1 **Proteomic network analysis of bronchoalveolar lavage fluid in ex-smokers to discover**

2 **implicated protein targets and novel drug treatments for chronic obstructive pulmonary**

3 **disease**

4 Manoj J. Mammen[1, 6], Chengjian Tu[2,3], Matthew C. Morris[5], Spencer Richman[5], William

5 Mangione[6], Zackary Falls[6], Jun Qu[2,3], Gordon Broderick[5], Sanjay Sethi[1,4], Ram Samudrala[6]

6

7 [1]Department of Medicine, Jacobs School of Medicine and Biological Sciences, State University

8 of New York at Buffalo, Buffalo, NY 14214 USA; [2]Department of Pharmaceutical Sciences,

9 State University of New York at Buffalo, Buffalo, NY 14260 USA; [3] New York State Center of

10 Excellence in Bioinformatics and Life Sciences, 701 Ellicott Street, Buffalo, NY 14203 USA;

11 [4]WNY VA Healthcare System, Buffalo, NY 14215 USA; [5]Center for Clinical Systems Biology,

12 Rochester General Hospital, Rochester, NY 14621 USA; Department of Biomedical Informatics,

13 Jacobs School of Medicine and Biological Sciences, State University of New York, Buffalo, NY

14 14214 USA; [6]Department of Biomedical Informatics, Jacobs School of Medicine and Biological

15 Sciences, State University of New York, Buffalo, NY 14214 USA

16

17    ***Corresponding Authors:**

18    Manoj J. Mammen, MD, MS

19    100 High Street, B-8

20    Buffalo, NY 14203

21    Phone: (716) 859-2271

22    Email: mammen@buffalo.edu

23    Ram Samudrala, PhD

24    77 Goodell St.

25    Buffalo, NY 14203

26    Phone: (206) 251-8852

27    Email: rams@buffalo.edu

28    **Authors and Contributors:** MJM, CT, JQ, MCM, WM, ZF, SR, GB, SS, RS

44    **MANDATORY DISCLAIMER**

45    The opinions and assertions contained herein are the private views of the authors and are not to
46    be construed as official or as reflecting the views of the Department of Defense.
47

48    To obtain the raw proteomic data on a DVD media, please contact Dr. Jun Qu, junqu@buffalo.edu.

49

50    **Running title:** Computational analysis of BALF proteome to repurpose drugs for COPD.

51    **Word Count:**

52    **Keywords:** proteomics, drug repurposing, translational bioinformatics, interaction signature,

53    bronchoalveolar lavage fluid, chronic obstructive pulmonary disease, label-free quantitation,

54    plasma, serum, biomarker

55

56    **Abbreviations:** BALF=bronchoalveolar lavage fluid; BANDOCK= bioanalytical docking;

57    CANDO=computational analysis of novel drug opportunities; COPD= chronic obstructive

58    pulmonary disease; DAVID=database for annotation visualization and integrated discovery;

59    GO= gene ontology; IPA= ingenuity pathway analysis; LTQ Orbitrap=linear ion trap combined

60    with an orbitrap analyzer mass spectrometer.

61

62

63    **Abstract**

64    **Rationale:** Bronchoalveolar lavage of the epithelial lining fluid can sample the profound

65    changes in the airway lumen milieu prevalent in  Chronic Obstructive Pulmonary Disease

66    (COPD). Characterizing the proteins in bronchoalveolar lavage fluid in COPD with advanced

67    proteomic methods will identify disease-related changes, provide insight into pathogenetic

68    mechanisms and potential therapeutics that will aid in the discovery of more effective

69    therapeutics for COPD.

70

71    **Objectives:** We compared  epithelial lining fluid proteome of ex-smokers with moderate COPD

72    who are  not in exacerbation status COPD, to non-smoking healthy control subjects using

73    advanced proteomics methods and applied proteome-scale translational bioinformatics

74    approaches to identify potential therapeutic protein targets and drugs that modulate these proteins

75    towards the treatment of COPD.

76

77    **Methods:** Proteomic profiles of bronchalveolar lavage fluid were obtained from 1) never-smoker

78    control subjects with normal lung function (n=10) or 2) individuals with stable moderate (GOLD

79    stage 2, $FEV_1$ 50% – 80% predicted) COPD who were ex-smokers for at least one year (n=10).

80    NIH's Database for Annotation, Visualization and Integrated Discovery (DAVID) and

81    Ingenuity's Ingenuity Pathway Analysis (IPA) were the two bioinformatics tools employed for

82    network analysis on the differentially expressed proteins to identify potential crucial hub

83    proteins. The drug-proteome interaction signature comparison and ranking approach

84    implemented in the Computational Analysis of Novel Drug Opportunities (CANDO) platform

85    for multiscale therapeutic discovery was utilized to identify potential repurposable drugs for the

86    treatment of COPD based on the BALF proteome. Subsequently, a literature-based knowledge

87    graph was utilized to rank combinations of drugs that would most likely ameloriate inflammatory

88    processes  by inhibition or activation of their functions.

89

90    **Results:** Proteomic network analysis demonstrated that 233 of the >1800 proteins identified in

91    the BALF were differentially expressed in COPD versus control, including proteins associated

92    with inflammation, structural elements, and energy metabolism. Functional annotation of the

93    differentially expressed proteins by their implicated biological processes, cellular localization,

94    and transcription factor interactions was accomplished via DAVID. Canonical pathways

95    containing the differential expressed proteins were detailed via the Ingenuity Pathway Analysis

96    application. Topological network analysis demonstrated that four proteins act as central node

97    proteins in the inflammatory pathways in COPD. The CANDO multiscale drug discovery

98    platform was used to analyze the behavioral similarity between the interaction signatures of all

99    FDA-approved drugs and the identified BALF proteins. The drugs with the signatures most

100   similar interaction signatures to approved COPD drugs were extracted with the CANDO

101   platform.  The analysis revealed 189  drugs that putatively target the proteins implicated in

102   COPD. The putative COPD drugs that were identified using CANDO were subsequently

103   analyzed using a knowledge based technique to identify an optimal two drug combination that

104   had the most appropriate effect on the central node proteins.

105

106   **Conclusion:** Analysis of the BALF proteome revealed novel differentially expressed proteins in

107   the epithelial lining fluid that elucidate COPD pathogenesis. Network analyses identified critical

108   targets that have critical roles in modulating COPD pathogenesis, for which we identified several

109    drugs that could be repurposed to treat COPD using a multiscale shotgun drug discovery

110    approach.

111

112

**Introduction**

Chronic obstructive pulmonary disease (COPD) is a leading cause of mortality and morbidity in the US.[1-5] Additionally, COPD results in millions of hospitalizations in the developing world.[1,6-11] The prevalence of cigarette smoking continues to rise in most developing countries around the world.[12-14] However, only 25-50% of tobacco smokers develop COPD, suggesting only a subset develops an exaggerated inflammatory process that leads to lung destruction.[12,13,15] Bronchoalveolar lavage fluid (BALF) and bronchial samples from ex-smokers reveal active inflammation long after smoking cessation.[16]

Although structural changes in the airways, parenchyma, and pulmonary vessels are typical in patients with COPD, the lower airways and the alveoli are the initial sites of the inflammatory process.[17,18] The inflammatory process initiated by smoking persists after cessation and is likely exaggerated by autoimmunity and infection.[19,20] Accurate and precise measurement of the molecular mediators in the airways should aid in rigorous analysis of their role in disease.

There has been a keen interest in understanding the genetic determinants of COPD, as the interaction between genes and environment leads to protein expression, ultimately resulting in either healthy or disease states. However, genomic data alone does not predict protein abundance or activity; proteins are the ultimate participants in integrated biological processes that lead to healthy physiological function or pathology. Proteome-based analysis of bronchoalveolar lavage

133    fluid (BALF) in COPD can identify tissue-specific markers of inflammation that can lead to

134    understanding the mechanisms of COPD progression.

135

136    We sought to determine an unbiased proteome-based analysis of BALF in COPD under stable

137    conditions (not in exacerbation status) to identify a broad series of molecules involved in COPD

138    pathogenesis. A label-free proteomics mass spectroscopy method was utilized. The differentially

139    expressed proteins were analyzed using multiple bioinformatics tools  to critical pathways that

140    were altered in these ex-smoker patients with COPD compared to healthy, never smoker

141    controls, proteins implicated in COPD etiology, and to identify putative drug candidates that can

142    be repurposed to treat COPD.

