a training set and a verification set. The

210  training set was used to train the model, and the verification set was used to evaluate the model.

211  In penalized logistic regression (PLR) , the weights of four proteins were optimized iteratively

212  using the least shrinkage and selection operator (Lasso, L1 regularization) penalty and the ridge

213  regression (L2 regularization) penalty. The combination with the highest AUC value was

214  selected. To simplify OBC, sets of any three of the four proteins were selected, resulting in four

215  combinations, and the AUC values of these combinations were compared with the AUC value

216  of OBC. The combination with an AUC value closest to that of OBC was selected as the final

217  combination. The PLR algorithm was implemented in R 4.1.2 with the glmnet package.

218

219  **Survival analysis**

220  Univariate Cox regression analysis was conducted to determine the relationship between the

221  expression of proteins and prognosis of IPF patients. Proteins with a p-value < 0.05 were

222  regarded as prognostic proteins. After that, patients were divided into high-risk and low-risk

223  groups by setting the median value of risk scores as cut-off value. The overall survival (OS) of

224  these two groups was calculated by the Kaplan-Meier method with log-rank test. All statistical

225  analyses were performed using Prism 8 software and the R package "survival", with statistical

226  significance defined as $p < 0.05$.

227

228  **Cell culture**

229  The human lung fibroblast cell line (MRC-5) was purchased from the ATCC (CCL-171). Cells

9

230  were cultured in DMEM supplemented with 10% fetal bovine serum and a 1% antibiotic-

231  antimycotic solution at 37°C in 5% CO2.

232

233  **Plasmids RNA interference and transfection**

234  The human CCT6A gene was cloned into the pCDNA3.1 plasmid (Generay Biotech, CN).

235  Fibroblasts grown to 80–90% confluence were transfected with this plasma using

236  Lipofectamine 3000 reagent according to the manufacturer's protocol. The CCT6A siRNA

237  transfection target sequence, 5′-GTGTCATTAGAGTATGAGA-3′, and a negative control were

238  purchased from RiboBio. The siRNAs (75 nM) were transfected into cells using INVI DNA

239  RNA Transfection Reagent (Invigentech) according to the manufacturer's instructions.

240

241  **Protein extraction and western blot analysis**

242  Mouse lung tissue samples and cultured cells were lysed in RIPA lysis buffer. Equal amounts

243  of protein were separated on SDS-PAGE and transferred to PVDF-membranes, which were

244  hybridized overnight with appropriate primary antibodies. The membranes were washed and

245  incubated with horseradish peroxidase-conjugated secondary antibodies, followed by

246  visualization using the Odyssey Fc Dual-Mode Imaging System (LI-COR, USA), according to

247  the manufacturer's instructions.

248

249  **Immunofluorescence staining**

250  Transfected fibroblasts were fixed with 4% paraformaldehyde and permeabilized with 0.3%

251  Triton X100/PBS. Cells were incubated with primary antibodies at 4°C overnight followed by

10

252 incubation with fluorescent-labeled secondary antibodies for 30 min at 37℃. Images were

253 visualized using an Axio Imager D2 (Zeiss, GER).

254

255 **Extracellular flux technology**

256 The extracellular acidification rate (ECAR) of fibroblasts was measured using a Seahorse XF96

257 Extracellular Flux Analyzer (Seahorse Bioscience, USA). All assays were performed using a

258 seeding density of 30,000 cells/well in 200 μL DMEM in an XF96 cell culture microplate

259 (Seahorse Bioscience). ECAR was measured after sequential addition of glucose, oligomycin,

260 and 2-DG, to reach working concentrations of 10 mM, 1 μM, and 50 mM, respectively.

261

262 **LDH activity**

263 LDH activity was assessed using LDH activity assay kits, according to the manufacturer's

264 instructions. Briefly, extract was added to the transfected cells, and the cells were disrupted by

265 ultrasound and centrifuged at 8,000 g for 10 min at 4℃. LDH activity was evaluated by

266 measuring the amount of pyruvate produced.

267

268 **Lactate assay**

269 The intracellular and tissue concentrations of lactate were determined using Lactate Assay Kits,

270 according to the manufacturer's instructions. Tissues or cells were homogenized in four volumes

271 of Lactate Assay Buffer and centrifuged at 13,000g for 10 minutes to remove insoluble material.

272 The samples were deproteinized with a 10 kDa MWCO spin filter to remove lactate

273 dehydrogenase, and the absorbance of the soluble fraction at 570 nm was measured.

274

**Immunoassays**

Serum protein concentrations were measured using commercially available ELISA kits, as described by the manufacturer. Measure the absorbance of each sample at 450nm with Microplate Reader（Thermo Fisher）. For immunohistochemistry (IHC) staining, paraffin-embedded tissue sections(5 μm thick) were de-paraffinized and dehydrated, followed by antigen retrieval according to standard procedures. Tissue samples were incubated with specific antibodies, with images captured by AxioScan.Z1 (Zeiss).

**Statistical analysis**

GraphPad Prism 8.0 and R was used for statistical analysis. The details of experiments can be found in the methods and figure legends. Genes with p-values < 0.05 and fold changes > 1.5 or other thresholds were visualized using R package heatmaps. Between-group analysis of DEPs was performed using paired two-class of the same R package with an FDR threshold of 0.05. Pathway enrichment to identify pathway alterations was analyzed using Reactome. Differential analysis of samples with different phenotypes was performed using Fisher's exact t-tests, with DEPs compared in groups of patients with IPF and healthy controls. Spearson rank analysis was used to analyze the correlation. GraphPad Prism 8.0 was used to analyze the quantitative results of the cell / animal experiments and ELISA results. Significant differences between groups were evaluated using the student's t-test or analysis of variance (ANOVA). p< 0.05 was considered to be statistically significant.

295

12

## Results

**Serum proteome profiling of IPF**

The serum proteomic landscape was investigated in 30 patients with IPF and 30 healthy subjects

differing in demographic and clinical characteristics, including by gender, age, smoking status,

features of HRCT, and others (Table 1, Table S1). A data-independent acquisition (DIA)

strategy was adopted (Figure S1), and the consistency of the MS performance of the whole

HEK293T cell extract was assessed using Spearman correlation coefficients (average

correlation coefficient; R = 0.89) (Figure S2A). The abundance profiles of the quality markers

FGA, FGB, and FGG indicated that the collection and handling of the samples were regular[19]

(Figure S2B). About 2,383 gene products were collected from the 30 healthy subjects and the

30 patients with IPF (Figure 1A), with the number of proteins per sample ranging from 703 to

1,014 (median 892) (Figure 1B, Table S2). The abundance of the identified proteins varied

widely, with APOA1 being most abundant and ATP6V1A being the least abundant (Figure 1C).

Sixty-seven significantly differentially expressed proteins (DEPs) (P < 0.05 and a differential

expression ratio [IPF/N] >1.5 or <0.67) were identified (Figure S2C, Table S3). Of the DEPs

3.7% upregulated, whereas 3.8% of the significantly downregulated proteins in patients with

IPF ( Figure S2D)


**Three molecular subtypes of IPF and their association with clinical features**

Consensus Cluster Plus (Table S4, Figure S3A) analysis of the top 1,190 DEPs identified three

distinct patient clusters (S-I,S-II,S-III) with differences in survival (Figure 1D). The 30 patients

with IPF were followed-up for a median 27.9 months (range, 1–58 months). Association

13

318    analysis between IPF subtypes and OS demonstrated that OS was longest in the S-II and shortest

319    in the S-III (log-rank $P$ = 0.026, Figure 1E). IPF patients in the three proteomic subgroups

320    showed distinct molecular features, including differences in subgroup-specific pathways and

321    expression of representative proteins (Figure 1F, G, Figure S3B, Table S5). Higher expression

322    of BMP2K, which has been implicated in endocytosis and cell differentiation[20], was

323    associated with a longer OS in the S-I; and a high level of PI16, a shear stress and inflammation-

324    regulated inhibitor of MMP2[21], increased OS in the S-II. By contrast, elevated expression of

325    ATP5A1, a subunit of mitochondrial ATP synthase, was associated with a poorer OS in the S-

326    III (Figure 1H). These specific protein signatures may enable classification of these IPF

327    subgroups. The associations between proteomic subtypes and clinical features were examined

328    using Fisher's exact tests for categorical data and Wilcoxon rank-sum tests for continuous data.

329    We found that younger age was closely associated with longer OS in the S-II (Figure 1I),

330    indicating that age affects the survival of patients with IPF[22].

331

332    **Aging-associated signatures highlighted in the sera of IPF patients**

333    Weighted gene correlation network analysis (WGCNA) of a single dataset composed of samples

334    from all 30 IPF patients with 686 proteomic variables and ten clinical traits yielded the global

335    correlation network heatmap shown in Figure 2A (Table S6). Module-trait relationships analysis

336    showed that the module MEturquoise was positively associated with age patterns (Figure 2B).

337    The signatures correlated with age were clustered, and the top altered proteins in this module

338    mainly belonged to the S-III subgroup with elder patients (Figure 2C, D). Cellular senescence-

339    associated proteins, such as KL (Klotho), HSP90AB1, and SERPINE1; mitochondrial

14

340  dysfunction-associated proteins, including HSPD1, ATP5A1, and SDPR; and several other

341  proteins associated with DNA repair and the cell cycle, such as HIST2H2BE, NCK1, S100A8,

342  and CDK10, were significantly upregulated in S-III subgroup. In line with our findings,

343  VCAM1 and POSTN expression correlated positively with age(23), whereas UBA, CD14,

344  ORM1, and ORM2, which are involved in inflammatory responses, and CREM and CAMKK1,

345  which are involved in cell apoptosis, correlated negatively with age in the S-III subgroup.

346  SERPINA4, an age-related marker in lung disease(24), was decreased in the S-III subgroup

347  (Figure 2D). Moreover, increased expression of HSP90AB1 and reduced expression of

348  CAMKK1 were associated with poor survival in patients with IPF (Figure 2E).These protein

349  correlation profiles reflect the complex relationships between age and cellular senescence,

350  mitochondrial dysfunction, DNA repair and replication, inflammatory response, and cell

351  apoptosis.

352

353  **Integration of specific molecular markers with high level of lactic acid for**

354  **multidisciplinary diagnosis of IPF**

355  Increased glycolysis contributes to IPF by regulating glucose metabolic enzymes; these

356  enzymes are secreted and can be measured in blood. HK1, PFKP, ENO1/3, GAPDH, LDHA,

357  and ALDOB were significantly differentially expressed in the IPF and control groups (Figure

358  3A, B). Lactate dehydrogenase (LDH) converts pyruvate to lactic acid during glycolysis, with

359  human LDH, consisting of two subunits, LDHA and LDHB, being a key glycolytic terminal

360  enzyme that catalyzes the interconversion of pyruvate and lactate in the anaerobic glycolytic

361  pathway. Compared with controls, LDHA and LDHB were altered in the sera of patients with

15

362    IPF (Figure 3C), with survival analysis showing that LDHA may be a significant predictor of

363    poor prognosis in these patients (Figure 3D). Proteomics data showed that the level of serum

364    LDHA was upregulated in IPF patients with high lactate content (>1.7 mmol/L, Table S1,Table

365    S2), based on routine blood tests by ELISA. In addition, the expression of CCT6A, which was

366    predicted to act through an interactive network of signaling pathways with LDHA, was

367    increased in the serum of patients with IPF (Figure 3E). To explore the association of CCT6A

368    with high serum lactic acid content, we measured the levels of CCT6A in IPF patients and the

369    bleomycin model of lung fibrosis in mice. ELISA analysis confirmed that the level of CCT6A

370    was higher in IPF patients in an independent cohort (Figure 3F, Table S7), and the increases

371    were in accordance with MS data (Figure S4). IHC staining of lung tissue from patients with

372    IPF showed that CCT6A was mainly expressed by macrophages and the alveolar epithelium

373    surrounding the fibrotic interstitium, but was weakly expressed in normal alveolar epithelium

374    (Figure 3G). CCT6A expression was also significantly increased in the bleomycin model of

375    lung fibrosis in mice (Figure 3H, I). Moreover, the downregulation of GAPDH observed in the

376    sera of patients with IPF (Figure 3 B) was also observed in fibrotic mouse lungs (Figure 3J).

377    The increased levels of serum CCT6A in patients with IPF were associated with elevated lactic

378    acid concentrations, which may lead to pulmonary fibrosis.

379    To demonstrate that the changes of CCT6A have a direct effect on fibroblast phenotype, CCT6A

380    was overexpressed or knocked down in MRC-5 cells. Overexpression of CCT6A significantly

381    enhanced the expression of α-SMA in MRC-5 cells (Figure 4A-C), whereas knockdown of

382    CCT6A reduced the levels of FN-1 and Col1A1 (Figure 4D, E), indicating that CCT6A

383    promotes the development of lung fibrosis. To further clarify the association of increased

16

384 CCT6A with the high content of lactic acid in the sera of patients with IPF, real-time

385 extracellular acidification rate (ECAR) was measured using the Seahorse XFe96 Analyser

386 (Agilent Technologies). Overexpression of CCT6A was associated with significant increases in

387 glycolysis rate and glycolytic capacity (Figure 4F), as was lactate production in the supernatants

388 of MRC-5 cells and in the lungs of bleomycin-treated mice (Figure 4G, H). Cells

389 overexpressing CCT6A also showed significant upregulation of the expression of LDHA in

390 mRNA and protein level, and decreased production of pyruvate (Fig 4I-L). Collectively, these

391 results show that CCT6A plays an important role in glycolysis through regulation of LDHA and

392 drives pulmonary fibrosis.

393

394 **Machine-learning-based selection of combinatorial biomarkers for classification of IPF**

395 A machine-learning algorithm involving potential combinatorial biomarkers was developed to

396 classify IPF patients and healthy subjects (Figures S5A). Candidate biomarkers were selected

397 from the significantly differentially expressed proteins using PLR for model training and

398 parameter optimization. This process generated a set of combinatorial biomarkers, including

399 serpin G1 (SERPING1), kininogen 1 (KNG1), ficolin 3 (FCN3), and transthyretin (TTR). The

400 5-fold cross-validation AUC value of this four-protein combinatorial that differentiated IPF

401 patients and healthy individuals was 0.826 (95% confidence interval [CI] = 0.700–0.800)

402 (Figure 5A, B). The corresponding matrix demonstrated that the training model could correctly

403 classify different samples with high accuracy (Figures 5C). The accuracy of the machine-

404 learning-based classification of IPF was validated in a public IPF cohort (PRIDE project

405 PXD010965) that included 19 healthy individuals and 17 IPF patients. The AUC value for the

406    diagnosis of IPF was 0.848 (95% CI = 0.684–0.941) (Figure 5 D, E), with the data matrix

407    showing promising accuracy in this independent cohort (Figures 5F). The combinatorial

408    biomarkers predicted poorer, but not significantly different, OS in our cohort (Figures 5G).

409    Lack of survival information prevented determination of the ability of the combinatorial

410    biomarkers to predict OS in the public dataset, but these markers exhibited significant

411    performance based on their relative abundances (Figures 5H, Figures S5B).

412    Our previous study showed that thyroid hormone inhibits lung fibrosis in mice(25). Because

413    TTR transports thyroid hormones in plasma and cerebrospinal fluid, the serum concentrations

414    of TTR were measured by ELISA in an independent cohort. Serum TTR concentrations were

415    significantly lower in IPF patients than in normal controls (Figure 5I). Although low

416    transthyretin levels were reported to correlate with age and stroke(26), serum TTR level did not

417    significantly correlate with age in our patient cohort (Figure 5J).