143

144    The raw proteomic data used in this manuscript was initially detailed in a previously published

145    methodology manuscript using strict criteria (2 peptide identification criteria for a protein, ≥1.5

146    fold change, and p-value<0.05) to identify 423 individual proteins with 76 proteins expressed

147    differently between COPD and controls.[21] In this analysis, we adopted a pragmatic approach to

148    the same raw proteomic data (1 peptide identification criterion, ≥1.5 fold change, and p-

149    value<0.05) that identified 1831 individual proteins and 233 differentially expressed proteins

150    between the two groups.  The latter, more practical, approach provides important additional

151    information for biomarker and therapeutic target discovery that may be utilized in future research

152    to discover useful interventions.

153  **Methods**

154

155  We analyzed the protein quantifications derived from the BALF of subjects with COPD and

156  healthy ex-smoker control subjects via liquid chromatography and mass spectroscopy.  We then

157  used pathway analysis tools to identify relevant cellular pathways associated with differentially

158  expressed proteins quantified from the BALF analysis.  We subsequently employed  the

159  Computational Analysis of Novel Drug Opportunities (CANDO) platform to identify FDA

160  approved drugs that could be repurposed to COPD, based on their putative interaction with the

161  differentially expressed proteins. Using topological network analysis, we identified putative hub

162  proteins that modulate the cellular pathways associated with COPD. Using the medical literature

163  to predict the repurposed drugs effects on the most important hub protein, we created a refined

164  list of drugs predicted to modulate the cellular pathway in order to impede COPD pathogenesis.

165  to generate proteomic interaction signatures for the compounds

166

167  **Recruitment of subjects**

168  BALF was obtained in a NHLBI funded study of innate lung defense in COPD.[22] All procedures

169  received approval from the Institutional Review Board (IRB), Veterans Affairs Western New

170  York Healthcare System (WNY-VA), and strictly adhered to institutional guidelines.

171

172  **Ethics statement**

173
174  This study is a sub-study of a larger group of patients with COPD and healthy controls to

175  understand biological determinants of exacerbation frequency and was approved by the

176  Institutional Review Boards of the Veterans Affairs Western New York Healthcare System and

177    University at Buffalo. The participants gave written consent to the study via an IRB-approved

178    consent form.

179

180    **Inclusion/exclusion criteria**

181
182    The inclusion criteria and procedures for this study have been described previously and are

183    provided in the supplementary material. [22] After informed consent, 116 volunteers were divided

184    into three groups: 1) healthy nonsmokers, 2) ex-smokers with COPD, and 3) active smokers with

185    COPD and underwent bronchoscopy and bronchoalveolar lavage. The methodology for

186    bronchoscopy, lavage, and sample processing is included in the supplementary material.

187

188    For this study, we selected BALF obtained from ten ex-smokers with moderate COPD and ten

189    healthy non-smoking controls for proteomic analysis, respectively.  To minimize variability due

190    to effects of acute smoking and disease severity, we confined this analysis to ex-smokers and

191    moderate  stage 2 disease per the Global Obstructive Lung Disease (GOLD)[23] criteria of the

192    forced expiratory volume in 1 second (FEV1) 50-80% predicted.  All ex-smokers had ceased

193    smoking for at least one year.

194

**Bronchoscopy and BALF sample preparation**

196    The research bronchoscopy and BALF sample preparation were performed as described

197    previously.[24]

198

**Protein identification/quantification.**

200    To investigate the soluble molecules in the epithelial lining fluid that may participate in COPD

201    pathogenesis, unbiased proteomic analysis of BALF commenced without protein depletion or

202    fractionation. Details of the methodology have been published[25] and are also provided in the

203    supplementary material.

204

**Long gradient nano-RPLC/mass spectrometry**

206    Complete separation of the complex peptide mixture utilized a nano-LC/nanospray setup;[26] the

207    ion-current long gradient approach with mass spectrometry and subsequent protein identification

208    was performed as described in Tu, et al. [25-27] All proteins identified with one or more peptide

209    hits,  fold change of ≥1.5, and p-value <0.05 are included as part of the differentially expressed

210    BALF proteome.

211

212 **Bioinformatics analyses**

213 *Manually curated pathway analysis*

214 Gene ontology, transcription factors, and expression locations were determined by uploading the

215 protein expression dataset onto a web-based tool, the NIH's Database for Annotation,

216 Visualization and Integrated Discovery (DAVID) v6.7 (http://david.abcc.ncifcrf.gov/). [28,29]

217 Biological networks were generated with Ingenuity Pathway Analysis (IPA, Ingenuity Systems),

218 a web-based relational database and network generator. Proteins overrepresented in the uploaded

219 datasets in biological networks, canonical pathways, and biological processes were identified.

220
221 *Literature informed protein-protein and protein-drug interaction network*

222 In addition to annotating differentially expressed proteins with the manually curated pathways

223 cataloged in IPA, a network of protein-protein interactions was created using known regulatory

224 relationships extracted from published scientific literature using the MedScan text-mining

225 engine[30] as well as protein-drug interactions cataloged in the Reaxsys medicinal chemistry

226 database (Elsevier, Amsterdam). These are embedded in the broader Elsevier Knowledge Graph

227 database [31] and were accessed via the Pathway Studio interface (Elsevier, Amsterdam).[32]

228
229 *Shotgun multiscale drug discovery platform*

230 We used the Computational Analysis of Novel Drug Opportunities (CANDO) platform[33-40] to

231 predict drugs that can be repurposed for the treatment of stable COPD. In CANDO, a

232 compound/drug is considered to be potentially repurposable for an indication when it is found to

233    have similar binding interactions with a specific proteome or library of proteins as a drug with

234    known approval for the indication of interest.

235

236    In this study, we calculated the interaction scores between 2,450 United States Federal Drug

237    Administration (FDA) approved drugs from the CANDO version 2.3 compound library and a

238    curated human library of 8,385 proteins, including 5,316 solved X-ray crystallography structures

239    and 3,069 computed protein structures modeled by I-TASSER[41,42].  The interaction scores were

240    calculated using the bioanalytic docking (BANDOCK) protocol in the CANDO  which utilizes

241    predicted binding site information and chemical similarity to determine an interaction score that

242    is a surrogate for the likelihood of interaction between a compound and protein.[33] Binding sites

243    were predicted for all human proteins using COACH [43], which uses the consensus of three

244    complementary methods utilizing structure and sequence information to find similarity to solved

245    structures in the Protein Data Bank (PDB).[44,45] For each binding site predicted by COACH, a

246    confidence score (PScore) and an associated co-crystallized ligand are output. The ligand is then

247    compared to the query compound/drug using chemical fingerprinting methods, which

248    enumerates the presence or absence of molecular substructures on the compound/drug.  The

249    Sorensen-Dice coefficient[46] between the protein-ligand and compound/drug fingerprints

250    (CScore) is also computed. The BANDOCK interaction score outputted for each compound-

251    protein pair  is the product of the Pscore and the Cscore.

252

253    For this analysis, we focused on the differentially expressed proteins in the BALF proteome (as

254    described), and drugs used to treat COPD ("MESH:D029424")  (Table S1). We selected proteins

255    in the CANDO human protein library that were also represented in the differentially expressed

256    BALF proteome. We then used the CANDO platform to  predict the top drug candidates that

257    could be repurposed to treat COPD based on compound-proteome interaction signature

258    similiarity to drugs currently approved/used to treat stable COPD. The protocol iterates through

259    34 known drugs used to treat stable COPD, counting the number of times drugs not associated

260    with COPD show up in the top30 most similar compounds to the known treatments, then outputs

261    the consensus  predictions ranked by the number of times each compound appeared across all

262    top30 lists. The similarity between a given drug and all other drugs in the library is determined

263    by comparing their proteomic interaction signatures using the cosine similarity metric, where

264    compounds with greater similarity scores rank stronger than those with low similarity. Thereby

265    drugs that were most similar (in terms of interaction signatures) to multiple drugs used to treat

266    COPD will be ranked highest.

267

268    *Network topological analysis*

269    Although not a complete descriptor, the topological location, and aspects of the connectivity

270    linking a node to a broader biological network can inform the node's function in mediating

271    network behavior. Among the measures of a node's importance or centrality, *betweenness*

272    centrality has been used to describe how a node might serve as an important mediator of

273    information flow in a regulatory network.  In this work, $C_b(n)$ for each node *n* of a network was

274    calculated using the Brandes algorithm.[47] The betweenness centrality of a node *n* reflects the

275    amount of control that this node exerts over the interaction between communities of neighboring

276    nodes in the network[48] and can be computed as follows:

277
$$C_b(n) = \sum_{s \neq t \neq n} \left( {\sigma_{s,t}(n)} / {\sigma_{s,t}} \right)$$
(1)

278    Where $s$ and $t$ are the source and target nodes in the network different from $n$, $\sigma_{s,t}$ denotes the

279    number of shortest paths from all $s$ to all $t$, and $\sigma_{s,t}(n)$ is the number of shortest paths from s to $t$

280    that must pass through node $n$. Here, unweighted betweenness centralities were calculated for

281    each node in the literature-informed protein-protein network. The betweenness centrality scores

282    for all nodes were expressed as fractions of the maximum betweenness centrality present in the

283    network. All calculations were conducted in R version 4.0.2.[49]

284

285    *Literature based drug enrichment analysis*

286    Using putative drugs ranked by CANDO and further analyzed via the Elsevier Knowledge

287    Graph,[31] a drug enrichment analysis was performed to predict which drugs can most closely mimic

288    an idealized intervention against the hub proteins identified in the network topological analysis.