418

**Discussion**

420    Poor molecular understanding of the heterogeneity of IPF can impede determination of its

421    pathogenesis, leading to inefficient treatment and an inability to predict its occurrence. To

422    address this problem, we sought to determine the serum protein profile in patients with IPF.

423    Analysis showed that IPF could be classified into three subtypes, which exhibited its

424    heterogeneity and diversity. This study also found that CCT6A was associated with the elevated

425    levels of lactic acid in IPF. A global correlation network was developed to identify the indicators

426    of senescence associated with IPF; a combinatorial predictive biomarker that can be used to

427    distinguish patients with IPF from healthy subjects.

428    Molecular subtyping can stratify patients into subtypes associated with clinical features,

18

429  responses to treatment, and biological characteristics(27). IPF could be classified into three

430  subtypes based on serum proteomes, with these proteomic subtypes differing in signaling

431  pathways and clinical outcome. Specifically, patients with the S-III subtype had a poorer

432  prognosis. In addition, a functional module related to senescence was found to be associated

433  with the S-III subtype. Because aging is a multifactorial series of molecular alterations that

434  result in progressive reduction of lung tissue function, the involvement of proteins associated

435  with various physiological processes related to aging was not surprising, serum proteins may

436  be candidate markers of aging. The altered-senescence-associated protein patterns in S-III were

437  related to aging rather than to a single biomarker, providing clear and direct evidence that aging

438  is a critical risk factor for IPF.

439  Coupling of altered proteins under defined conditions could exploit the information content of

440  serum and identify biomarkers likely to be of clinical value. MS-based proteomics can enable

441  assessment of the roles of blood proteins in clinical diagnoses, as well as identifying new

442  biomarkers and biomarker panels. Analysis of serum proteomes can result in the detection of

443  secreted metabolic enzymes, including those involved in enhancing glycolysis, upregulation of

444  the key metabolic enzyme LDHA was indicative of poorer clinical outcomes. Therefore, the

445  presence of high levels of CCT6A and LDHA and high serum lactic acid concentrations may

446  be diagnostic of IPF.

447  Use of machine learning to explore the ability of combined biomarkers to predict disease

448  outcomes and prognosis is a promising strategy to improve the accuracy of diagnostic

449  performance. Intriguingly, SERPING1 itself is a candidate biomarker in patients with

450  tuberculosis(28), downregulation of KNG1 expression was observed in patients with sepsis-

19

451    induced ALI(29), and TTR is a specific biomarker for the clinical diagnosis of non-small cell

452    lung carcinoma(30). In the present study, these three proteins were selected by the machine-

453    learning algorithm as the most important indicators for classification of IPF, showing high

454    specificity and sensitivity in two independent patient cohorts. Furthermore, the combinatorial

455    biomarker panel and clinical data was found to be prognostic in this patient cohort.

456    The present study had several limitations. The number of patients included in the study cohort

457    was small, as were the numbers in each of the subgroups, suggesting the need for studies in

458    larger patient cohorts, as well as validation of these biomarkers by methods other than serum

459    proteome analysis. Moreover, the kits used to process serum samples can lead to the depletion

460    of highly abundant proteins. For example, EDTA could interfere with the precise determination

461    of MMPs, such as MMP7 and CCL18, previously shown to be markers of IPF(31). Taken

462    together, our data characterized the molecular subtypes of IPF and identified a biomarker panel

463    associated with the pathophysiology of IPF. These results strongly suggest that measuring

464    CCT6A and LDHA, along with high serum levels of lactic acid, could be diagnostic of IPF.

465    Additional studies in larger patient cohorts are needed to determine whether the combination of

466    these three biomarkers could accurately predict IPF.

467

468    **Data and materials availability**

469    The raw mass spectrometry (MS) proteomics data generated in this study have been deposited

470    in the ProteomeXchange Consortium via the iProX partner repository (http://www.iprox.

471    cn/)(32) under Project ID IPX0004334000, and can be accessed with a direct link

472    https://www.iprox.cn/page/PSV023.html;?url=1664089052598znXd with the password: ASQd.

473

**Author Contributions**

G.Y. and C.D.: Designed the research plan. L. W.: Data curation, Writing- Original draft preparation, L.W. Y. L., X. Ch.: Proteomics experiments, Z. L.,H.Z: Statistical analysis, Data analysis and data visualization, S.Y.: Performed cell and mouse assay and related data visualization, IHC staining of tissue samples, J.Y and X.P.: Performed ELISA assay, H.Y. and M.Z. :Consulted on clinical questions. I.R. Writing – review & editing. All authors discussed the results and commented on the manuscript.

**References**

1. Lederer, D. J., and Martinez, F. J. (2018) Idiopathic Pulmonary Fibrosis. *N Engl J Med* 378, 1811-1823

2. Hutchinson, J., Fogarty, A., Hubbard, R., and McKeever, T. (2015) Global incidence and mortality of idiopathic pulmonary fibrosis: a systematic review. *Eur Respir J* 46, 795-806

3. Selman, M., Lopez-Otin, C., and Pardo, A. (2016) Age-driven developmental drift in the pathogenesis of idiopathic pulmonary fibrosis. *European Respiratory Journal* 48, 538-552

495    4.    Chakravarti, D., LaBella, K. A., and DePinho, R. A. (2021) Telomeres: history, health, and

496    hallmarks of aging. *Cell* 184, 306-322

497    5.    Su, Y. J., Wang, P. W., and Weng, S. W. (2021) The Role of Mitochondria in Immune-Cell-

498    Mediated Tissue Regeneration and Ageing. *International Journal of Molecular Sciences* 22

499    6.    Henderson, N. C., Rieder, F., and Wynn, T. A. (2020) Fibrosis: from mechanisms to

500    medicines. *Nature* 587, 555-566

501    7.    Le, A., Cooper, C. R., Gouw, A. M., Dinavahi, R., Maitra, A., Deck, L. M., Royer, R. E.,

502    Vander Jagt, D. L., Semenza, G. L., and Dang, C. V. (2010) Inhibition of lactate dehydrogenase

503    A induces oxidative stress and inhibits tumor progression. *Proc Natl Acad Sci U S A* 107, 2037-

504    2042

505    8.    Kottmann, R. M., Kulkarni, A. A., Smolnycki, K. A., Lyda, E., Dahanayake, T., Salibi, R.,

506    Honnons, S., Jones, C., Isern, N. G., Hu, J. Z., Nathan, S. D., Grant, G., Phipps, R. P., and Sime,

507    P. J. (2012) Lactic acid is elevated in idiopathic pulmonary fibrosis and induces myofibroblast

508    differentiation via pH-dependent activation of transforming growth factor-β. *Am J Respir Crit*

509    *Care Med* 186, 740-751

510    9.    Raghu, G., Remy-Jardin, M., Myers, J. L., Richeldi, L., Ryerson, C. J., Lederer, D. J., Behr,

511    J., Cottin, V., Danoff, S. K., Morell, F., Flaherty, K. R., Wells, A., Martinez, F. J., Azuma, A.,

512    Bice, T. J., Bouros, D., Brown, K. K., Collard, H. R., Duggal, A., Galvin, L., Inoue, Y., Jenkins,

513    R. G., Johkoh, T., Kazerooni, E. A., Kitaichi, M., Knight, S. L., Mansour, G., Nicholson, A. G.,

514    Pipavath, S. N. J., Buendia-Roldan, I., Selman, M., Travis, W. D., Walsh, S., Wilson, K. C., Soc,

515    A. T., Soc, E. R., Soc, J. R., and Soc, L. A. T. (2018) Diagnosis of Idiopathic Pulmonary Fibrosis

516    An Official ATS/ERS/JRS/ALAT Clinical Practice Guideline. *Am J Resp Crit Care* 198, E44-

517  E68

518  10. Geyer, P. E., Holdt, L. M., Teupser, D., and Mann, M. (2017) Revisiting biomarker

519  discovery by plasma proteomics. *Mol Syst Biol* 13

520  11. Saraswat, M., Joenväärä, S., Tohmola, T., Sutinen, E., Vartiainen, V., Koli, K., Myllärniemi,

521  M., and Renkonen, R. (2020) Label-free plasma proteomics identifies haptoglobin-related

522  protein as candidate marker of idiopathic pulmonary fibrosis and dysregulation of complement

523  and oxidative pathways. *Sci Rep* 10, 7787

524  12. Raghu, G., Remy-Jardin, M., Myers, J. L., Richeldi, L., Ryerson, C. J., Lederer, D. J., Behr,

525  J., Cottin, V., Danoff, S. K., Morell, F., Flaherty, K. R., Wells, A., Martinez, F. J., Azuma, A.,

526  Bice, T. J., Bouros, D., Brown, K. K., Collard, H. R., Duggal, A., Galvin, L., Inoue, Y., Jenkins,

527  R. G., Johkoh, T., Kazerooni, E. A., Kitaichi, M., Knight, S. L., Mansour, G., Nicholson, A. G.,

528  Pipavath, S. N. J., Buendia-Roldan, I., Selman, M., Travis, W. D., Walsh, S., Wilson, K. C.,

529  American Thoracic Society, E. R. S. J. R. S., and Latin American Thoracic, S. (2018) Diagnosis

530  of Idiopathic Pulmonary Fibrosis. An Official ATS/ERS/JRS/ALAT Clinical Practice Guideline.

531  *Am J Respir Crit Care Med* 198, e44-e68

532  13. Chen, Y. M., Zheng, Y., Yu, Y., Wang, Y., Huang, Q., Qian, F., Sun, L., Song, Z. G., Chen,

533  Z., Feng, J., An, Y., Yang, J., Su, Z., Sun, S., Dai, F., Chen, Q., Lu, Q., Li, P., Ling, Y., Yang, Z.,

534  Tang, H., Shi, L., Jin, L., Holmes, E. C., Ding, C., Zhu, T. Y., and Zhang, Y. Z. (2020) Blood

535  molecular markers associated with COVID-19 immunopathology and multi-organ damage.

536  *Embo j* 39, e105896

537  14. Schwanhausser, B., Busse, D., Li, N., Dittmar, G., Schuchhardt, J., Wolf, J., Chen, W., and

538  Selbach, M. (2011) Global quantification of mammalian gene expression control. *Nature* 473,

539     337-342

540     15. Wilkerson, M. D., and Hayes, D. N. (2010) ConsensusClusterPlus: a class discovery tool

541     with confidence assessments and item tracking. *Bioinformatics* 26, 1572-1573

542     16. Langfelder, P., and Horvath, S. (2008) WGCNA: an R package for weighted correlation

543     network analysis. *BMC Bioinformatics* 9, 559

544     17. Zhang, B., and Horvath, S. (2005) A general framework for weighted gene co-expression

545     network analysis. *Stat Appl Genet Mol Biol* 4, Article17

546     18. Langfelder, P., and Horvath, S. (2007) Eigengene networks for studying the relationships

547     between co-expression modules. *Bmc Syst Biol* 1

548     19. Niu, L., Geyer, P. E., Wewer Albrechtsen, N. J., Gluud, L. L., Santos, A., Doll, S., Treit, P.

549     V., Holst, J. J., Knop, F. K., Vilsboll, T., Junker, A., Sachs, S., Stemmer, K., Muller, T. D.,

550     Tschop, M. H., Hofmann, S. M., and Mann, M. (2019) Plasma proteome profiling discovers

551     novel proteins associated with non-alcoholic fatty liver disease. *Mol Syst Biol* 15, e8793

552     20. Cendrowski, J., Kaczmarek, M., Mazur, M., Kuzmicz-Kowalska, K., Jastrzebski, K.,

553     Brewinska-Olchowik, M., Kominek, A., Piwocka, K., and Miaczynska, M. (2020) Splicing

554     variation of BMP2K balances abundance of COPII assemblies and autophagic degradation in

555     erythroid cells. *Elife* 9

556     21. Hazell, G. G. J., Peachey, A. M. G., Teasdale, J. E., Sala-Newby, G. B., Angelini, G. D.,

557     Newby, A. C., and White, S. J. (2016) PI16 is a shear stress and inflammation-regulated

558     inhibitor of MMP2. *Sci Rep-Uk* 6

559     22. Thannickal, V. J. (2013) Mechanistic links between aging and lung fibrosis. *Biogerontology*

560     14, 609-615

561   23. O'Dwyer, D. N., and Moore, B. B. (2017) The role of periostin in lung fibrosis and airway

562   remodeling. *Cell Mol Life Sci* 74, 4305-4314

563   24. Kim, Y. I., Ahn, J. M., Sung, H. J., Na, S. S., Hwang, J., Kim, Y., and Cho, J. Y. (2016)

564   Meta-markers for the differential diagnosis of lung cancer and lung disease. *J Proteomics* 148,

565   36-43

566   25. Yu, G., Tzouvelekis, A., Wang, R., Herazo-Maya, J. D., Ibarra, G. H., Srivastava, A., de

567   Castro, J. P. W., DeIuliis, G., Ahangari, F., Woolard, T., Aurelien, N., Arrojo, E. D. R., Gan, Y.,

568   Graham, M., Liu, X., Homer, R. J., Scanlan, T. S., Mannam, P., Lee, P. J., Herzog, E. L., Bianco,

569   A. C., and Kaminski, N. (2018) Thyroid hormone inhibits lung fibrosis in mice by improving

570   epithelial mitochondrial function. *Nat Med* 24, 39-49

571   26. Qiu, H., Song, J., Hu, J., Wang, L., Qiu, L., Liu, H., Lin, G., Luan, X., Liu, Y., and He, J.

572   (2022) Low serum transthyretin levels predict stroke-associated pneumonia. *Nutr Metab*

573   *Cardiovasc Dis* 32, 632-640

574   27. Ge, S., Xia, X., Ding, C., Zhen, B., Zhou, Q., Feng, J., Yuan, J., Chen, R., Li, Y., Ge, Z., Ji,

575   J., Zhang, L., Wang, J., Li, Z., Lai, Y., Hu, Y., Li, Y., Li, Y., Gao, J., Chen, L., Xu, J., Zhang, C.,

576   Jung, S. Y., Choi, J. M., Jain, A., Liu, M., Song, L., Liu, W., Guo, G., Gong, T., Huang, Y., Qiu,

577   Y., Huang, W., Shi, T., Zhu, W., Wang, Y., He, F., Shen, L., and Qin, J. (2018) A proteomic

578   landscape of diffuse-type gastric cancer. *Nat Commun* 9, 1012

579   28. Lubbers, R., Sutherland, J. S., Goletti, D., de Paus, R. A., Dijkstra, D. J., van Moorsel, C.

580   H. M., Veltkamp, M., Vestjens, S. M. T., Bos, W. J. W., Petrone, L., Malherbe, S. T., Walzl, G.,

581   Gelderman, K. A., Groeneveld, G. H., Geluk, A., Ottenhoff, T. H. M., Joosten, S. A., and Trouw,

582   L. A. (2020) Expression and production of the SERPING1-encoded endogenous complement

583 regulator C1-inhibitor in multiple cohorts of tuberculosis patients. *Mol Immunol* 120, 187-195

584 29. Hu, Q., Wang, Q., Han, C. G., and Yang, Y. (2020) Sufentanil attenuates inflammation and

585 oxidative stress in sepsis-induced acute lung injury by downregulating KNG1 expression. *Mol*

586 *Med Rep* 22, 4298-4306

587 30. Wang, D. B., Li, X., Lu, X. K., Sun, Z. Y., Zhang, X., Chen, X., Ma, L., and Xia, H. G.