289    Drugs are represented as vectors with a length equal to the empirically derived number protein

290    entities in the network model.  Each index value is listed as 0 if there is no interaction between the

291    drug and the corresponding model entity, a 1 if the drug promotes that entity, or a –1 if the drug

292    inhibits that entity. Next, the cosine similarity, $S_c$, between each drug vector and the idealized

293    intervention vector is calculated.[50] Cosine similarity is calculated as:

294
$$S_c(D, M) = \frac{D \cdot M}{\|D\| \|M\|}$$

295    Where $D$ is the drug vector and $M$ is the idealized intervention. Higher $S_c$ indicates a closer match

296    between the drug vector and the idealized vector. A $S_c$ of 1 means the two vectors are identical,

297     and -1 indicates that the two are exactly opposed. For multidrug combinations, the net-effect of

298     the individual drug vectors is calculated as:

299
$$sgn\left(\sum_{i=1}^{n} \overrightarrow{d_i} D_i\right)$$

300     Where $n$ is the total number of drugs in the combination, $D_i$ is the vector corresponding to the $i$th

301     drug, and $sgn$ is the sign function. The cosine similarity of the net-effect vector and idealized vector

302     is then calculated.

303         The statistical significance of these enrichment scores is determined empirically from an

304     estimated null distribution of cosine similarities. This null distribution uses a set of model-relevant

305    background drugs for which each interacts with at least one entity in the network. All CANDO

306    drugs of interest were included in the background. Empirical p-values are estimated as

307
$$\hat{p} = \frac{(r + 1)}{(n + 1)}$$

308    Where $r$ is the number of null $S_c$ values greater than the observed $S_c$ and $n$ is the total number of

309    null $S_c$ values.

310

311    **Statistical analysis**

312     Statistical analysis was performed with SPSS/19. Demographic values were depicted as mean ±

313    SEM.

314

315    **Results**

316    ***Study population characteristics***

317    Characteristics for subjects included in the BALF study are shown in ***Table 1***, with the only

318    significant differences between the two groups in tobacco smoke exposure and lung function.

319

320    ***BALF proteome characteristics***

321    A total of 1831 unique proteins were identified in the BALF proteome.  A total of 233 proteins

322    (>1.5-fold absolute change, p-value <0.05) had a significant differential expression in BALF

323    samples from patients with COPD versus healthy ex-smokers, 138 proteins were decreased in

324    COPD while 95 proteins were increased (***Table S2 and Table S3).***

325

326    ***Manually curated pathway analysis***

327    *Functional annotation of differential expressed proteins and transcription factor interactions*

328   The 233 differentially quantified proteins were characterized by their biological processes,

329   transcription factor interactions, and cellular localization by employing NIH's DAVID.[28,29] The

330   proteins involved in several biological processes implicated in COPD pathogenesis (total number

331   of proteins, number upregulated, number downregulated) such as proteolysis[51] (20,4,16),

332   extracellular matrix[52] (13,6,7), cell adhesion[53] (11,2,9), cytoskeleton[54] (32,14,18), defense

333   response[55] (16, 7,9), cell migration[56] (12,4,8), and oxidation-reduction[21] (11,2,9) were altered in

334   COPD. As expected with examining the lung lining fluid, the largest single group of

335   differentially expressed proteins was associated with the extracellular space (49, 30, 19).

336

337   Transcription factors (**Table S4**) associated with the differentially expressed proteins (total

338   number of proteins associated with the transcription factor) included serum response factor-SRF

339   (148), transcription factor 8-AREB6 (166), signal transducer and activator of transcription factor

340   1-STAT1 (69), zinc finger protein-GFI1 (97), signal transducer and activator of transcription

341   factor 3-STAT3 (101), nuclear factor kappa-light-chain-enhancer of activated B cells-NF-**κ**B

342   (79), CCAAT/enhancer-binding protein **β**-CEBPB(109), paired box gene 2-PAX2(113), and

343   activating transcription factor 2-CREBP1(95).

344

345   *Bioinformatic pathway analysis of BALF proteomic data*

346   The protein expression datasets were imported into IPA (Ingenuity Systems) and projected onto

347   the relevant biological pathways; processes linked to the differentially expressed proteins were

348   analyzed with IPA's manually curated knowledge database.  Of the 233 differentially expressed

349   proteins, 217 matched to the IPA curated database and were analyzed. Sixteen pathways were

350   noted to have several proteins associated with the differentially expressed BALF dataset (**Table**

351   **S5**), including proteins implicated in cellular movement, cellular death and survival, cell

352   morphology, immune cell trafficking, and cell cycle. Appendix **Figures S1 to S4** depict IPA

353   networks of selected pathways with the highest number of differentially expressed proteins.

354

355   *Computational drug prediction*

356   130 out of 233 BALF differentially expressed proteins were identified in the CANDO human

357   protein library. This subset of proteins within the CANDO platform was used to predict 189

358   putative drug candidates that have the most similar protein interaction signatures to the set of

359   known drugs used to treat COPD ( **Figure 1** and **Table S6**). Many of the drugs were

360   corticosteroids; however other putative drugs included tezacaftor[57], a recently developed drug to

361   potentiate sodium channel activity in the treatment of cystic fibrosis; two additional drugs

362   predicted to treat COPD, gemfibrozil[58], and pioglitazone[59], are drugs currently used to treat

363   hyperlipidemia and diabetes, respectively.

364

365   *Candidate key mediators of COPD pathology based on literature derived drug enrichment*

366   *Literature informed protein-protein and protein-drug interaction network*

367   A total of 233 proteins were identified as differentially expressed between COPD patients and

368   healthy controls by mass spectrometry. Of these, 214 were represented in the Elsevier

369   Knowledge Graph[31], with the remainder comprising specific immunoglobulin chain proteins. A

370   query of the Knowledge Graph for documented regulatory interactions between these protein

371   entities yielded 206 regulatory edges supported by 807 references (with a median of 1 reference

372   per edge). 112 of the 214 identified proteins could not be connected to the broader network

373   circuit by a documented interaction. The protein entities in this network were then assessed in

374     terms of their importance as mediators of signal transfer based on their betweenness centrality.

375     (**Figure 2**).

376

377     *Network topological analysis*

378     Four nodes representing proteins in the network stood out based on the normalized betweenness

379     centrality values representing a greater than linear increase from the next lower ranking node:

380     fibronectin, vimentin, intercellular adhesion molecule 1 (ICAM1), and galectin-3. These

381     potentially key signaling mediators had a betweenness centrality of at least 25% of the

382     maximum.

383

384     Analysis of the initial data reveals fibronectin and ICAM1 are reduced in COPD patients relative

385     to healthy controls; thus, any candidate therapeutic should target an increase in their activity. The

386     reverse is true for vimentin and galectin-3. We, therefore, sought drugs or combinations of drugs

387     predicted to accomplish the appropriate activation or inhibition of the four most central nodes,

388     specifically drugs that will lead to the promotion of central node proteins that were

389     downregulated in the COPD cohort and inhibition of central node proteins that were

390     overabundant in COPD. The idealized drug vector, therefore, constitutes interactions leading to

391     desirable modulation of the central hub protein.   CANDO identified 189 distinct drugs (**Figure**

392     **1, Table S7**S6) with relevance for COPD; 39 of these represented in the Elsevier Knowledge

393     Graph [31] were analyzed for their enrichment for the desired agonist and antagonist effects on the

394     most central entities in the protein regulatory network. Highly enriched drugs or drug pairs were

395     are predicted to be more likely than randomly selected drugs to exert appropriate inhibition or

396     promotion of the most central proteins. Two single drugs (fluocinolone acetonide and

397     dexrazoxane) and 57 two-drug combinations were significantly enriched. Fluocinolone acetonide

398     and dexrazoxane appeared in 54% and 46% of all significantly enriched 2-drug combinations,

399     respectively, far greater than the other drugs appearing in these combinations (Figure 3). The

400     combination of fluocinolone acetonide and dexrazoxane is the most enriched two-drug

401     combination leading to an idealized drug vector.