588 (2021) Transthyretin Suppressed Tumor Progression in Nonsmall Cell Lung Cancer by

589 Inactivating MAPK/ERK Pathway. *Cancer Biother Radio*

590 31. Hamai, K., Iwamoto, H., Ishikawa, N., Horimasu, Y., Masuda, T., Miyamoto, S.,

591 Nakashima, T., Ohshimo, S., Fujitaka, K., Hamada, H., Hattori, N., and Kohno, N. (2016)

592 Comparative Study of Circulating MMP-7, CCL18, KL-6, SP-A, and SP-D as Disease Markers

593 of Idiopathic Pulmonary Fibrosis. *Dis Markers* 2016, 4759040

594 32. Ma, J., Chen, T., Wu, S. F., Yang, C. Y., Bai, M. Z., Shu, K. X., Li, K. L., Zhang, G. Q., Jin,

595 Z., He, F. C., Hermjakob, H., and Zhu, Y. P. (2019) iProX: an integrated proteome resource.

596 *Nucleic Acids Res* 47, D1211-D1217

597 **Figure legends**

598 Fig. 1 Proteomic features of the IPF subgroups. Molecular subtyping of IPF was based on

599 altered proteomes and their correlations with clinical features.

600 A. Cumulative number of proteins identified in serum samples from 30 healthy controls (blue

601     dots) and 30 patients with IPF (red dots).

602 B. Numbers of identified proteins in serum samples from 30 healthy controls (blue dots) and

603     30 IPF patients (red dots).

604 C. Relative abundance of 2,314 serum proteins. Several proteins ranged widely in abundance

605 (black dots).

606 D. Consensus clustering analysis of the proteomic profiling identifying three subtypes in the

607 IPF cohort.

608 E. Kaplan–Meier analyses of overall survival (OS) of patients in the S-I (n=16), S-II (n=4),

609 and S-III (n=10) subgroups. (P-values calculated by two-sided log-rank tests).

610 F. Heat map of the over-represented proteins in the three IPF subtypes.

611 G. Proteins differentially expressed in the three IPF subtypes.

612 H. Associations between expression of BMP2K, PI16, and ATP5A1 proteins, and overall

613 survival (Kaplan–Meier analysis, P-value from log-rank test, high means IPF/N >median

614 value).

615 I. Age with the three IPF proteomic subtypes (P-values calculated by Fisher's exact tests).

616 Fig. 2 WGCNA identification of modules of highly correlated genes and assessment of their

617 relationships to clinical variables.

618 A. Heatmap of the weighted gene co-expression network. The plot indicates the TOM among

619 all genes analyzed. Genes in columns and their corresponding rows are hierarchically

620 clustered by cluster dendrograms, which are presented along the top and left side of the plot.

621 B. Module-trait relationships between six modules and ten clinical traits.

622 C. Heatmap of the change in genes in the module of age.

623 D. Heatmap of the age-related genes in the three subgroups.

624 E. Associations of HSP90AB and CAMKK1 expression with clinical outcomes in 30 IPF

625 patients.

626 Fig. 3 Aberrantly expressed metabolic enzymes involved in enhanced glycolysis in serum

27

627      proteomes of patients with IPF.

628    A. Pathway schematic showing DEPs (t-test, $p < 0.05$) mapped onto glucose metabolism

629       pathways.

630    B. Boxplots showing proteins differentially expressed by IPF patients with normal and above-

631       normal levels of serum lactate (P-values calculated by t-test).

632    C. Violin plots of LDHA and LDHB expression in 30 healthy controls (blue dots) and 30 IPF

633       patients (red dots).

634    D. Associations of LDHA expression with clinical outcomes in IPF patients (p-values

635       calculated by log-rank tests).

636    E. Violin plots of CCT6A expression in 30 healthy controls (blue dots) and 30 IPF patients

637       (red dots).

638    F. ELISA validation of CCT6A expression in IPF patients (P-values calculated by t-tests).

639    G. IHC staining showing CCT6A expression in lungs from healthy controls and IPF patients.

640    H. IHC staining showing CCT6A expression in the bleomycin model of lung fibrosis in mice.

641    I. Representative immunoblots of whole lung lysates of mice incubated with antibodies

642       against CCT6A and GAPDH.

643    J. Western blots of CCT6A expression normalized to β-actin. * $P < 0.05$, as determined by

644       ANOVA.

645    Fig. 4 Association of changes in CCT6A expression and high lactic acid concentrations with

646    the fibroblast phenotype.

647    A. Representative immunoblots showing CCT6A and α-SMA expression in MCR5 cells

648    transfected with control plasmid and plasmid overexpressing CCT6A.

28

649    B. Western blots of CCT6A expression normalized to β-actin. * P < 0.05, as determined by

650    ANOVA.

651    C. Representative images of α-SMA immunofluorescence staining of MRC5 cells. Original

652    magnification, ×100. Scale bars: 5 μm.

653    D. Representative immunoblots showing CCT6A, COLA1, and FN expression in MCR5 cells

654    transfected with control and CCT6A siRNAs.

655    E. Western blots of CCT6A expression normalized to β-actin. * P < 0.05, ** P < 0.01, as

656    determined by ANOVA.

657    F. ECAR of control and CCT6A-overexpressing MRC5 cells.

658    G-H. Lactate production in the supernatants of MRC5 cells and in the lungs of bleomycin mice.

659    I. Pyruvate production in MRC5 cells.

660    J. Expression of LDHA mRNA in MRC5 cells overexpressing CCT6A.

661    K. Representative immunoblots showing LDHA expression in MRC5 cells overexpressing

662    CCT6A.

663    L. Western blots of LDHA expression normalized to β-actin. * P < 0.05, as determined by

664    ANOVA.

665    Fig. 5 Machine-learning-based selection of biomarker combinations for classification of IPF.

666    A.  Receiver operating characteristic (ROC) curve for the classification model. Calculation of

667       AUC values in the patient cohort by 5-fold cross-validation. Confusion matrix of the four-

668       protein combination in the patient cohort.

669    B.  ROC curve for the test model Calculation of AUC values in the public cohort by 5-fold

670       cross-validation. Confusion matrix of the four-protein combination in the public cohort.

29

671     C. Associations between the protein combinations and clinical outcomes in 30 IPF patients of

672        the classification model.

673     D. Heatmap of the combination biomarkers in the public cohort (PRIDE project PXD010965).

674     E. ELISA determination of TTR expression in an independent cohort. (P-values calculated by

675        t-tests).

676     F. Correlation between TTR expression and patient age in the study cohort.

677

678

679 **Supplemental figures legends**

680 Figure S1 Schematic of the proteomic analyses of serum samples from 30 IPF patients and 30

681 healthy controls.

682 Figure S2 Profiling of serum proteomics of IPF patients and healthy controls.

683 A. Quality control of mass spectrometry using a tryptic digest of HEK293T cells. The

684 bottom-left half of the panel shows the pairwise Spearman's correlation coefficients

685 of the samples, and the top-right half of the panel depicts the pairwise scatter plots

686 from the same comparisons.

687 B. Assessment of study quality by analysis of the protein markers FGA, FGB, and

688 FGG.

689 C. Heatmap of the altered proteins in healthy controls and IPF samples.

690 D. Volcano plot of differentially expressed genes differing significantly in non-IPF and IPF

691 samples. The log2 differential expression ratio and the -log10 (p-value) were plotted for each

692 gene. Proteins with differential expression ratios >2 or <2 were defined as those significantly

693 up- and downregulated, respectively.

694

695 Figure S3 Proteomic subtypes of IPF with their molecular characteristics.

696 A. Consensus clustering plus identification of three serum proteomic subtypes of IPF

697 samples. The panel shows a consensus matrix of 30 IPF samples from k=2 to k=6,

698 with k=3 considered the ideal value based on visual inspection of the consensus

699 matrix and the change in area under the CDF.

700 B. Top 30 exclusively expressed proteins in S-I, S-II, and S-III patients.

701 Figure S4 Correlation of FOT value for CCT6A(MS data) with ELISA value of CCT6A in the

702 same sera of IPF patients; Pearson correlation, P = 0.0056. n = 17 human samples.

703 Figure S5 Identification of combination of biomarkers by machine-learning.

704 A. Workflow of the machine-learning.

705 B. TTR, KNG1 and FCN3 expression in the public cohort (left) and the study cohort (right).

31

Table 1 Information of IPF cohort and healthy control cohort

| Characteristics | Control | IPF |
|---|---|---|
| Number | 30 | 30 |
| Age | 61.07 ± 9.85 | 64.50 ± 10.58 |
| Gender | | |
| Male | 17 | 22 |
| Female | 13 | 8 |
| Smoking | 8 | 11 |

# Figure 1

Figure 2

A



B

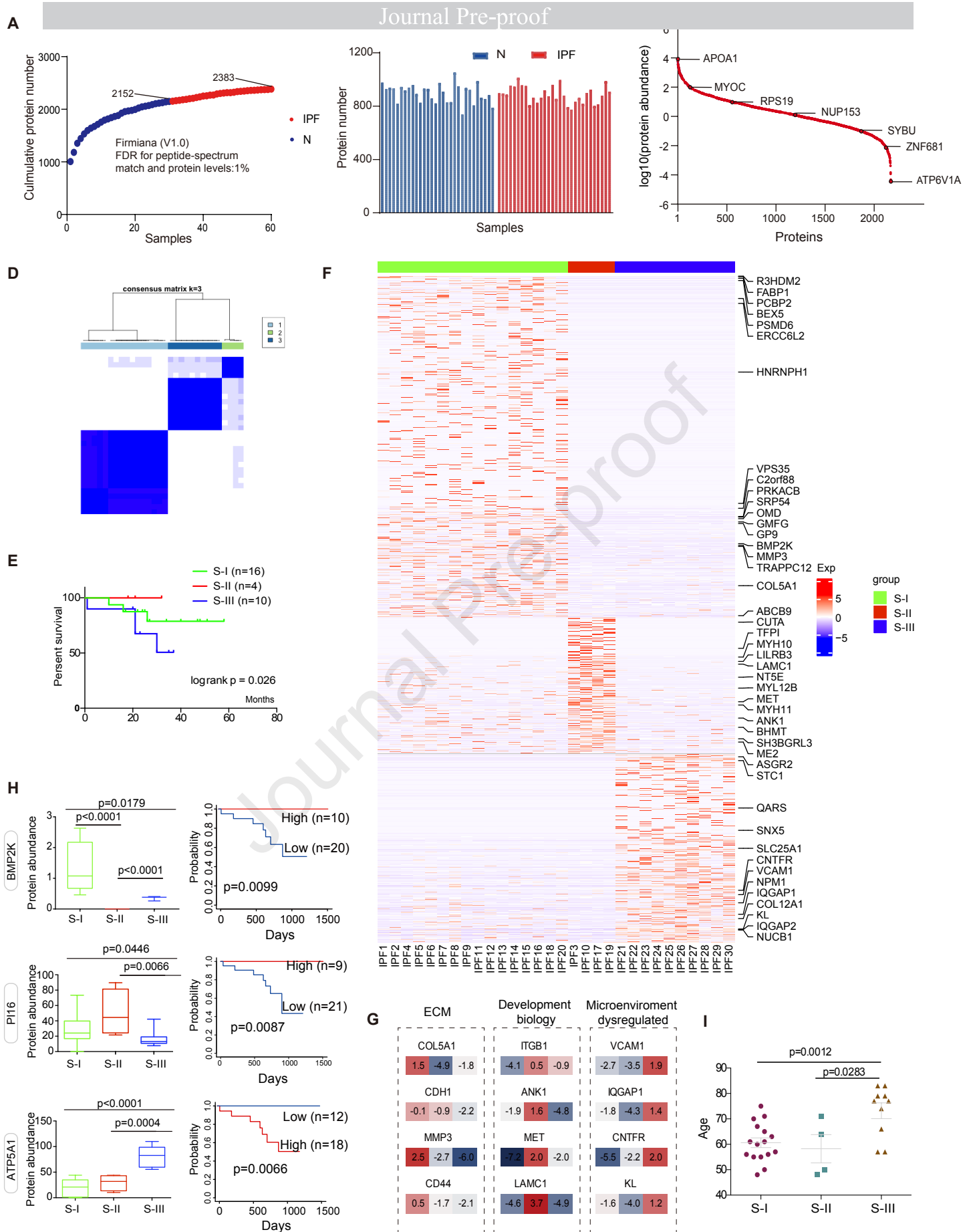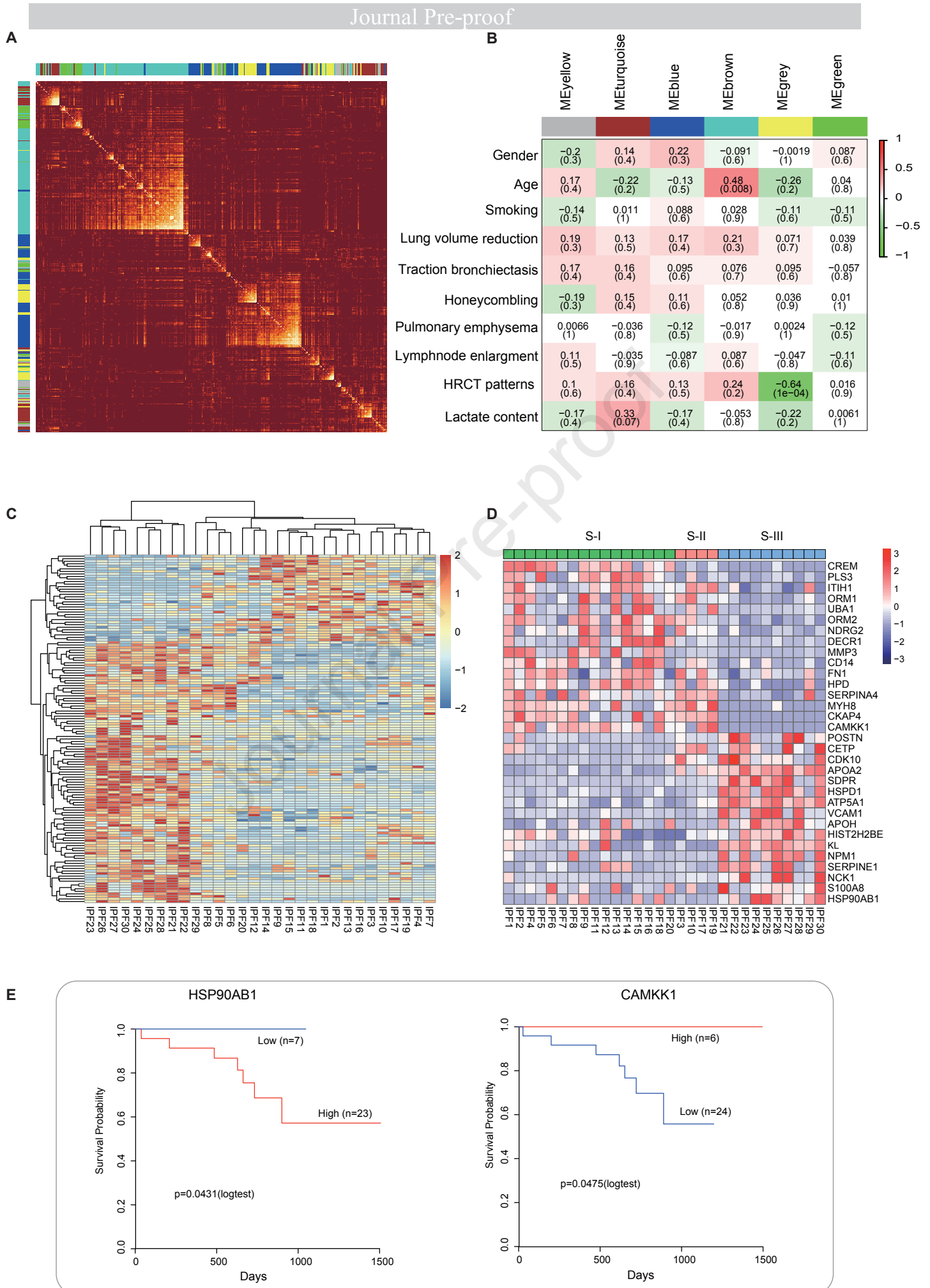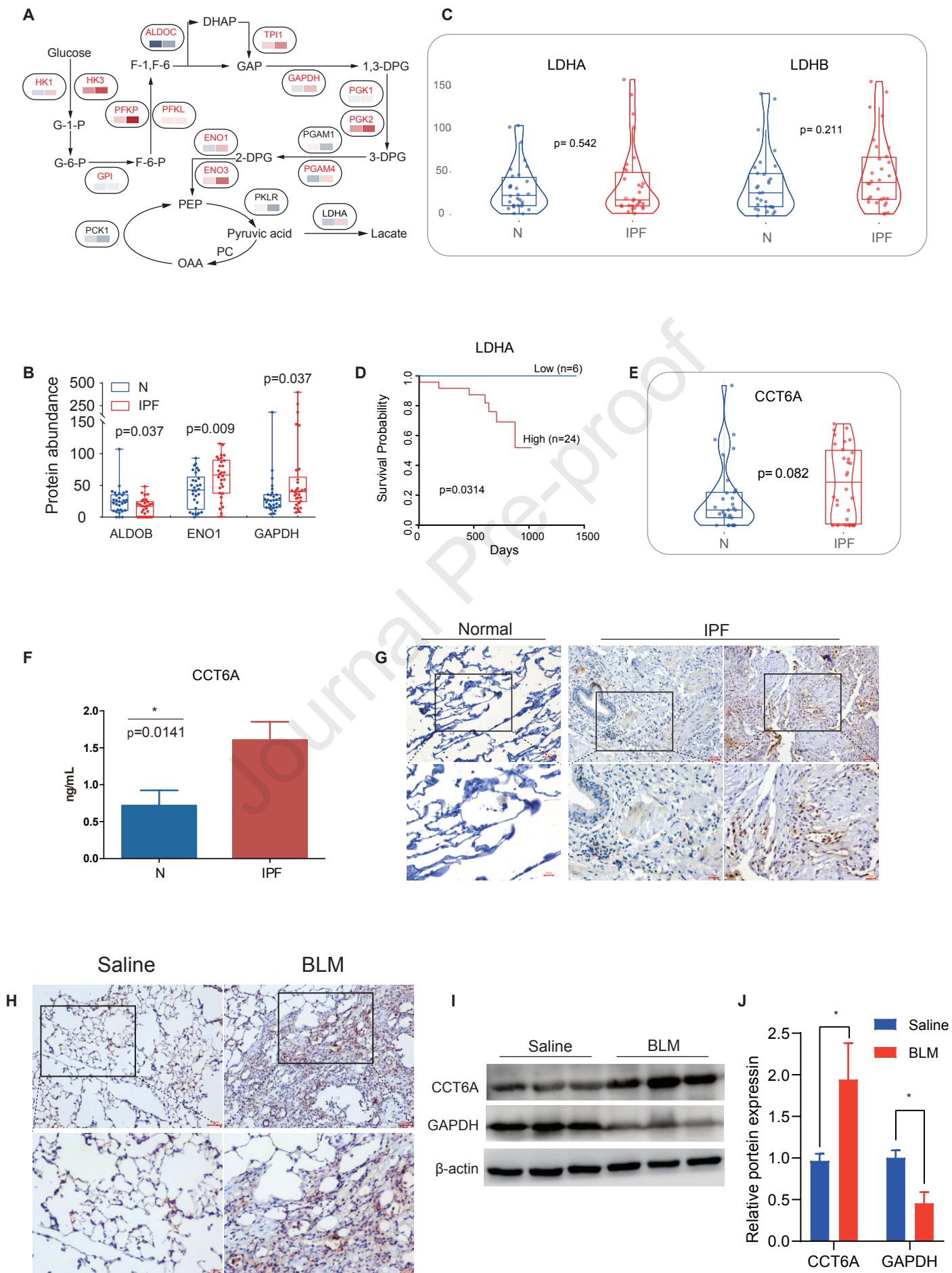|  | MEyellow | MEturquoise | MEblue | MEbrown | MEgrey | MEgreen |
|---|---|---|---|---|---|---|
| Gender | -0.2 (0.3) | 0.14 (0.4) | 0.22 (0.3) | -0.091 (0.6) | -0.0019 (1) | 0.087 (0.6) |
| Age | 0.17 (0.4) | -0.22 (0.2) | -0.13 (0.5) | 0.48 (0.008) | -0.26 (0.2) | 0.04 (0.8) |
| Smoking | -0.14 (0.5) | 0.011 (1) | 0.088 (0.6) | 0.028 (0.9) | -0.11 (0.6) | -0.11 (0.5) |
| Lung volume reduction | 0.19 (0.3) | 0.13 (0.5) | 0.17 (0.4) | 0.21 (0.3) | 0.071 (0.7) | 0.039 (0.8) |
| Traction bronchiectasis | 0.17 (0.4) | 0.16 (0.4) | 0.095 (0.6) | 0.076 (0.7) | 0.095 (0.6) | -0.057 (0.8) |
| Honeycombling | -0.19 (0.3) | 0.15 (0.4) | 0.11 (0.6) | 0.052 (0.8) | 0.036 (0.9) | 0.01 (1) |
| Pulmonary emphysema | 0.0066 (1) | -0.036 (0.8) | -0.12 (0.5) | -0.017 (0.9) | 0.0024 (1) | -0.12 (0.5) |
| Lymphnode enlargment | 0.11 (0.5) | -0.035 (0.9) | -0.087 (0.6) | 0.087 (0.6) | -0.047 (0.8) | -0.11 (0.6) |
| HRCT patterns | 0.1 (0.6) | 0.16 (0.4) | 0.13 (0.5) | 0.24 (0.2) | -0.64 (1e-04) | 0.016 (0.9) |
| Lactate content | -0.17 (0.4) | 0.33 (0.07) | -0.17 (0.4) | -0.053 (0.8) | -0.22 (0.2) | 0.0061 (1) |