402
403     We additionally conducted a targeted query to assess the predicted effects of drugs commonly

404     applied in pulmonary disease treatment on the most central proteins of this regulatory network

405     (Figure 4 and Table S7).  While some of these have been documented to have the desired effect

406     on two of the central proteins, fibronectin or vimentin, all have been documented to have the

407     opposite effect on at least one of the most central proteins. Therefore, they were not significantly

408     enriched out of the set of all possible candidate drugs.

409

410    **Discussion**

411    Our investigation of the COPD BALF proteome utilizing novel bioinformatic techniques

412    revealed significant differences in proteins involved in multiple biological processes, including

413    lung-specific mechanisms, protease/anti-protease homeostasis, immunoregulation, and the

414    extracellular matrix. Proteomic profiling of the complex pathways implicated in COPD provides

415    broader physiological exploration not provided by studying one entity at a time.  We identified

416    several differentially expressed proteins in COPD versus controls that, based on a review of

417    published literature, have not been previously implicated in COPD etiology. This preliminary

418    analysis illustrates how our BALF proteomic analysis represents a powerful approach to

419    elucidate COPD pathogenesis and identify novel biomarkers.

420

421    Employing the bioinformatics tool DAVID and IPA, putative pathway networks were

422    constructed based on the differentially expressed proteins in the BALF proteome that implicated

423    multiple transcription factor pathways and disparate biological processes, such as extracellular

424    space, proteolysis, extracellular matrix, cell adhesion, cytoskeleton, defense response, cell

425    migration, and oxidation-reduction.

426

427    The CANDO platform identified 189 drug candidates that had similar protein interaction

428    signatures based on the BALF proteome when compared to known drugs that are used to treat

429    COPD. However, while most putative drug and protein interactions are likely inhibitors, the

430    induction or inhibition of a target protein is indeterminable with solely the binding potential

431    between drug and protein pairs.

432

433    Topological analysis of the interaction network connecting 233 proteins differentially expressed

434    in COPD through regulatory interactions documented in the literature suggested that ICAM1[60]

435    and galectin-3[61] are important central mediators of inflammation while both fibronectin[62,63] and

436    vimentin[64] are central mediators of inflammation and fibrogenesis. This corroborates the results

437    of the pathway enrichment analysis described above, and points to fibrosis and innate

438    inflammation as important processes governing the pathogenesis and progression of COPD. A

439    literature knowledge-based query (Elsevier Knowledge Graph) of drugs with desired drug-target

440    interactions (generated using CANDO) identified putative drugs, such as anti-neoplastic, anti-

441    fibrotic drugs, and regulators of inflammation, that would restore key central proteins to the

442    levels characteristic of healthy controls. Our results also suggest currently utilized medications

443    for COPD have disparate effects on the identified central node proteins that are key putative

444    mediators of COPD pathogenesis and progression.  In contrast,  the corticosteroid fluocinolone

445    acetonide[65] and the cardioprotective agent dexrazoxane[66] were highly enriched for the desired

446    effects on central network entities, both individually and in combination. Fluocinolone acetonide

447    is a stronger potentiator than other corticosteroids of the TGF-β pathway[67] which is noted to be

448    dysregulated in COPD[68], and fluocinolone acetonide may be more effective than comparable

449    corticosteroids in improved homeostasis in that pathway. Dexrazoxane[66] is used to reduce

450    cardiac toxicity associated with anthracycline-based chemotherapy agents by binding to iron and

451    reducing reactive oxygen species; with oxidative stress as a significant factor in COPD

452    pathogenesis[69], antioxidative therapy may be beneficial.

453

454    The documented actions of these immunomodulators were predicted here to substantially

455    counteract the observed dysregulation of centrally-connected proteins in COPD patients. The

456    relatively high representation of immunomodulators among the candidate agents and the

457    increased centrality of fibrosis-related proteins is consistent with the paradigm of airway

458    remodeling as central to COPD pathology.[70] With additional data, this regulatory circuit could be

459    used as a testbed for computational evaluation of these and other candidate drug effects using

460    network topological methods. [71]

461

462    *Limitations and strengths*

463    Our approach does has some limitations.  The variability in how much BALF is recovered from

464    each aliquot of saline infused in to the lower airway in COPD vs. control subjects are inherent in

465    most BALF proteomic analyses. However, the BALF proteins were normalized to albumin

466    BALF concentrations to account for the variability. The examination of protein levels without

467    accounting for post-translational modifications, such as phosphorylation, may neglect important

468    differences in protein interactions and activity, despite no significant differences in protein

469    levels. Also, the BALF samples were from subjects in the COPD group who were ex-smokers.

470    This exclusion limits the generalizability of our findings particularly current smokers, since the

471    acute effects of tobacco smoke were excluded in our study design.

472

473

474    However, we confined our analysis to ex-smokers with moderate COPD to obtain some

475    uniformity of the COPD phenotype and to avoid the acute inflammatory effects of current

476    smoking.  Future work on proteomic profiles will inform us of the difference between such

477    profiles in current smokers and different stages of COPD.

478

479    *Comparison to previously published studies*

480    A sputum proteomics study endeavored to identify COPD severity biomarkers by employing 2D

481    gel electrophoresis and revealed 15 proteins that were significantly differentially expressed

482    between healthy smoker controls and subjects with GOLD stage II; subsequently, 9 of the 15

483    candidate proteins were validated with Western Blot. Of the nine candidate proteins validated

484    with Western Blot, seven were statistically significantly different between groups, specifically

485    albumin, alpha-2-HS glycoprotein, transthyretin, PSP94, apolipoprotein A1, lipocalin-1, and

486    PLUNC. [72] Employing quantitative ELISA data normalized for protein content, the investigators

487    identified apolipoprotein A1 and lipocalin-1 as statistically differentially expressed in COPD.

488    Although apolipoprotein A1 and lipocalin-1 were identified in our study of the BALF proteome,

489    the proteins were not significantly differentially expressed, likely due to the differences in

490    expression in the different biocompartments of sputum vs. bronchoalveolar lumen.

491    A 2D differential gel electrophoresis study and subsequent mass spectroscopy were performed

492    by Ohlmeier et al., which compared healthy smokers, non-smokers, and smokers with GOLD

493    stage II COPD and revealed a different set of 15 proteins that were differentially expressed

494    between the groups.[73] Of these proteins, polymeric immunoglobulin receptor levels in lung tissue

495    and blood between the three groups were correlated with airflow obstruction.

496

497     In Lee et al., tumor-free lung tissue harvested from patients with lung cancer resection, when

498     examined via 2D gel electrophoresis/MALDI-TOF-MS, revealed eight proteins that were

499     upregulated in subjects with COPD compared to nonsmokers and ten significantly differentially

500     expressed proteins between subjects with COPD and smoking subjects without COPD.[74] Two of

501     the identified proteins, matrix metalloprotease 13 (MMP13) and thioredoxin-like 2, were

502     confirmed to be increased in COPD subjects with Western Blot and immunohistochemical

503     staining, with MMP13 localized to the alveolar macrophage and type II pneumocytes and

504     thioredoxin-like 2 found in the bronchial epithelium. Thioredoxin-like 2, which contains

505     thioredoxin, was found in the BALF proteome but not significantly differentially expressed.

506     However, MMP13 was not identified in our BALF study, likely due to differences in study

507     populations and variable biocompartments.

508

509     **Conclusion**

510     In summary, our work provides a valuable pipeline for identifying many proteins associated

511     with COPD pathogenesis that illustrate the complexity of the development of this disease, as

512     well as identifying putative therapeutic treatment options using cutting-edge bioinformatics

513     approaches.  Identifying differentially expressed proteins will form the basis for future

514     mechanistic studies of critical pathways and novel treatment discovery. Validation of our

515     proposed therapeutic approach in animal models and pilot human studies are important next

516     steps.

517

520
521

522

## REFERENCES

1.  Brown DW, Croft JB, Greenlund KJ, Giles WH. Deaths From Chronic Obstructive Pulmonary Disease-- United States, 2000-2005. *JAMA: Journal of the American Medical Association.* 2009;301(13):1331.

2.  From the Global Strategy for the Diagnosis, Management and Prevention of COPD. Global Initiative for Chronic Obstructive Lung Disease (GOLD). *Available from:* http://www.goldcopd.org2010.

3.  Murray CJ, Lopez AD. Mortality by cause for eight regions of the world: Global Burden of Disease Study. *Lancet.* 1997;349(9061):1269-1276.

4.  Kobayashi S, Shibata H, Yokota K, et al. FHL2, UBC9, and PIAS1 are novel estrogen receptor alpha-interacting proteins. *Endocrine research.* 2004;30(4):617-621.

5.  Croft JB, Wheaton AG, Liu Y, et al. Urban-Rural County and State Differences in Chronic Obstructive Pulmonary Disease - United States, 2015. *MMWR Morb Mortal Wkly Rep.* 2018;67(7):205-211.

6.  Geitona M, Hatzikou M, Steiropoulos P, Alexopoulos EC, Bouros D. The cost of COPD exacerbations: a university hospital--based study in Greece. *Respir Med.* 2011;105(3):402-409.

7.  Dalal AA, Shah M, D'Souza AO, Rane P. Costs of COPD exacerbations in the emergency department and inpatient setting. *Respiratory Medicine.* 2011;105(3):454-460.

8.  Dalal AA, Christensen L, Liu F, Riedel AA. Direct costs of chronic obstructive pulmonary disease among managed care patients. *Int J Chron Obstruct Pulmon Dis.* 2010;5:341-349.

9.  Dalal AA, Shah M, D'Souza AO, Rane P. Costs of inpatient and emergency department care for chronic obstructive pulmonary disease in an elderly Medicare population. *J Med Econ.* 2010;13(4):591-598.

10. Hutchinson A, Brand C, Irving L, Roberts C, Thompson P, Campbell D. Acute care costs of patients admitted for management of chronic obstructive pulmonary disease exacerbations: contribution of disease severity, infection and chronic heart failure. *Intern Med J.* 2010;40(5):364-371.

11. Ciapponi A, Alison L, Agustina M, Demian G, Silvana C, Edgardo S. The epidemiology and burden of COPD in Latin America and the Caribbean: systematic review and meta-analysis. *COPD.* 2014;11(3):339-350.

12. Rennard SI, Vestbo J. COPD: the dangerous underestimate of 15%. *Lancet.* 2006;367(9518):1216-1219.

13. Rennard SI, Vestbo J. Natural histories of chronic obstructive pulmonary disease. *Proc Am Thorac Soc.* 2008;5(9):878-883.

14. Schubert C. Anti-tobacco efforts going up in smoke. *Nat Med.* 2006;12(8):866-866.

15. Mammen MJ, Sethi S. COPD and the microbiome. *Respirology.* 2016;21(4):590-599.

16. Rutgers SR, Postma DS, ten Hacken NH, et al. Ongoing airway inflammation in patients with COPD who do not currently smoke. *Thorax.* 2000;55(1):12-18.

17. Hogg JC. Pathophysiology of airflow limitation in chronic obstructive pulmonary disease. *Lancet.* 2004;364(9435):709-721.

18. Hogg JC, Chu F, Utokaparch S, et al. The nature of small-airway obstruction in chronic obstructive pulmonary disease. *N Engl J Med.* 2004;350(26):2645-2653.

19. Cosio MG, Saetta M, Agusti A. Immunologic aspects of chronic obstructive pulmonary disease. *N Engl J Med.* 2009;360(23):2445-2454.

20. Sethi S. Infection as a comorbidity of COPD. *Eur Respir J.* 2010;35(6):1209-1215.

21. Tu C, Mammen MJ, Li J, et al. Large-scale, ion-current-based proteomics investigation of

568  bronchoalveolar lavage fluid in chronic obstructive pulmonary disease patients. *J Proteome Res.*
569  2014;13(2):627-639.

570  22.  Berenson CS, Garlipp MA, Grove LJ, Maloney J, Sethi S. Impaired phagocytosis of nontypeable
571      Haemophilus influenzae by human alveolar macrophages in chronic obstructive pulmonary
572      disease. *J Infect Dis.* 2006;194(10):1375-1384.

573  23.  Singh D, Agusti A, Anzueto A, et al. Global strategy for the diagnosis, management, and
574      prevention of chronic obstructive lung disease: the GOLD science committee report 2019. *Eur.*
575      *Resp. J.* 2019;53(5).

576  24.  Berenson CS, Wrona CT, Grove LJ, et al. Impaired alveolar macrophage response to Haemophilus
577      antigens in chronic obstructive lung disease. *Am J Respir Crit Care Med.* 2006;174(1):31-40.

578  25.  Tu C, Mammen MJ, Li J, et al. Large-Scale, Ion-Current-Based Proteomics Investigation of
579      Bronchoalveolar Lavage Fluid in Chronic Obstructive Pulmonary Disease Patients. *Journal of*
580      *Proteome Research.* 2014;13(2):627-639.

581  26.  Tu C, Li J, Young R, et al. Combinatorial peptide ligand library treatment followed by a dual-
582      enzyme, dual-activation approach on a nanoflow liquid chromatography/orbitrap/electron
583      transfer dissociation system for comprehensive analysis of swine plasma proteome. *Analytical*
584      *chemistry.* 2011;83(12):4802-4813.

585  27.  Tu C, Li J, Bu Y, Hangauer D, Qu J. An ion-current-based, comprehensive and reproducible
586      proteomic strategy for comparative characterization of the cellular responses to novel anti-
587      cancer agents in a prostate cell model. *Journal of proteomics.* 2012;7:187-201.

588  28.  Huang da W, Sherman BT, Lempicki RA. Systematic and integrative analysis of large gene lists
589      using DAVID bioinformatics resources. *Nat Protoc.* 2009;4(1):44-57.

590  29.  Huang da W, Sherman BT, Lempicki RA. Bioinformatics enrichment tools: paths toward the
591      comprehensive functional analysis of large gene lists. *Nucleic acids research.* 2009;37(1):1-13.

592  30.  Novichkova S, Egorov S, Daraselia N. MedScan, a natural language processing engine for
593      MEDLINE abstracts. *Bioinformatics.* 2003;19(13):1699-1706.

594  31.  Kamdar MR, Stanley CE, Carroll M, et al. Text snippets to corroborate medical relations: an
595      unsupervised approach using a knowledge graph and embeddings. *AMIA Summits on*
596      *Translational Science Proceedings.* 2020;2020:288.

597  32.  Cheadle C, Cao H, Kalinin A, Hodgkinson J. Advanced literature analysis in a Big Data world. *Ann*
598      *N Y Acad Sci.* 2017;1387(1):25-33.

599  33.  Falls Z, Mangione W, Schuler J, Samudrala R. Exploration of interaction scoring criteria in the
600      CANDO platform. *BMC research notes.* 2019;12(1):318.

601  34.  Schuler J, Samudrala R. Fingerprinting CANDO: Increased Accuracy with Structure- and Ligand-
602      Based Shotgun Drug Repurposing. *ACS Omega.* 2019;4(17):17393-17403.

603  35.  Mangione W, Falls Z, Chopra G, Samudrala R. cando.py: Open Source Software for Predictive
604      Bioanalytics of Large Scale Drug-Protein-Disease Data. *J Chem Inf Model.* 2020;60(9):4131-4136.

605  36.  Hudson ML, Samudrala R. Multiscale Virtual Screening Optimization for Shotgun Drug
606      Repurposing Using the CANDO Platform. *Molecules.* 2021;26(9).

607  37.  Schuler J, Falls Z, Mangione W, Hudson ML, Bruggemann L, Samudrala R. Evaluating the
608      performance of drug-repurposing technologies. *Drug Discov Today.* 2021.

609  38.  Horst JA, Pieper U, Sali A, et al. Strategic protein target analysis for developing drugs to stop
610      dental caries. *Adv Dent Res.* 2012;24(2):86-93.

611  39.  Minie M, Chopra G, Sethi G, et al. CANDO and the infinite drug discovery frontier. *Drug Discov*
612      *Today.* 2014;19(9):1353-1363.

613  40.  Sethi G, Chopra G, Samudrala R. Multiscale modelling of relationships between protein classes
614      and drug behavior across all diseases using the CANDO platform. *Mini Rev Med Chem.*
615      2015;15(8):705-717.

616  41.  Yang J, Yan R, Roy A, Xu D, Poisson J, Zhang Y. The I-TASSER Suite: protein structure and function
617       prediction. *Nature methods.* 2015;12(1):7-8.
618  42.  Yang J, Zhang Y. I-TASSER server: new development for protein structure and function
619       predictions. *Nucleic acids research.* 2015;43(W1):W174-181.
620  43.  Yang J, Roy A, Zhang Y. Protein-ligand binding site recognition using complementary binding-
621       specific substructure comparison and sequence profile alignment. *Bioinformatics.*
622       2013;29(20):2588-2595.
623  44.  Dutta S, Burkhardt K, Young J, et al. Data deposition and annotation at the worldwide protein
624       data bank. *Mol Biotechnol.* 2009;42(1):1-13.
625  45.  Gore S, Sanz Garcia E, Hendrickx PMS, et al. Validation of Structures in the Protein Data Bank.
626       *Structure.* 2017;25(12):1916-1927.
627  46.  Dice LR. Measures of the Amount of Ecologic Association Between Species. *Ecology.*
628       1945;26(3):297-302.
629  47.  Brands U. A faster algorithm for betweenness cen-trality [J]. *Journal of Mathematical Sociology.*
630       2001;25(2):163-177.
631  48.  Gyorgy A. A Practical Step-by-Step Guide for Quantifying Retroactivity in Gene Networks.
632       *Synthetic Gene Circuits*: Springer; 2021:293-311.
633  49.  Team RC. R: A language and environment for statistical computing. 2013.
634  50.  Thada V, Jaglan V. Comparison of jaccard, dice, cosine similarity coefficient to find best fitness
635       value for web retrieved documents using genetic algorithm. *International Journal of Innovations*
636       *in Engineering and Technology.* 2013;2(4):202-205.
637  51.  Shapiro SD. Proteolysis in the lung. *The European respiratory journal. Supplement.* 2003;44(44
638       suppl):30s-32s.
639  52.  Annoni R, Lancas T, Yukimatsu Tanigawa R, et al. Extracellular matrix composition in COPD. *Eur*
640       *Respir J.* 2012;40(6):1362-1373.
641  53.  Riise GC, Larsson S, Lofdahl CG, Andersson BA. Circulating cell adhesion molecules in bronchial
642       lavage and serum in COPD patients with chronic bronchitis. *Eur Respir J.* 1994;7(9):1673-1677.
643  54.  Yang M, Kohler M, Heyder T, et al. Proteomic profiling of lung immune cells reveals
644       dysregulation of phagocytotic pathways in female-dominated molecular COPD phenotype.
645       *Respir Res.* 2018;19(1):39.
646  55.  Sethi S. Bacterial infection and the pathogenesis of COPD. *Chest.* 2000;117(5 Suppl 1):286S-
647       291S.
648  56.  Stockley RA. Neutrophils and the pathogenesis of COPD. *Chest.* 2002;121(5 Suppl):151S-155S.
649  57.  Davies JC, Moskowitz SM, Brown C, et al. VX-659–tezacaftor–ivacaftor in patients with cystic
650       fibrosis and one or two Phe508del alleles. *New england journal of medicine.* 2018;379(17):1599-
651       1611.
652  58.  Frick MH, Elo O, Haapa K, et al. Helsinki Heart Study: primary-prevention trial with gemfibrozil in
653       middle-aged men with dyslipidemia. *New England Journal of Medicine.* 1987;317(20):1237-
654       1245.
655  59.  Eckland D, Danhof M. Clinical pharmacokinetics of pioglitazone. *Experimental and clinical*
656       *endocrinology & diabetes.* 2000;108(Sup. 2):234-242.
657  60.  Walter RE, Wilk JB, Larson MG, et al. Systemic inflammation and COPD: the Framingham Heart
658       Study. *Chest.* 2008;133(1):19-25.
659  61.  Mammen MJ, Sands MF, Abou-Jaoude E, et al. Role of Galectin-3 in the pathophysiology
660       underlying allergic lung inflammation in a tissue inhibitor of metalloproteinases 1 knockout
661       model of murine asthma. *Immunology.* 2018;153(3):387-396.
662  62.  Muro AF, Moretti FA, Moore BB, et al. An essential role for fibronectin extra type III domain A in
663       pulmonary fibrosis. *Am J Respir Crit Care Med.* 2008;177(6):638-645.

63. To WS, Midwood KS. Plasma and cellular fibronectin: distinct and independent functions during tissue repair. *Fibrogenesis & tissue repair.* 2011;4(1):21.

64. Li FJ, Surolia R, Li H, et al. Autoimmunity to vimentin is associated with outcomes of patients with idiopathic pulmonary fibrosis. *The Journal of Immunology.* 2017;199(5):1596-1605.

65. Mills J, Bowers A, Djerassi C, Ringold H. Steroids. CXXXVII. 1 Synthesis of a New Class of Potent Cortical Hormones. 6α, 9α-Difluoro-16α-hydroxyprednisolone and its Acetonide. *Journal of the American Chemical Society.* 1960;82(13):3399-3404.

66. Wiseman LR, Spencer CM. Dexrazoxane. A review of its use as a cardioprotective agent in patients receiving anthracycline-based chemotherapy. *Drugs.* 1998;56(3):385-403.

67. Hara ES, Ono M, Pham HT, et al. Fluocinolone acetonide is a potent synergistic factor of TGF-β3–associated chondrogenesis of bone marrow–derived mesenchymal stem cells for articular surface regeneration. *Journal of Bone and Mineral Research.* 2015;30(9):1585-1596.

68. Di Stefano A, Sangiorgi C, Gnemmi I, et al. TGF-beta Signaling Pathways in Different Compartments of the Lower Airways of Patients With Stable COPD. *Chest.* 2018;153(4):851-862.

69. Kirkham PA, Barnes PJ. Oxidative stress in COPD. *Chest.* 2013;144(1):266-273.

70. Moisieieva NV, Burya LV, Kapustianskaya AA, Kolenko IA, Rumyantseva MA, Shumeiko OH. Comprehensive patterns of comorbidity: copd and depression. Aspects of treatment. *Wiad Lek.* 2018;71(3 pt 1):588-591.

71. Morris MC, Richman S, Lyman CA, et al. Hacking the Immune Response to Infection in Chronic Obstructive Pulmonary Disease. Paper presented at: 2020 IEEE 20th International Conference on Bioinformatics and Bioengineering (BIBE)2020.

72. Nicholas BL, Skipp P, Barton S, et al. Identification of lipocalin and apolipoprotein A1 as biomarkers of chronic obstructive pulmonary disease. *Am J Respir Crit Care Med.* 2010;181(10):1049-1060.

73. Ohlmeier S, Mazur W, Linja-aho A, et al. Sputum Proteomics Identifies Elevated PIGR levels in Smokers and Mild-to-Moderate COPD. *Journal of Proteome Research.* 2011;11(2):599-608.

74. Lee EJ, In KH, Kim JH, et al. Proteomic analysis in lung tissue of smokers and COPD patients. *Chest.* 2009;135(2):344-352.
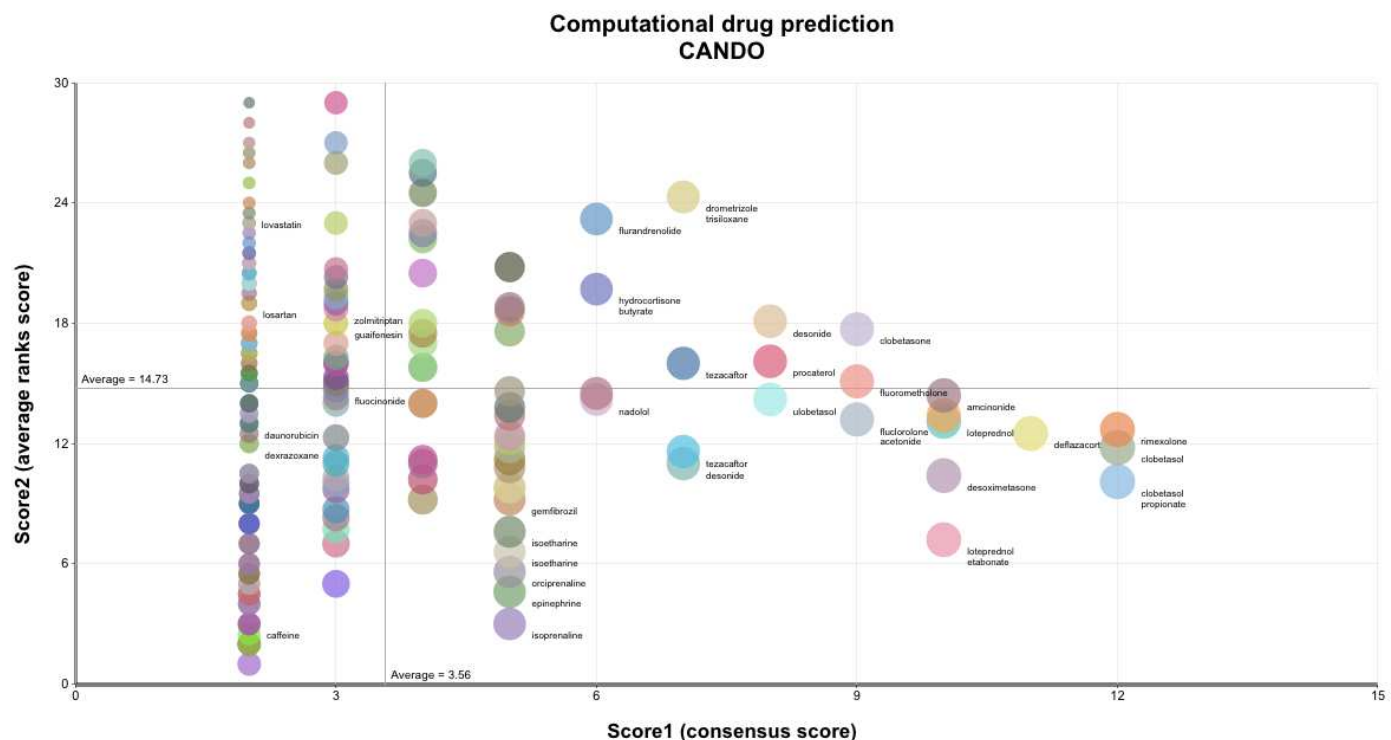
694