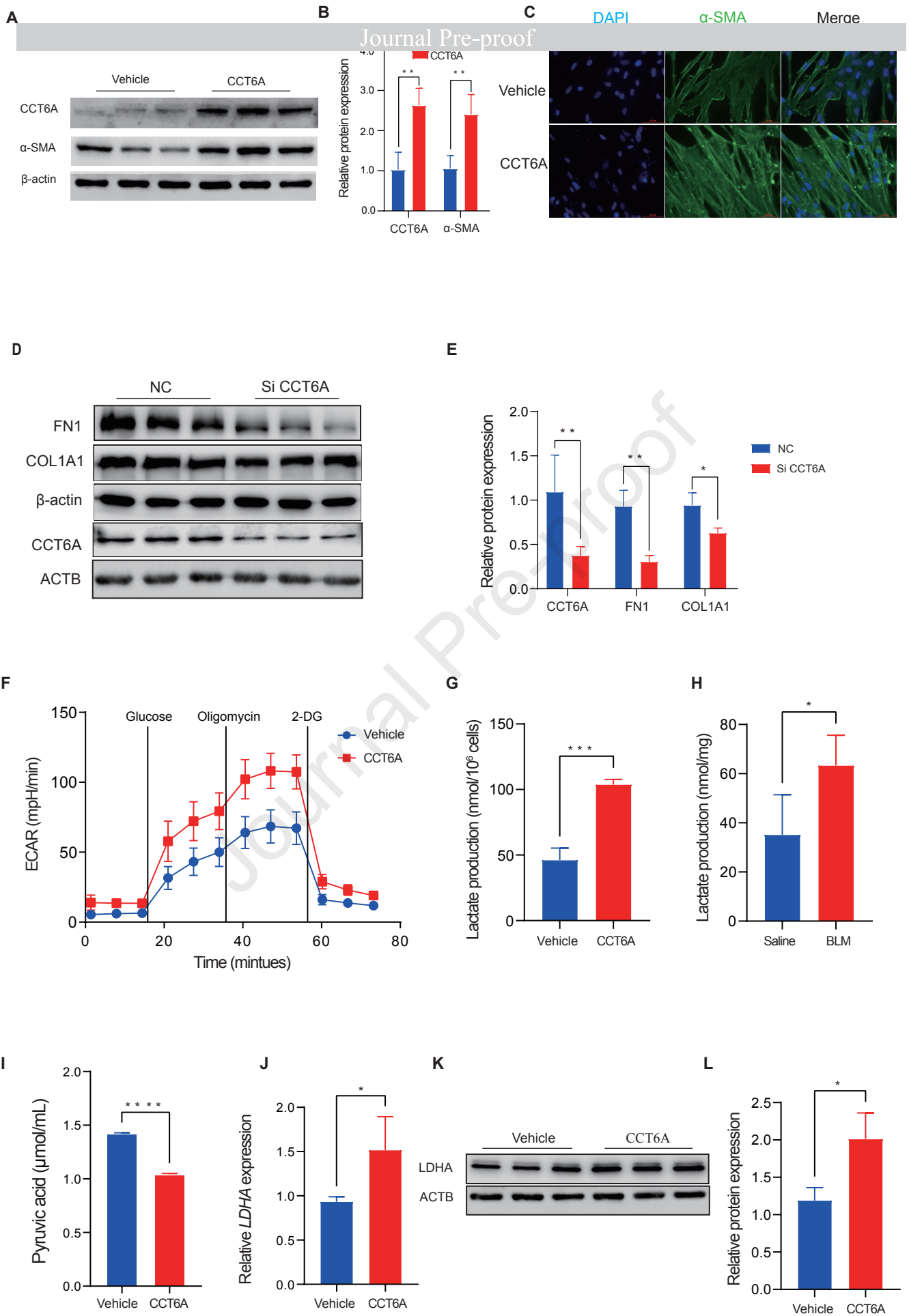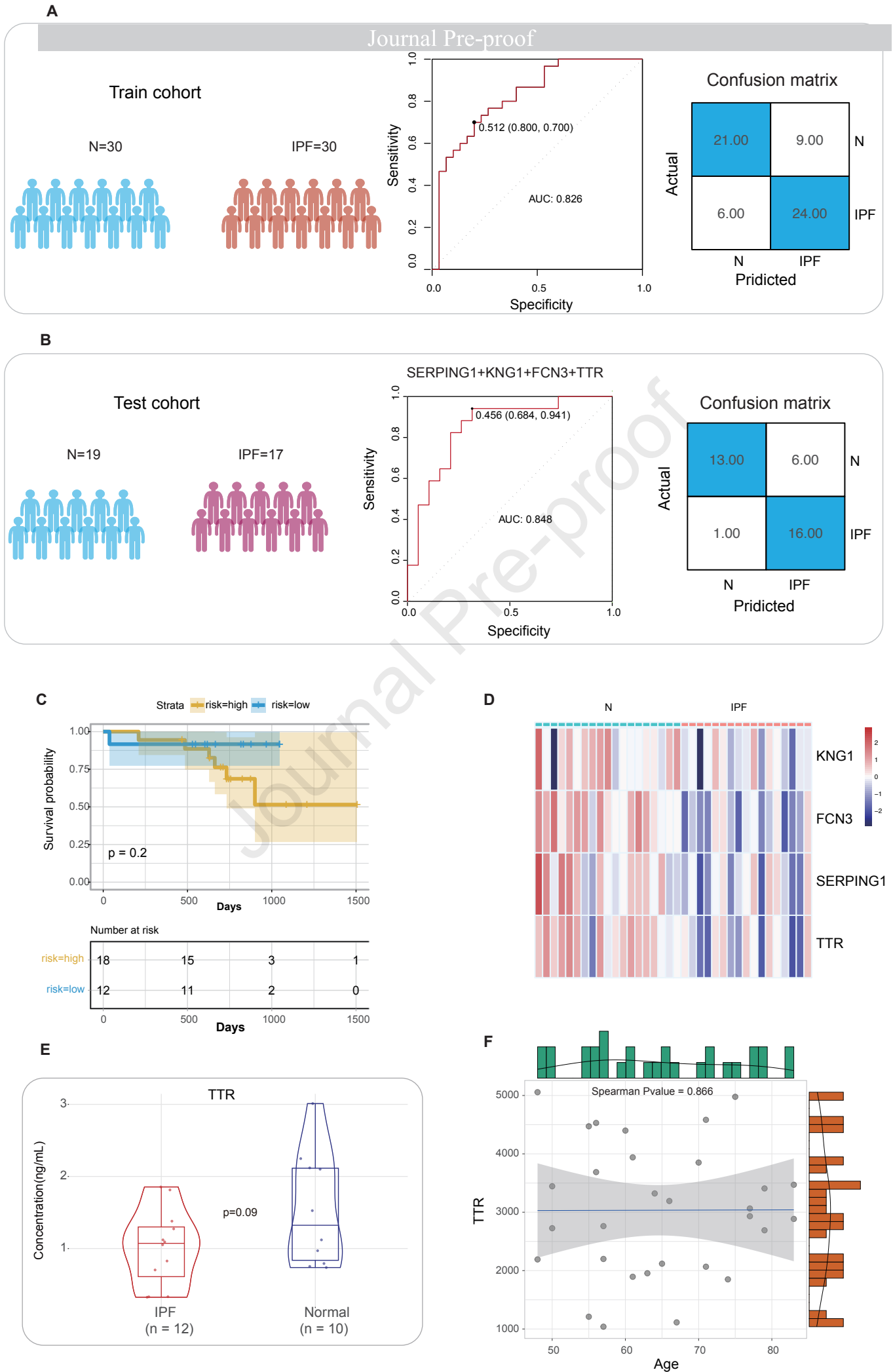C