Figure/ Figure Legends.



**Figure 1 Putative drug candidates for treating COPD generated using the CANDO platform.** A subset of 130 proteins from the CANDO human protein library were identified from 233 differentially expressed proteins in the BALF. These 130 proteins were utilized to generate BALF-specific interaction signatures for 2,450 FDA-approved drugs via our in-house docking protocol BANDOCK (see methods). These drug-proteome interaction signatures were compared to those of 34 known drugs used to treat COPD to predict 189 most similar putative drug candidates. The 189 drugs are represented by colored circles, with the diameter of the circles decreasing with descending overall rank. Drug name labels are depicted for a selection of the 189 drugs shown by the colored circles. The horizontal axis plots the consensus score count or the number of times the particular drug is listed within the top30 most similar drugs to those known to treat COPD based on interaction signature similarity. The vertical axis plots the average of the cumulative ranks of the consensus scores for the putative drug. The overall rank

of a putative drug is determined by initially sorting the drug by the consensus score, as noted

above, and then additional sorting by the average rank. Many of the drug candidates were

corticosteroids not used to treat COPD; however other putative drugs included tezacaftor, a drug

to potentiate sodium channel activity in the treatment of cystic fibrosis; two additional drugs

predicted to treat COPD, gemfibrozil, and pioglitazone, are drugs currently used to treat

hyperlipidemia and diabetes, respectively.  This analysis indicates that the CANDO platform

applied to the BALF proteome is able to generate putative drug candidates for COPD treatment.

BANDOCK= bioanalytical docking
CANDO=computational analysis of novel drug opportunities
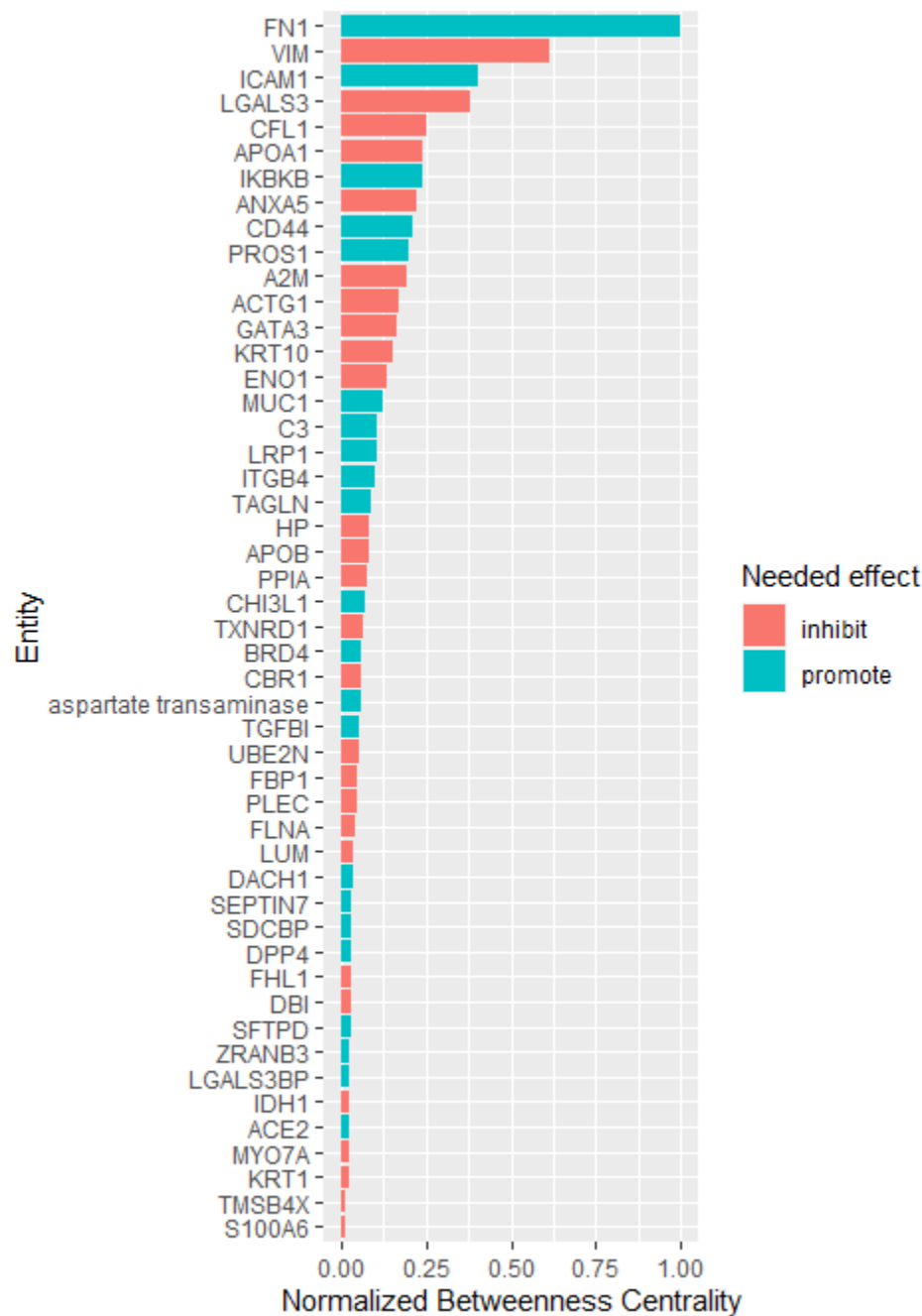COPD=chronic obstructive pulmonary disease.

Figure 2 **BALF network centrality nodes ranked by betweenness centrality.**
Betweenness centrality quantitatively describes how a node (in this case, a differentially

expressed protein in the BALF proteome) mediates the interaction between communities of

neighboring nodes in the network. Shown are 44 network entities with betweenness centrality

>0.01, normalized to the maximum betweenness centrality present in the network. The

betweenness centrality scores for all nodes were expressed as fractions of the maximum

betweenness centrality present in the network. The (red and blue) colors indicate the needed effect (inhibition/induction) to restore these entities from COPD levels to the normal levels in healthy control subjects. The four nodes with ≥25% of the maximum betweenness centrality (fibronectin) with normalized betweenness centrality values representing a greater than the linear increase from the next lower ranking node are fibronectin, vimentin, intercellular adhesion molecule1 (ICAM1), and galectin-3. These potential key signaling mediators had a betweenness centrality of at least 25% of the maximum. Topological analysis of the interaction network regulatory interactions documented in the literature suggests that these proteins were central mediators of COPD.[57-61]

Colors indicate the needed effect to restore these entities to the normal levels in healthy control subjects.
CANDO=computational analysis of novel drug opportunities
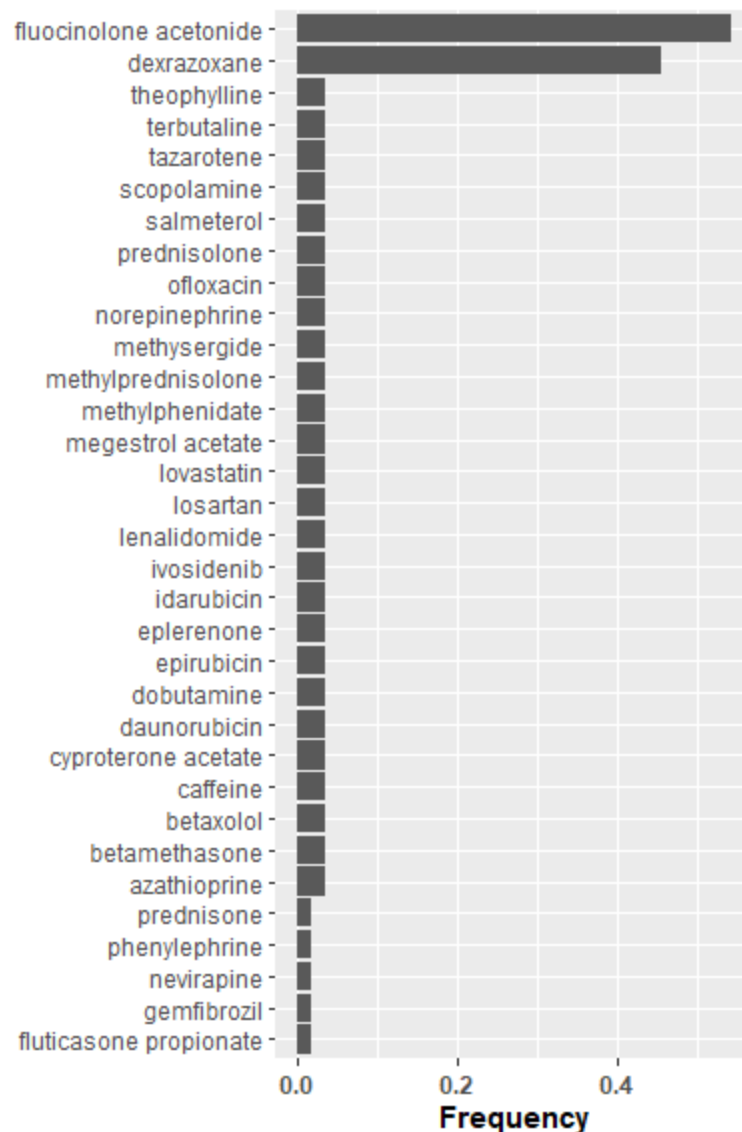COPD=chronic obstructive pulmonary disease

**Figure 3. Drug frequency amongst idealized drug combinations predicted to modulate central proteins in COPD.** The combinations of drugs initially identified by CANDO, which were predicted to activate or inhibit the four central nodes with the highest maximum betweenness centrality are listed (Figure 2). Drugs that lead to the promotion of central node proteins that were downregulated in the COPD cohort and inhibition of central node proteins (identified by network topological graph) that were overabundant in COPD. The

idealized drug vector constitutes interactions leading to desirable modulation of the central hub protein. Representation of individual drug frequency among the 57 significantly enriched two-drug combinations (idealized drug vectors) out of the 39 proteins represented in the Elsevier Knowledge Graph are listed in descending order. Fluocinolone acetonide and dexrazoxane appeared in 54% and 46% of all significantly enriched two-drug combinations respectively, far greater than other drugs appearing in these combinations. The combination of fluocinolone acetonide and dexrazoxane is the most enriched two-drug combination leading to an idealized drug vector that most likely reverses the protein levels of the four central nodes to levels found in healthy control subjects.

CANDO=computational analysis of novel drug opportunities
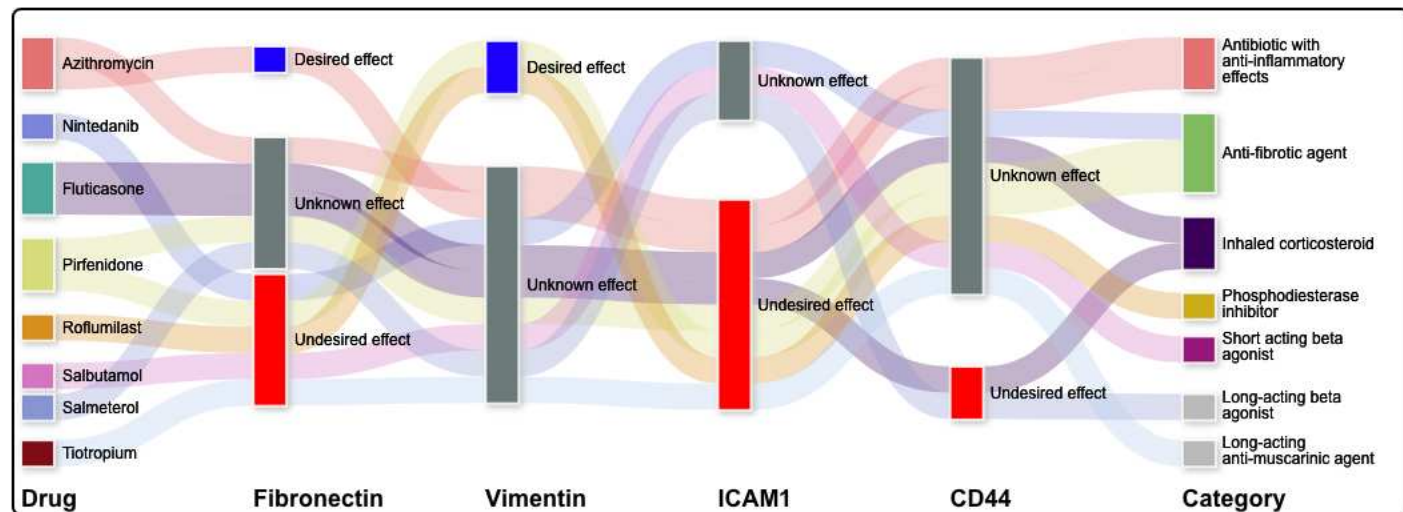COPD=chronic obstructive pulmonary disease

**Figure 4: Commonly used pulmonary drugs and their putative effects on four central node proteins in COPD.** A Sankey diagram categorizing the likely effects of putative drugs on four central node proteins in COPD is depicted. There are nine drugs on the left of the diagram used to treat different pulmonary diseases, with the corresponding drug classes displayed on the right side of the diagram. The effects of these drugs on four nodes (fibronectin, vimentin, intercellular adhesion molecule1 (ICAM1), and cd44) are detailed in the middle of the diagram, with broad lines connecting the proteins in the right to the putative effect (desired, unknown, undesired. While some of these have been documented to have the desired effect on fibronectin (promotion) or vimentin (inhibition), all have been reported to have the opposite effect on at least one of the most central proteins. This suggests using drugs commonly used to treat pulmonary disease, if repurposed for COPD, may have contrary effects on the mediators of the pathways involved in COPD, reinforcing the need to have a more nuanced approach to drug repurposing.

CANDO=computational analysis of novel drug opportunities

COPD=Chronic obstructive pulmonary disease

TABLES:

Table 1:

**Clinical parameters of never-smoking healthy subjects and ex-smokers with stable COPD**

**in BALF study**

| | Control subjects n=10 | COPD subjects (GOLD stage 2) n=10 | P-value |
|---|---|---|---|
| **Age** (years) | 63.4 ± 11.7 | 67.8 ± 8.5 | 0.15 |
| **Sex** | | | 0.31 |
| Male | 6 | 7 | |
| Female | 4 | 3 | |
| **Race** | | | 0.083 |
| Caucasian | 8 | 10 | |
| African-American | 2 | 0 | |
| **BMI** (kg/m2) | 28.5±4.2 | 32±9.7 | 0.32 |
| **Years patient quit smoking** | NA | 12.9 ± 4.4 | |
| **Tobacco smoking, Pack years** | NA | 56.6 ± 17.2 | <0.001 |
| **FEV₁** (% predicted) | 96.3 ± 14.8 | 65.9 ± 8.1 | <0.001 |
| **FVC** (% predicted) | 95.6 ± 13.4 | 87.6± 13.1 | 0.19 |
| **FEV₁/FVC** | 77.6± 3.8 | 57.8 ± 8.6 | <0.001 |

Table 1:

$FEV_1$: forced expiratory volume in 1 second.  FVC: forced vital capacity. Years quit:  Years subjects quit tobacco smoking. Pack-Years: The average number of packs of cigarettes smoked per week multiplied by the years the subject smoked cigarettes.