D



E

Figure 3

Figure 4

Figure 5

Highlights

● A serum proteome profiling by DIA-MS identified 2833 gene products from IPF and normal subjects, three subgroups were distinguished in IPF patients in signal pathways and overall survival.

● Aging-associated signatures in module MEturquoise were identified by WGCNA coincidently falling into S-III which provided clear and direct evidence that aging is a critical risk factor for IPF rather than to a single biomarker

● LDHA and CCT6A expression, which were associated with glucose metabolic reprogramming, were correlated with high serum lactic acid content in the patients with IPF.

● Cross-model analysis and machine learning showed that a combinatorial biomarker accurately distinguished IPF patients from healthy subjects and validated from another cohort and ELISA assay.

# AUTHOR CONTRIBUTION STATEMENT

G.Y. and C.D.: Designed the research plan. L. W.: Data curation, Writing- Original draft preparation, L.W. Y. L., X. Ch.: Proteomics experiments, Z. L.,H.Z: Statistical analysis, Data analysis and data visualization, S.Y.: Performed cell and mouse assay and related data visualization, IHC staining of tissue samples, J.Y and X.P.: Performed ELISA assay, H.Y. and M.Z. :Consulted on clinical questions. I.R. Writing – review & editing.

All authors discussed the results and commented on the manuscript.

Brief

Wang et el (2022) performed serum proteomics by DIA-MS and identified 2833 gene products from IPF and normal subjects, and distinguished in IPF patients into three subgroups in signal pathways and overall survival. Aging-associated signatures by WGCNA coincidently provided clear and direct evidence that aging is a critical risk factor for IPF rather than to a single biomarker. LDHA and CCT6A expression, which were associated with glucose metabolic reprogramming, were correlated with high serum lactic acid content in the patients with IPF. Cross-model analysis and machine learning showed that a combinatorial biomarker accurately distinguished IPF patients from healthy subjects and validated from another cohort and ELISA assay.

**Declaration of interests**

☒ The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

☐The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: