

1 **Whole exome sequencing analysis in severe chronic obstructive pulmonary disease**

2 Dandi Qiao^{1*}, Asher Ameli^{1,3}, Dmitry Prokopenko¹, Han Chen^{4,5}, Alvin T. Kho⁶, Margaret M. Parker¹,
3 Jarrett Morrow¹, Brian D. Hobbs^{1,2}, Yanhong Liu⁷, Terri H. Beaty⁸, James D. Crapo⁹, Kathleen C. Barnes¹⁰,
4 Deborah A. Nickerson¹¹, Michael Bamshad¹², Craig P Hersh^{1,2}, David A. Lomas¹³, Alvar Agusti¹⁴, Barry J.
5 Make⁹, Peter M.A. Calverley¹⁵, Claudio F. Donner¹⁶, Emiel F. Wouters¹⁷, Jørgen Vestbo¹⁸, Peter D. Paré¹⁹,
6 Robert D. Levy¹⁹, Stephen I. Rennard^{20, 21}, Ruth Tal-Singer²², Margaret R. Spitz⁷, Amitabh Sharma¹, Ingo
7 Ruczinski²³, Christoph Lange²⁴, Edwin K. Silverman^{1,2}, Michael H. Cho^{1,2*}; NHLBI Exome Sequencing
8 Project, University of Washington Center for Mendelian Genomics, Lung GO, COPD Gene Investigators

9

10 ¹ Channing Division of Network Medicine and ²Division of Pulmonary and Critical Care Medicine,
11 Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston,
12 Massachusetts 02115, United States of America

13 ³ Department of Physics, Northeastern University, Boston, Massachusetts 02115, United States of America

14 ⁴ Human Genetics Center, Department of Epidemiology, Human Genetics and Environmental Sciences,
15 School of Public Health, The University of Texas Health Science Center at Houston, Houston, Texas 77030,
16 United States of America

17 ⁵ Center for Precision Health, School of Public Health and School of Biomedical Informatics, The
18 University of Texas Health Science Center at Houston, Houston, Texas 77030, United States of America

19 ⁶ Boston Children's Hospital and Harvard Medical School, Boston, Massachusetts 02115, United States of
20 America

21 ⁷ Dan L. Duncan Comprehensive Cancer Center, Department of Medicine, Baylor College of Medicine,
22 Houston, Texas 77030, United States of America

- 1 ⁸ Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Johns Hopkins
2 University, Baltimore, Maryland 21205, United States of America
- 3 ⁹ National Jewish Health, Denver, Colorado 80206, United States of America
- 4 ¹⁰ Division of Allergy and Clinical Immunology, Department of Medicine, Johns Hopkins University,
5 Baltimore, Maryland 21224, United States of America
- 6 ¹¹ Department of Genome Sciences, University of Washington, Seattle, Washington 98195, United States of
7 America
- 8 ¹² Division of Genetic Medicine, Department of Pediatrics, University of Washington and Seattle Children's
9 Hospital, Seattle, Washington 98195, United States of America
- 10 ¹³ University College London, London WC1E 6BT, United Kingdom
- 11 ¹⁴ Respiratory Institute, Hospital Clinic, IDIBAPS, University of Barcelona, CIBERES, Barcelona 08007,
12 Spain
- 13 ¹⁵ University of Liverpool, Liverpool L69 3BX, United Kingdom
- 14 ¹⁶ Mondo Medico di I.F.I.M. srl, Multidisciplinary and Rehabilitation Outpatient Clinic, Borgomanero,
15 Novara 28021, Italy
- 16 ¹⁷ Department of Respiratory Medicine, Maastricht University Medical Center, 6202 AZ Maastricht,
17 Netherland
- 18 ¹⁸ University of Manchester, Manchester M13 9PL, United Kingdom
- 19 ¹⁹ Respiratory Division, Department of Medicine, University of British Columbia, Vancouver, British
20 Columbia V6T 1Z4, Canada
- 21 ²⁰ University of Nebraska Medical Center, Omaha, Nebraska 68198, United States of America
- 22 ²¹ AstraZeneca, Cambridge CB2 0RE, United Kingdom

1 ²² GSK Research and Development, King Of Prussia, Pennsylvania 19426, United States of America

2 ²³ Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, Maryland
3 21205, United States of America

4 ²⁴ Department of Biostatistics, Harvard School of Public Health, Boston, Massachusetts 02115, United
5 States of America

6

7

8

9 * **Corresponding author**

10 **Address: Channing Division of Network Medicine, 181 Longwood Avenue, Boston, MA 02115**

11 **Fax: 617 525 0958**

12 **Phone: 617 525 2113**

13 **Email: redaq@channing.harvard.edu**

14

15 **ABSTRACT**

16 Chronic obstructive pulmonary disease (COPD), one of the leading causes of death worldwide, is
17 substantially influenced by genetic factors. Alpha-1 antitrypsin deficiency demonstrates that rare coding
18 variants of large effect can influence COPD susceptibility. To identify additional rare coding variants in
19 patients with severe COPD, we conducted whole exome sequencing analysis in 2,543 subjects from two
20 family-based studies (Boston Early-Onset COPD Study and International COPD Genetics Network) and one
21 case-control study (COPDGene). Applying a gene-based segregation test in the family-based data, we
22 identified significant segregation of rare loss of function variants in *TBC1D10A* and *RFPL1* ($P < 2 \times 10^{-6}$),

1 but were unable to find similar variants in the case-control study. In single variant, gene-based, and
2 pathway association analyses, we were unable to find significant findings that replicated or were significant
3 in meta-analysis. However, we found that the top results in the two datasets were in proximity to each other
4 in the protein-protein interaction network ($p=0.014$), suggesting enrichment of these results for similar
5 biological processes. A network of these association results and their neighbors was significantly enriched
6 in the transforming growth factor beta-receptor binding and cilia-related pathways. Finally, in a more
7 detailed examination of candidate genes, we identified individuals with putative high-risk variants,
8 including patients harboring homozygous mutations in genes associated with cutis laxa and Niemann-Pick
9 Disease Type C. Our results likely reflect heterogeneity of genetic risk for COPD along with limitations of
10 statistical power and functional annotation, and highlight the potential of network analysis to gain insight
11 into genetic association studies.

12

13 **Introduction**

14 Chronic obstructive pulmonary disease (COPD) is a heterogeneous and complex disease, with a significant
15 genetic component to its susceptibility(1). Using genome-wide association analysis (GWAS), a number of
16 COPD susceptibility loci have been identified, including *FAM13A* (2), *HHIP* (3, 4), and
17 *CHRNA3/CHRNA5/IREB2* (5-7). Complementary studies have identified more than a hundred loci
18 associated with lung function, many of which likely also affect risk of COPD(8, 9). However, identified loci
19 only explain 5-10% of the heritability of COPD or quantitative measures of lung function traits (8, 10).
20 GWAS effectively tests common variants, but the well-known examples of alpha-1 antitrypsin
21 deficiency(11), cutis laxa (12-16), and the more recently described association between telomere-related
22 genes (17-19) indicate that, as has been shown for other diseases(20-24), rare coding variants also contribute
23 to COPD risk.

24 We previously analyzed exome sequencing data of 49 families with severe, early-onset COPD, and,
25 although we found several candidate genes, none showed convincing evidence of replication (25, 26). We

1 further showed using simulations that genetic heterogeneity may be a major contributor to this failure to
2 replicate(26). In this study, we applied additional sequencing and analytic strategies to increase the sample
3 size and the power of the analysis. We applied a recently developed family-based method, GENE-based
4 SEgregation (GESE) (25), to a larger family-based dataset enriched for severe COPD, and also performed
5 single-variant, as well as set-based tests using SKAT-O for both genes and pathways in the family-based
6 and in an additional case-control study. We tested for enrichment of our results in gene expression and
7 monogenic models of disease, and examined the overlap between case-control and family-based results
8 using network analysis. Finally, we investigated a set of candidate genes identified in previous genetic
9 studies, including Mendelian syndromes, for potentially deleterious rare variants.

10 **Results**

11 **Gene-based segregation test on the ICGN and Boston EOCOPD pedigrees**

12 Baseline characteristics of the studied subjects are shown in **Table 1**. Additional information on the
13 probands can be found in **Supplementary Table S1**. To identify causal variants in our exome sequencing
14 data with the characteristics of Mendelian variants for COPD (e.g. alpha-1 antitrypsin deficiency), we
15 applied our recently described gene-based segregation (GESE) test (25) to the family-based data. We
16 focused on ultra-rare ($MAF < 0.1\%$) predicted loss-of-function variants. Two genes were significant after
17 Bonferroni correction for the total of 18268 genes: *RFPL1* ($p = 1.60e-06$) and *TBC1D10A* ($p=1.10e-06$).
18 *RFPL1* segregated in 4 families, including two singleton families and two families with affected sibling
19 pairs of severe COPD. *TBC1D10A* segregated in a parent-offspring pair and a singleton family. *TBC1D10A*
20 is intolerant to loss-of-function variants (ExAC intolerance probability = 0.98 (27)). The top 10 genes from
21 this analysis are shown in **Table 2**. All 10 of these genes are expressed in the adult lung (see Methods,
22 enrichment $p = 0.17$), and the expression of 5 out of 9 of those genes was associated with forced expiratory
23 volume in 1 second (FEV_1) % predicted, a measure of COPD severity, in our lung tissue data (enrichment p
24 = 0.024). We further sought supportive evidence for association of these genes in the COPDGene case-
25 control dataset. However, no subjects harbored loss of function variants in these genes. We additionally

1 tested for evidence of higher burden of rare (MAF < 0.1%), non-synonymous variants in the cases, and did
2 not find convincing evidence of association (*RFPL1* p=0.576; *TBC1D10A* p=0.081).

3 **Single-variant association analysis in the case-control and family data**

4 Next, we performed single variant association analysis. We tested both rare coding variants (moderate effect
5 by SNPEff, and MAF < 5%) as well as all variants. We found no significant results (**Supplementary Table**
6 **S2 and S3**) in either our primary analysis using COPDGene as the discovery cohort (using a Bonferroni
7 significance level of 1.32e-06 for non-synonymous variants with MAF < 5%, and 5.07e-07 for all variants),
8 or using the family-based data (3.55e-07 for non-synonymous variants with MAF < 5% and significance
9 level 1.86e-07 for all variants). However, top variants in the case-control analysis included rs8040868
10 (MAF = 0.41) and rs1051730 (MAF = 0.35) in *CHRNA3*(28) with p = 5.05e-05 and 7.39e-05 respectively,
11 which reside at a previously described GWAS locus (**Supplementary Table S2**). Top variants in the
12 family-based analysis included rs2232710 (MAF = 0.012; p = 4.05e-05) in *SERPINA10* (in high D' with the
13 alpha-1 Z allele, which causes alpha-1 antitrypsin deficiency – note that severe alpha-1 antitrypsin
14 deficiency, including ZZ homozygosity, was an exclusion criteria for these studies) and rs10507051 (MAF
15 = 0.063; p = 1.28e-04) in *VEZT* (**Supplementary Table S3**), near a locus associated with COPD in a recent
16 GWAS of lung function(8). We also considered whether any variants were significant in meta-analysis by
17 combining results from the two studies (case-control status in the COPDGene data, and lung function in the
18 BEOCPD-ICGN data) using the Stouffer method. Meta-analysis did not identify significant variants
19 among the rare coding variants (significance level 2.82e-06) or among all variants (significance level 7.08e-
20 07) (**Supplementary Table S4, Supplementary Table S5**); top results overall included variants in
21 *CHRNA3* and *SERPINA10*.

22 **Gene-and pathway-based analyses in case-control and family data**

23 Next, we performed gene-based analyses. In the analysis using SKAT-O and predicted deleterious variants
24 with a MAF < 1%, we found no significant genes in the COPDGene data. The top 10 genes are shown in
25 **Supplementary Table S6**. We found no significant enrichment of genes expressed in lung (enrichment p =

1 0.97) among the top 10 genes. However, four genes have expression associated with FEV₁ % predicted (p =
2 0.087), including the top two genes *VNN1* and *PLA1A*. *EGFL8*, the 3rd ranked gene in the list, is located near
3 the *AGER* locus which was previously associated with risk of COPD (8, 9). In the pathway analysis using
4 the KEGG(29) database, we found one significant pathway using the burden test, the Jak-STAT signaling
5 pathway (p = 6.78e-05)(30). However, association with this pathway was not replicated in the family-based
6 analysis (p = 0.54) using the burden test. Top results from family-based analyses can be found in
7 **Supplementary Table S7**. We also conducted a meta-analysis of the COPDGene dataset and the
8 BEOCOPD-ICGN dataset; however, no gene achieved significance (**Supplementary Table S8**).

9 **Enrichment and network-based approach to overlap**

10 Given our lack of significant associations using standard association tests, we sought evidence that our top
11 case-control and family-based results were enriched for associations with COPD. We tested for enrichment
12 of overlap of genes yielding nominal significance (i.e. p < 0.01) between the case-control and the family-
13 based association results using a standard hypergeometric approach. The enrichment p-value was 1, which
14 was consistent with our lack of overlap and meta-analysis findings.

15 While we did not observe overlap between the top results in the case-control analysis and the family-based
16 analysis using a simple hypergeometric test, we were interested in studying common biological pathways
17 shared by the two sets of top genes. Recently, network-based methods have demonstrated the ability to
18 identify related diseases in the protein-protein interactome (31). We hypothesized that application of this
19 method to two independent association results for COPD would a) identify whether there were overlapping
20 association signals, despite the lack of replication; and b) identify genes or pathways of highest priority. We
21 computed the network-based separation (31) defined as the normalized average shortest path between
22 members from the two modules to see whether top genes from the case-control and family-based analysis
23 were close to each other in the protein-protein interaction (PPI) network. For the analysis of rare and
24 deleterious variants, we found genes with p < 0.01 from the case-control analysis and the family-based
25 analysis had significantly overlapping neighborhoods with negative separation score (score = -2.29, p =
26 0.014). To explore the neighborhood of these genes and the common pathways that connect the top genes,

1 we added the first neighbors of the top genes in the PPI. These genes (top genes, along with all of their first-
2 degree neighbors - a total of 522 genes) formed a largest connected component (LCC) of 513 genes, which
3 means almost all the top genes and their first neighbors were connected. **Figure 1** shows the network
4 module containing the largest connected component formed by the top genes from the two analyses and
5 their first neighbors. There were 19 genes with $p < 0.01$ in the family-based data, which had 274 first-degree
6 neighbors in the LCC network; there were 14 genes with $p < 0.01$ in the COPDGene data, which had 216
7 first-degree neighbors. Between the two groups of 274 and 216 first-degree neighbors, 10 overlapped, thus
8 these genes together formed a network module of 513 genes. 14 genes at loci previously associated with
9 COPD or lung function (out of 329 genes in the curated set, see Methods) were in this set (enrichment $p =$
10 0.065, Supplementary Material). Additional examination of these genes in murine models showed that the
11 513 genes were significantly enriched for genes associated with the respiratory system (enrichment $p =$
12 0.045) and were enriched for genes involved in normal murine lung development in three common inbred
13 strains of mice (enrichment $p = 1.35e-02, 1.96e-03$ and $2.40e-03$, respectively) (See Methods). From this
14 result, we postulate that there is a large disease network module exists likely including a subset of these 513
15 genes for severe COPD, and only part of this disease module was observed using either analysis alone due
16 to limited power. However, since the two sections of disease module share similar function and pathways,
17 they were significantly close to each other in the PPI network.

18 To further explore the functions of this network, we also looked at the pathways enriched for these 513
19 genes using ToppFun in the ToppGene Suite(32), and found a large number of Gene-ontology pathways
20 were significantly enriched. To examine more specific pathways, we examined GO pathways with fewer
21 than 100 genes in total. The top two pathways meeting these criteria were “GO:0005160: transforming
22 growth factor beta (TGFB) receptor binding” and “GO:0030991: intraciliary transport particle A”. Fourteen
23 out of fifty three genes in the TGFB receptor binding pathway were present in the network module. Multiple
24 lines of evidence, including genetic association (which has identified *TGFB2*) and other genomic and
25 mechanistic studies have implicated this pathway in risk to COPD (33-35). Twelve of these fourteen genes
26 are expressed in human lung tissue (two genes have missing data). The right panel in **Figure 1** shows the

1 small network formed by genes in this pathway and *ACVR2B*, which was the top-ranked gene from this set
2 of genes in the association analysis and was also the second largest hub in the network. For the “intraciliary
3 transport particle A pathway” (36), seven out of eight genes were in the network, which are shown in the
4 left panel in **Figure 1**. Six out of seven genes were expressed in human lung tissue (one gene had missing
5 data). *TULP3* was in the top ranked genes from the family-based analysis and was the largest hub in the
6 module. *TULP3* is a known target of the Hedgehog pathway. Notably, genome-wide association studies and
7 follow-up functional studies have identified an important role for *HHIP* in the development of COPD(7);
8 *TULP3* has been shown to change expression after *HHIP* silencing (37). Also, *WDR35* and *IFT140* were
9 associated with respiratory system abnormalities in mouse models (*WDR35* leads to lung hypoplasia and
10 mutations in *IFT140* produces severely misshapen lungs). Additional top results from this GO analysis can
11 be found in **Supplementary Table S9**. Thus, our network results highlight *ACVR2B* and *TULP3*, which
12 may be prioritized for further examination of functional rare variants.

13 **Evidence of association for candidate genes**

14 A substantial proportion of rare variants identified for complex disease are located at loci that also harbor
15 common risk variants (38, 39). In addition, several Mendelian syndromes have COPD, emphysema, or
16 obstructive lung disease as a manifestation of disease. Therefore, in addition to looking at exome-wide
17 results, we examined a list of the 329 curated genes (see Methods, **Supplementary Table S10**) (1, 7-9, 12-
18 19, 40-45). This included regions identified from 105 SNPs from GWAS analyses(8, 9) and 29 Mendelian
19 genes with manifestations that include COPD or emphysema in their resulting syndromes (**Supplementary**
20 **Material**). We examined functional and rare variants with MAF < 5%, and found multiple genes to be
21 nominally associated with COPD status or FEV₁ % predicted value, including *CHRNA5*, *AGER*, and
22 *CYP2A6* (**Supplementary Table S11**). To identify whether there was any independent evidence of rare
23 variant effects at these loci in the COPDGene cohort, we conditioned on the risk allele for the 104 SNPs
24 identified by GWAS. Several genes were still nominally significant after conditioning on the GWAS SNPs,
25 including *CYP2A6* (full results shown in **Table 3**); whether these rare variants have independent effects on
26 COPD susceptibility at these loci will likely need to be addressed by additional, larger studies.

1 We also looked closely at the 29 genes causing Mendelian syndromes including emphysema or obstructive
2 lung disease as part of their syndrome. To determine whether there was enrichment in these genes in our
3 dataset, we performed a burden test including only variants with MAF < 0.1% in ExAC and predicted
4 deleterious by FATHMM, SIFT, and CADD (>15). We found that the burden-based tests gave a $p = 0.80$ in
5 the COPDGene case-control study, and a $p = 0.018$ for the family-based EOCOPD and ICGN data. Thus,
6 we observed some significant accumulation of deleterious variants in these genes in the family-based data,
7 suggesting that ultra-rare variants in these Mendelian genes contributing to lung function may be related to
8 severe COPD risk in our family-based datasets.

9 To examine these variants individually, we intersected variants in these genes with Clinvar, using an
10 annotation of significance level 4 (likely pathogenic) and above, and additionally included variants in
11 published reports associated with respiratory disease in *TERT* (17-19, 41-43). We found 47 of these variants
12 in our datasets. These variants are listed in **Supplementary Table S12** along with their counts among cases
13 and controls separately. Given the strong evidence of pathogenicity for variants in *SERPINA1* and
14 telomere-related genes, these findings are shown in **Table 4**. We also assessed the carriers of these rare
15 variants using a recessive model of inheritance, and those variants with homozygous genotypes present in
16 any dataset are listed in **Table 5**. Among our findings for Mendelian genes were two previously identified
17 cases from COPDGene with heterozygous *TERT* variants (19), and evidence for an increased burden (cases
18 > controls) for the *SERPINA1* Z and PI P (Lowell) (rs121912714) (**Table 4**). For recessive variants, we
19 identified rs140130028, a splice-donor variant in *NPC2*, which is a gene for Niemann-Pick disease type C2,
20 a disease previously associated with emphysema (46) (**Table 5**). One pair of sibs with severe COPD in the
21 ICGN study was homozygous for this variant; two of their half-siblings carried one copy, one with severe
22 COPD. None of these subjects had known Neimann-Pick disease. Also, variant rs61748181 in *TERT* was
23 present as homozygous in 7 unrelated cases in the datasets (**Table 5**). While this association did not reach
24 candidate-wide significance ($p = 0.167$), this variant was experimentally demonstrated to induce telomere
25 dysfunction (47) and predicted to be disease causing by Mutation Taster (48). For variants not annotated by
26 Clinvar or annotated with a significance level of 3 (uncertain significance) or below, we filtered based on

1 MAF < 0.1% in ExAC v0.3 non-Finnish Europeans and predicted deleterious effects by FATHMM, SIFT,
2 or CADD (>15). There were in total 346 such variants in our datasets. One of these variants occurred in
3 homozygous form in a proband with severe COPD in the BEOCOPD study. This variant is an ultra-rare
4 splice-acceptor variant in *ATP6V0A2* (novel in ExAC database) (**Table 4**), a Mendelian gene for cutis laxa.
5 A chest CT scan of this subject showed severe emphysema, however, no phenotypic information related to
6 dermatological characteristics was available. In addition, 66 variants were predicted to be deleterious by all
7 three annotations: FATHMM, SIFT and CADD (>15), and had supportive evidence in our datasets (with
8 greater counts in cases than in controls, **Supplementary Table S13**). Multiple variants have supportive
9 evidence in both case-control and family-based datasets. For example, rs141310608 in *EFEMP2* is present
10 in 2 cases in COPDGene study and 2 cases in BEOCOPD-ICGN study, while none in controls. Also, there
11 are multiple ultra-rare variants in *COL3A1* are carried by cases and none by controls. We have also listed the
12 variants that are predicted to be deleterious by all annotations, but are present in more controls than cases;
13 these variants are less likely to be high penetrance COPD susceptibility variants (**Supplementary Table**
14 **S14**). We also applied a more liberal filtering criteria (MAF < 0.05, CADD > 10 or predicted to be
15 deleterious by SIFT and FATHMM) for *TERT*, *RTEL1*, *CFTR* and *SERPINA1*. Detailed information about
16 these genes can be found in Supplementary **Tables S15, S16** and **S17** respectively.

17 **Discussion**

18 COPD is a common and heterogeneous disease; under the common-disease-common-variants hypothesis,
19 we expect that multiple common variants should contribute to a large proportion of COPD risk. However,
20 even though a number of COPD GWAS loci have been discovered through large-scale collaborative efforts,
21 most of the estimated heritability remains unexplained. Examples such as alpha-1 antitrypsin deficiency,
22 cutis laxa, and more recently, telomeropathies are associated with COPD and emphysema (17-19). These
23 results motivated us to search the entire exome for large effect variants that could represent a Mendelian
24 subtype of COPD, in the hope of finding new treatment strategies for a subset of the patients. In this study,
25 we examined multiple cohorts representing the largest exome sequencing study of COPD to date ascertained
26 under an extreme phenotype approach (where samples were enriched for severe COPD and normal controls

1 heavily exposed to smoking but with normal pulmonary function), to screen through the entire exome to
2 identify rare coding variants controlling risk to COPD. Results failed to identify new genes, pathways, or
3 variants consistently significant across all of our analyses, suggesting that single-variant or single-gene
4 effects of a contribution as large as alpha-1 antitrypsin deficiency are unlikely to exist (26). Yet, a network-
5 based analysis identified a significant relationship between the two modules formed by the top results of the
6 two analyses. These two sets of top genes, along with their first neighbors in the protein-protein interaction
7 network, form a well-connected network component. This largest connected component was significantly
8 enriched in genes involved in fetal lung development in mouse models (49). Additionally, this module sheds
9 light on related functions or pathways where such rare variants may be contributing to risk to COPD. For
10 example, multiple studies have suggested the transforming growth factor beta pathway is associated with
11 COPD(50), and the *TGFB2* locus was associated with COPD in genome-wide association studies (GWAS)
12 (51). Our study identified *ACVR2B* as a potential candidate; of interest, *ACVR1B*, an activin receptor which
13 interacts with *ACVR2B* (52) was identified in a network-informed genetic association study of COPD (53)
14 and in a integrative analysis of emphysema distribution (54).

15 Our finding lends further support to the transforming growth factor beta pathway and also suggests that rare
16 variants related to *ACVR2B* may contribute to COPD risk. Similarly, the identification of *TULP3* lends
17 further support to the identification of *HHIP* as a causal gene at this GWAS locus and the importance of the
18 hedgehog pathway in the development of COPD. The identification of cilia-related pathways is intriguing
19 given the importance of cilia to lung function (55), including reports from a smaller exome study of resistant
20 smokers (56) and reports of shortened cilia in smokers and in COPD patients (36).

21 Finally, we identified subjects carrying homozygous genotypes of rare and deleterious variants in Mendelian
22 genes for cutis laxa and Niemann-Pick disease, which are themselves intriguing candidates for causing
23 severe COPD. These findings illustrate the potential relevance of using filtering-based technique for
24 identifying syndromic forms of COPD. While we do not have enough power to individually test these or
25 other individual rare variants here, our results may provide support for future studies in these recognized
26 candidate genes.

1 COPD is known to be a highly heterogeneous disease, with varying contributions of emphysema and small
2 airway disease. We did not examine specific subsets of COPD, as detailed phenotyping was not available in
3 all cohorts. Multiple analysis methods are available for rare variant analysis ((57)), and the optimal methods
4 are still not clear. Our sequencing of a large number of affected individuals in families was appropriate for
5 methods such as GESE, which leverages a large reference dataset (ExAC); an alternative approach using
6 association would require large scale exome harmonization of controls with normal lung function,
7 preferably with heavy cigarette smoke exposure. Our results highlight the importance of integration with
8 other types of data (e.g., gene expression, protein-protein interaction) to better understand the results from
9 one data type. However, our analysis does not attempt to identify the confidence of individual genes in this
10 network; we cannot rule out the possibility that this network includes many genes that are false positives,
11 and our pathway analysis should be considered descriptive and exploratory. Additional investigation,
12 including genetic studies, integration of multi-omics data, and careful functional studies will be needed to
13 further infer biological mechanisms and potential disease causality for our identified genes.

14 In summary, in an exome sequencing study of COPD, we were unable to identify exome-wide significant
15 associations, but through network analysis we identified candidate genes in related pathways and a disease
16 module driven by rare variants. Our study is consistent with a potential contribution of multiple,
17 heterogeneous rare variants in COPD, and demonstrates the insight that network-based methods can offer.

18 **Materials and Methods**

19 **The COPDGene study**

20 The COPDGene study is a multi-center epidemiologic and genetic study of 10,192 current or ex-smokers,
21 which has been previously described (58). COPDGene subjects were sequenced in two sets. The first set
22 sequenced as part of the NHLBI Exome Sequencing Project (ESP; named COPDGene ESP) included severe
23 COPD cases with GOLD (Global Initiative for Chronic Obstructive Lung Disease) Grade 3 or 4 (post-
24 bronchodilator FEV1 < 50% predicted and FEV1/FVC < 0.70), and age < 65 years old, with substantial
25 emphysema (> 15% at -950 HU) by quantitative chest CT scan. Controls were selected to be resistant

1 smokers with frequency-matched pack-years of cigarette smoking, normal lung function ($FEV_1 > 80\%$
2 predicted and $FEV_1/FVC > 70\%$), age > 65 years old, and no significant emphysema ($< 5\%$ at -950 HU).
3 The second set sequenced at Baylor (named COPDGene Baylor) included severe COPD cases (GOLD
4 Grade 3 or 4) with no age requirement. Controls were selected to be resistant smokers with normal lung
5 function with age > 55 .

6 **The Boston Early-Onset COPD study and the International COPD Genetics Network study**

7 The family-based data contained samples selected from the Boston Early-Onset COPD study (BEOCOPD)
8 (59) and the International COPD Genetics Network (ICGN) study(45). Probands from BEOCOPD were
9 selected to be physician-diagnosed COPD cases with $FEV_1 \leq 40\%$ predicted, age ≤ 53 . All first-degree
10 relatives, older second-degree relatives and additional affected family members were enrolled. Probands in
11 the ICGN study were subjects with known COPD and were required to have $FEV_1 < 60\%$ predicted,
12 $FEV_1/FVC < 90\%$ predicted at age 45-65, pack-years ≥ 5 , and have at least one eligible sibling. An initial
13 set of 49 pedigrees selected from the Boston Early-Onset COPD Study were described and analyzed
14 previously (26). To this sample we added 147 families from BEOCOPD and 462 families from the ICGN
15 study. The COPDGene, BEOCOPD, and ICGN studies all excluded subjects with severe alpha-1 antitrypsin
16 deficiency.

17 **Exome sequencing**

18 We sequenced all subjects using Nimblegen capture and Illumina platforms. The COPDGene ESP,
19 BEOCOPD, and ICGN subjects were all sequenced at the University of Washington, using Nimblegen V2
20 exome capture; COPDGene Baylor samples used VChrome capture. Alignment, variant calling, and quality
21 control were performed using bwa, GATK, and in-house pipelines respectively. As COPDGene ESP and
22 COPDGene Baylor used slightly different capture platforms, calling was performed on these datasets
23 separately. All BEOCOPD and ICGN subjects were called together (joint calling) and went through the
24 same quality control steps together to provide the final family-based data (named BEOCOPD-ICGN) for

1 analysis. Baseline characteristics of the subjects in each of the cleaned datasets are shown in **Table 1** and
2 our overall study design is shown in **Figure 2**. More details can be found in the **Supplementary Material**.

3 **Analysis strategy**

4 *Loss of function variants using the Gene-based Segregation test (GESE)*

5 We first performed the gene-based segregation test (GESE)(25) on loss of function variants (defined by
6 SnpEff (60)) with MAF < 0.1% in the family-based BEOCPD-ICGN data using COPD affection status as
7 the outcome. We included only the most severe COPD subjects (GOLD spirometry grade 3 or 4) and
8 resistant smoking control subjects (normal spirometry, age > 40 yr, with at least 5 pack-years of cigarette
9 smoking). This analysis took advantage of the unique properties of a family-based strategy, including
10 having multiple copies of rare variants, and assumes a Mendelian model with a few rare variants with very
11 large effects. We sought supportive evidence for identified causal genes in COPDGene dataset by
12 attempting to identify similarly deleterious variants.

13 *Association analyses*

14 Second, we performed single variant, gene-based, and pathway-based association analyses. For all
15 association analyses, we used Bonferroni correction based on the number of genes, pathways, or variants
16 tested. For the COPDGene case-control data, COPD affection status was used as the outcome, which was
17 adjusted for pack-years, gender, age, and ancestry-based principal components (PCs) in the COPDGene
18 Baylor data, and the top PCs alone in the COPDGene ESP data due to the selection criteria, and as
19 performed previously (Supplementary Material). For the family-based data, due to the low number of
20 controls with normal lung function, but a wider range of FEV₁ available through family members, we
21 analyzed FEV₁ (forced expiratory volume in one second), a lung function measure highly correlated with
22 COPD (9) instead of COPD affection status itself. The outcome in the family-based association tests was the
23 rank of the residuals from regressing raw post-bronchodilator FEV₁ value on height, pack-years, sex, age,
24 top 5 genetic ancestry PCs and batch indicator variable.

1

2 ***Single-variant association analysis***

3 For single-variant analyses, we applied the Stouffer method to meta-analyze the results from the hybrid
4 method in SKATBinary_Single function (SKAT package) in the COPDGene case-control data, since the
5 two cohorts selected from the COPDGene study were sequenced and called separately. The hybrid method
6 in SKATBinary_Single function selects the most appropriate approach to compute p-values for each variant.
7 For single variant analysis in family-based data, we applied the variant-based generalized linear mixed
8 model association test (GMMAT(61)). In addition to using COPDGene as discovery and BEOCPD-ICGN
9 as replication, we also examined using BEOCPD-ICGN and both datasets as discovery by meta-analyzing
10 the results from the COPDGene case-control data and the BEOCPD-ICGN data using the Stouffer method.
11 For single variant analyses, we tested all variants, and also the subset with moderate effect with $MAF < 5\%$.

12 ***Gene- and pathway-based association analysis***

13 For both the gene-based and pathway-based analyses, we applied SKAT-O tests. In the COPDGene case-
14 control datasets, we applied the hybrid method in the SKATBinary function, implemented in the SKAT
15 package to each of the datasets, and meta-analyzed the two datasets using Fisher's method. For the
16 BEOCPD-ICGN family-based data, we applied MONSTER(62), which is a generalized version of SKAT-
17 O for family-based studies. We also meta-analyzed all results (case-control and family-based results) using
18 Fisher's method. Our primary gene- and pathway-based association analyses focused on deleterious variants
19 defined using FATHMM (57, 63) with $MAF < 1\%$ in the association analysis. In one study of amyotrophic
20 lateral sclerosis (ALS), FATHMM was found to give the best power to identify known causal genes for
21 ALS in gene-based association tests(57). Our secondary analyses included association testing on functional
22 variants with moderate effects (defined by SnpEff (60)) with $MAF < 5\%$. This is a less stringent filtering
23 criterion on the variants to prevent missing signals in this set of variants. Pathways were defined using
24 KEGG pathways(29) and the c2 collection of curated gene sets from the Molecular Signatures Database
25 (MsIGdb) in GSEA(64).

1 *Identification of enrichment in gene expression*

2 To help determine whether the identified genes were relevant for our phenotypes, we used publicly available
3 FPKM (per kilobase of gene model per million mapped reads) results from gene expression data from the
4 Lung Genomics Research Consortium (LGRC, <http://www.lung-genomics.org>) to identify whether any gene
5 was expressed in the lung (using a cutoff of 0.5(26)). We also used the results of differential expression for
6 lung function and COPD case-control status in an independent set of lung tissue from severe COPD subjects
7 and controls (65). In addition, enrichment for genes associated with respiratory system in mouse was carried
8 out using a curated set of genes associated with respiratory phenotype in the Mouse Genome Database
9 (<http://www.informatics.jax.org/marker>) (66). Gene expression information in human and normal murine
10 lung development for three common inbred strains of mice were obtained from the GEO dataset (GSE14334
11 and GSE74243), and genes involved in fetal lung development were obtained using methods described in
12 (49).

13 *Network-based analysis*

14 Finally, we applied the network-based separation measure defined in (31) to examine how closely connected
15 the top genes from the two independent analyses are in the protein-protein interaction (PPI) network. This
16 measure has been shown to predict pathobiological similarity of two sets of disease genes(31). In our
17 application here, since the two outcomes analyzed for the COPDGene and BEOCOPD-ICGN dataset are
18 highly correlated, genes that are causal for these outcomes should have much shorter network-based
19 distance. Therefore, a significant result tells us that at least a subset of the top genes from the two analyses
20 are topologically overlapping and exert some effect on risk of COPD.

21 *Examination of previously identified genetic associations with COPD*

22 To examine loci previously described to be associated with risk of COPD or lung function itself in GWAS
23 or harboring Mendelian variants related to COPD, we curated a set of 329 genes for closer examination
24 (**Supplementary Table S12**)(8, 9). At COPD GWAS loci, we identified all variants in a European reference
25 population with an $r^2 > 0.8$ with the lead variant, and then expanded these borders by 100kb. For Mendelian

1 syndromes, we included connective tissue disorders such as cutis laxa (12-16), as well as telomere-related
2 genes including *TERT*, *TERC*, *RTEL1*, and *NAF1* (17, 18, 41-43). We looked for supportive evidence of
3 association for these genes using several methods. First, we examined the association results in both primary
4 and secondary analyses as described above. Since 104 of the previously described lead SNPs based on
5 GWAS of lung function or COPD were also available for the COPDGene subjects, we additionally
6 performed conditional analyses for these genes by conditioning on the GWAS SNPs in proximity in an
7 attempt to identify independent rare variants contributing to COPD susceptibility. For both the marginal
8 association analyses and conditional analyses, COPD affection status was the outcome in the COPDGene
9 case-control analyses and FEV₁ was the outcome in the family-based analyses. Finally, we examined
10 Mendelian genes for evidence of pathogenic variants using Clinvar and other public annotation resources.

11

12 **Acknowledgements**

13 This work was supported by NHLBI R01 HL089856 (EKS), R01 HL089897 (JDC), R01 HL113264 (MHC
14 and EKS), P01 HL105339 (EKS) and PO1 114501 (EKS), K01 HL129039 (DQ), K07 CA181480 (YL). The
15 COPDGene project (NCT00608764) is also supported by the COPD Foundation through contributions made
16 to an Industry Advisory Board comprised of AstraZeneca, Boehringer Ingelheim, GlaxoSmithKline,
17 Novartis, Pfizer, Siemens and Sunovion. The ICGN study was funded by GlaxoSmithKline. Sequencing for
18 the Boston Early-Onset COPD Study was provided by the University of Washington Center for Mendelian
19 Genomics (UW CMG) and was funded by the National Human Genome Research Institute and the National
20 Heart, Lung and Blood Institute grant 1U54HG006493.

21 Authors would like acknowledge Victor M Pinto-Plata (Baystate Medical Center, Springfield, MA),
22 Nathaniel Marchetti (Temple University School of Medicine, Philadelphia, PA), Raphael Bueno (Brigham
23 and Women's Hospital and Harvard Medical School, Boston, MA), Bartolome R Celli (Brigham and
24 Women's Hospital and Harvard Medical School, Boston, MA), Gerald J Criner (Temple University School
25 of Medicine, Philadelphia, PA), and Dawn Demeo (Brigham and Women's Hospital and Harvard Medical

1 School, Boston, MA) for providing the lung tissue samples and their support of the project. The content is
2 solely the responsibility of the authors and does not necessarily represent the official views of the National
3 Heart, Lung, and Blood Institute or the National Institutes of Health,

4

5 **Conflict of Interest Statement**

6 Dr. Silverman has received honoraria and consulting fees from Merck, grant support and consulting fees
7 from GlaxoSmithKline, and honoraria from Novartis. Dr. Hersh has been a consultant for CSL Behring and
8 Mylan. Dr. Cho has received grant support from GSK.

9

10

11 **References**

- 12 1 Hersh, C.P., Demeo, D.L. and Silverman, E.K. (2005) Silverman EK, S.S., Lomas DA,
13 Weiss ST, editors (ed.), In *Respiratory genetics*. Hodder Arnold, New York, in press.
- 14 2 Cho, M.H., Boutaoui, N., Klanderman, B.J., Sylvia, J.S., Ziniti, J.P., Hersh, C.P., DeMeo, D.L.,
15 Hunninghake, G.M., Litonjua, A.A., Sparrow, D. *et al.* (2010) Variants in FAM13A are
16 associated with chronic obstructive pulmonary disease. *Nat. Genet.*, **42**, 200-202.
- 17 3 Zhou, X., Baron, R.M., Hardin, M., Cho, M.H., Zielinski, J., Hawrylkiewicz, I., Sliwinski, P.,
18 Hersh, C.P., Mancini, J.D., Lu, K. *et al.* (2012) Identification of a chronic obstructive pulmonary
19 disease genetic determinant that regulates HHIP. *Hum. Mol. Genet.*, **21**, 1325-1335.
- 20 4 Wilk, J.B., Chen, T.H., Gottlieb, D.J., Walter, R.E., Nagle, M.W., Brandler, B.J., Myers, R.H.,
21 Borecki, I.B., Silverman, E.K., Weiss, S.T. *et al.* (2009) A genome-wide association study of
22 pulmonary function measures in the Framingham Heart Study. *PLoS Genet.*, **5**, e1000429.
- 23 5 Hardin, M., Zielinski, J., Wan, E.S., Hersh, C.P., Castaldi, P.J., Schwinder, E.,
24 Hawrylkiewicz, I., Sliwinski, P., Cho, M.H. and Silverman, E.K. (2012) CHRNA3/5, IREB2, and
25 ADCY2 are associated with severe chronic obstructive pulmonary disease in Poland. *Am. J.*
26 *Resp. Cell. Mol.*, **47**, 203-208.
- 27 6 DeMeo, D.L., Mariani, T., Bhattacharya, S., Srisuma, S., Lange, C., Litonjua, A., Bueno, R.,
28 Pillai, S.G., Lomas, D.A., Sparrow, D. *et al.* (2009) Integration of genomic and genetic
29 approaches implicates IREB2 as a COPD susceptibility gene. *Am. J. Hum. Genet.*, **85**, 493-502.
- 30 7 Pillai, S.G., Ge, D., Zhu, G., Kong, X., Shianna, K.V., Need, A.C., Feng, S., Hersh, C.P., Bakke,
31 P., Gulsvik, A. *et al.* (2009) A genome-wide association study in chronic obstructive
32 pulmonary disease (COPD): identification of two major susceptibility loci. *PLoS Genet.*, **5**,
33 e1000421.

1 8 Wain, L.V., Shrine, N., Artigas, M.S., Erzurumluoglu, A.M., Noyvert, B., Bossini-Castillo,
2 L., Obeidat, M., Henry, A.P., Portelli, M.A., Hall, R.J. *et al.* (2017) Genome-wide association
3 analyses for lung function and chronic obstructive pulmonary disease identify new loci and
4 potential druggable targets. *Nat. Genet.*, **49**, 416-425.

5 9 Hobbs, B.D., de Jong, K., Lamontagne, M., Bosse, Y., Shrine, N., Artigas, M.S., Wain, L.V.,
6 Hall, I.P., Jackson, V.E., Wyss, A.B. *et al.* (2017) Genetic loci associated with chronic
7 obstructive pulmonary disease overlap with loci for lung function and pulmonary fibrosis.
8 *Nat. Genet.*, **49**, 426-432.

9 10 Zhou, J.J., Cho, M.H., Castaldi, P.J., Hersh, C.P., Silverman, E.K. and Laird, N.M. (2013)
10 Heritability of chronic obstructive pulmonary disease and related phenotypes in smokers.
11 *Am. J. Respir. Crit. Care Med.*, **188**, 941-947.

12 11 Laurell, C.B. and Eriksson, S. (2013) The electrophoretic alpha1-globulin pattern of
13 serum in alpha1-antitrypsin deficiency. 1963. *COPD*, **10 Suppl 1**, 3-8.

14 12 Urban, Z., Gao, J., Pope, F.M. and Davis, E.C. (2005) Autosomal dominant cutis laxa with
15 severe lung disease: synthesis and matrix deposition of mutant tropoelastin. *J. Invest.*
16 *Dermatol.*, **124**, 1193-1199.

17 13 Huchtagowder, V., Sausgruber, N., Kim, K.H., Angle, B., Marmorstein, L.Y. and Urban, Z.
18 (2006) Fibulin-4: a novel gene for an autosomal recessive cutis laxa syndrome. *Am. J. Hum.*
19 *Genet.*, **78**, 1075-1080.

20 14 Huchtagowder, V., Morava, E., Kornak, U., Lefeber, D.J., Fischer, B., Dimopoulou, A.,
21 Aldinger, A., Choi, J., Davis, E.C., Abuelo, D.N. *et al.* (2009) Loss-of-function mutations in
22 ATP6V0A2 impair vesicular trafficking, tropoelastin secretion and cell survival. *Hum. Mol.*
23 *Genet.*, **18**, 2149-2165.

24 15 Hu, Q., Loeys, B.L., Coucke, P.J., De Paepe, A., Mecham, R.P., Choi, J., Davis, E.C. and
25 Urban, Z. (2006) Fibulin-5 mutations: mechanisms of impaired elastic fiber formation in
26 recessive cutis laxa. *Hum. Mol. Genet.*, **15**, 3379-3386.

27 16 Callewaert, B., Su, C.T., Van Damme, T., Vlummens, P., Malfait, F., Vanakker, O., Schulz,
28 B., Mac Neal, M., Davis, E.C., Lee, J.G. *et al.* (2013) Comprehensive clinical and molecular
29 analysis of 12 families with type 1 recessive cutis laxa. *Hum. Mutat.*, **34**, 111-121.

30 17 Stanley, S.E., Merck, S.J. and Armanios, M. (2016) Telomerase and the Genetics of
31 Emphysema Susceptibility. Implications for Pathogenesis Paradigms and Patient Care. *Ann.*
32 *Am. Thorac. Soc.*, **13**, S447-S451.

33 18 Stanley, S.E., Gable, D.L., Wagner, C.L., Carlile, T.M., Hanumanthu, V.S., Podlevsky, J.D.,
34 Khalil, S.E., DeZern, A.E., Rojas-Duran, M.F., Applegate, C.D. *et al.* (2016) Loss-of-function
35 mutations in the RNA biogenesis factor NAF1 predispose to pulmonary fibrosis-emphysema.
36 *Sci. Transl. Med.*, **8**, 351ra107.

37 19 Stanley, S.E., Chen, J.J., Podlevsky, J.D., Alder, J.K., Hansel, N.N., Mathias, R.A., Qi, X.,
38 Rafaels, N.M., Wise, R.A., Silverman, E.K. *et al.* (2015) Telomerase mutations in smokers with
39 severe emphysema. *J. Clin. Invest.*, **125**, 563-570.

40 20 Do, R., Stitzel, N.O., Won, H.H., Jorgensen, A.B., Duga, S., Angelica Merlini, P., Kiezun, A.,
41 Farrall, M., Goel, A., Zuk, O. *et al.* (2015) Exome sequencing identifies rare LDLR and APOA5
42 alleles conferring risk for myocardial infarction. *Nature*, **518**, 102-106.

43 21 Ellinghaus, D., Zhang, H., Zeissig, S., Lipinski, S., Till, A., Jiang, T., Stade, B., Bromberg, Y.,
44 Ellinghaus, E., Keller, A. *et al.* (2013) Association between variants of PRDM1 and NDP52 and
45 Crohn's disease, based on exome sequencing and functional studies. *Gastroenterology*, **145**,
46 339-347.

1 22 Yu, T.W., Chahrour, M.H., Coulter, M.E., Jiralerspong, S., Okamura-Ikeda, K., Ataman, B.,
2 Schmitz-Abe, K., Harmin, D.A., Adli, M., Malik, A.N. *et al.* (2013) Using whole-exome
3 sequencing to identify inherited causes of autism. *Neuron*, **77**, 259-273.

4 23 Gonzaga-Jauregui, C., Harel, T., Gambin, T., Kousi, M., Griffin, L.B., Francescato, L., Ozes,
5 B., Karaca, E., Jhangiani, S.N., Bainbridge, M.N. *et al.* (2015) Exome Sequence Analysis Suggests
6 that Genetic Burden Contributes to Phenotypic Variability and Complex Neuropathy. *Cell*
7 *Rep.*, **12**, 1169-1183.

8 24 Lange, L.A., Hu, Y., Zhang, H., Xue, C., Schmidt, E.M., Tang, Z.Z., Bizon, C., Lange, E.M.,
9 Smith, J.D., Turner, E.H. *et al.* (2014) Whole-exome sequencing identifies rare and low-
10 frequency coding variants associated with LDL cholesterol. *Am. J. Hum. Genet.*, **94**, 233-245.

11 25 Qiao, D., Lange, C., Crapo, J.D., Beaty, T.H., Laird, N.M., Won, S., Silverman, E.K. and Cho,
12 M.H. (2016), In *American Thoracic Society International Conference 2016*, San Francisco, in
13 press.

14 26 Qiao, D., Lange, C., Beaty, T.H., Crapo, J.D., Barnes, K.C., Bamshad, M., Hersh, C.P.,
15 Morrow, J., Pinto-Plata, V.M., Marchetti, N. *et al.* (2016) Exome Sequencing Analysis in Severe,
16 Early-Onset Chronic Obstructive Pulmonary Disease. *Am. J. Resp. Crit. Care. Med.*, **193**, 1353-
17 1363.

18 27 Lek, M., Karczewski, K.J., Minikel, E.V., Samocha, K.E., Banks, E., Fennell, T., O'Donnell-
19 Luria, A.H., Ware, J.S., Hill, A.J., Cummings, B.B. *et al.* (2016) Analysis of protein-coding genetic
20 variation in 60,706 humans. *Nature*, **536**, 285-291.

21 28 Matsson, H., Soderhall, C., Einarsdottir, E., Lamontagne, M., Gudmundsson, S., Backman,
22 H., Lindberg, A., Ronmark, E., Kere, J., Sin, D. *et al.* (2016) Targeted high-throughput
23 sequencing of candidate genes for chronic obstructive pulmonary disease. *BMC Pulm. Med.*,
24 **16**, 146.

25 29 Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. and Morishima, K. (2017) KEGG: new
26 perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res.*, **45**, D353-D361.

27 30 Yew-Booth, L., Birrell, M.A., Lau, M.S., Baker, K., Jones, V., Kilty, I. and Belvisi, M.G.
28 (2015) JAK-STAT pathway activation in COPD. *Eur. Respir. J.*, **46**, 843-845.

29 31 Menche, J., Sharma, A., Kitsak, M., Ghiassian, S.D., Vidal, M., Loscalzo, J. and Barabasi,
30 A.L. (2015) Disease networks. Uncovering disease-disease relationships through the
31 incomplete interactome. *Science*, **347**, 1257601.

32 32 Chen, J., Xu, H., Aronow, B.J. and Jegga, A.G. (2007) Improved human disease candidate
33 gene prioritization using mouse phenotype. *BMC Bioinformatics*, **8**, 392.

34 33 Zandvoort, A., Postma, D.S., Jonker, M.R., Noordhoek, J.A., Vos, J.T., van der Geld, Y.M.
35 and Timens, W. (2006) Altered expression of the Smad signalling pathway: implications for
36 COPD pathogenesis. *Eur. Respir. J.*, **28**, 533-541.

37 34 Baraldo, S., Bazzan, E., Turato, G., Calabrese, F., Beghe, B., Papi, A., Maestrelli, P., Fabbri,
38 L.M., Zuin, R. and Sietta, M. (2005) Decreased expression of TGF-beta type II receptor in
39 bronchial glands of smokers with COPD. *Thorax*, **60**, 998-1002.

40 35 Ezzie, M.E., Crawford, M., Cho, J.H., Orellana, R., Zhang, S., Gelinas, R., Batte, K., Yu, L.,
41 Nuovo, G., Galas, D. *et al.* (2012) Gene expression networks in COPD: microRNA and mRNA
42 regulation. *Thorax*, **67**, 122-131.

43 36 Hessel, J., Heldrich, J., Fuller, J., Staudt, M.R., Radisch, S., Hollmann, C., Harvey, B.G.,
44 Kaner, R.J., Salit, J., Yee-Levin, J. *et al.* (2014) Intraflagellar transport gene expression
45 associated with short cilia in smoking and COPD. *PLoS One*, **9**, e85453.

46 37 Zhou, X., Qiu, W., Sathirapongsasuti, J.F., Cho, M.H., Mancini, J.D., Lao, T., Thibault, D.M.,
47 Litonjua, A.A., Bakke, P.S., Gulsvik, A. *et al.* (2013) Gene expression analysis uncovers novel

1 hedgehog interacting protein (HHIP) effects in human bronchial epithelial cells. *Genomics*,
2 **101**, 263-272.

3 38 Peloso, G.M., Auer, P.L., Bis, J.C., Voorman, A., Morrison, A.C., Stitzel, N.O., Brody, J.A.,
4 Khetarpal, S.A., Crosby, J.R., Fornage, M. *et al.* (2014) Association of low-frequency and rare
5 coding-sequence variants with blood lipids and coronary heart disease in 56,000 whites and
6 blacks. *Am. J. Hum. Genet.*, **94**, 223-232.

7 39 Beaudoin, M., Goyette, P., Boucher, G., Lo, K.S., Rivas, M.A., Stevens, C., Alikashani, A.,
8 Ladouceur, M., Ellinghaus, D., Torkvist, L. *et al.* (2013) Deep resequencing of GWAS loci
9 identifies rare variants in CARD9, IL23R and RNF186 that are associated with ulcerative
10 colitis. *PLoS Genet.*, **9**, e1003723.

11 40 Stuart, B.D., Choi, J., Zaidi, S., Xing, C., Holohan, B., Chen, R., Choi, M., Dharwadkar, P.,
12 Torres, F., Girod, C.E. *et al.* (2015) Exome sequencing links mutations in PARN and RTEL1
13 with familial pulmonary fibrosis and telomere shortening. *Nat. Genet.*, **47**, 512-517.

14 41 Diaz de Leon, A., Cronkhite, J.T., Katzenstein, A.L., Godwin, J.D., Raghu, G., Glazer, C.S.,
15 Rosenblatt, R.L., Girod, C.E., Garrity, E.R., Xing, C. *et al.* (2010) Telomere lengths, pulmonary
16 fibrosis and telomerase (TERT) mutations. *PLoS One*, **5**, e10680.

17 42 Petrovski, S., Todd, J.L., Durham, M.T., Wang, Q., Chien, J.W., Kelly, F.L., Frankel, C.,
18 Mebane, C.M., Ren, Z., Bridgers, J. *et al.* (2017) An Exome Sequencing Study to Assess the Role
19 of Rare Genetic Variation in Pulmonary Fibrosis. *Am. J. Respir. Crit. Care Med.*, **196**, 82-93.

20 43 Tsakiri, K.D., Cronkhite, J.T., Kuan, P.J., Xing, C., Raghu, G., Weissler, J.C., Rosenblatt, R.L.,
21 Shay, J.W. and Garcia, C.K. (2007) Adult-onset pulmonary fibrosis caused by mutations in
22 telomerase. *Proc. Natl. Acad. Sci. U. S. A.*, **104**, 7552-7557.

23 44 Bertuch, A.A. (2016) The molecular genetics of the telomere biology disorders. *RNA*
24 *Biology*, **13**, 696-706.

25 45 Zhu, G., Warren, L., Aponte, J., Gulsvik, A., Bakke, P., Anderson, W.H., Lomas, D.A.,
26 Silverman, E.K., Pillai, S.G. and International, C.G.N.I. (2007) The SERPINE2 gene is associated
27 with chronic obstructive pulmonary disease in two large populations. *Am. J. Respir. Crit. Care*
28 *Med.*, **176**, 167-173.

29 46 Elleder, M., Houstkova, H., Zeman, J., Ledvinova, J. and Poupetova, H. (2001)
30 Pulmonary storage with emphysema as a sign of Niemann-Pick type C2 disease (second
31 complementation group). Report of a case. *Virchows Arch.*, **439**, 206-211.

32 47 Zhang, Y., Calado, R., Rao, M., Hong, J.A., Meeker, A.K., Dumitriu, B., Atay, S., McCormick,
33 P.J., Garfield, S.H., Wangsa, D. *et al.* (2014) Telomerase variant A279T induces telomere
34 dysfunction and inhibits non-canonical telomerase activity in esophageal carcinomas. *PLoS*
35 *One*, **9**, e101010.

36 48 Dumanski, J.P., Rasi, C., Bjorklund, P., Davies, H., Ali, A.S., Gronberg, M., Welin, S.,
37 Sorbye, H., Gronbaek, H., Cunningham, J.L. *et al.* (2017) A MUTYH germline mutation is
38 associated with small intestinal neuroendocrine tumors. *Endocr. Relat. Cancer*, **24**, 427-443.

39 49 Kho, A.T., Chhabra, D., Sharma, S., Qiu, W., Carey, V.J., Gaedigk, R., Vyhldal, C.A., Leeder,
40 J.S., Tantisira, K.G. and Weiss, S.T. (2016) Age, Sexual Dimorphism, and Disease Associations
41 in the Developing Human Fetal Lung Transcriptome. *Am. J. Respir. Cell Mol. Biol.*, **54**, 814-821.

42 50 Beghe, B., Bazzan, E., Baraldo, S., Calabrese, F., Rea, F., Loy, M., Maestrelli, P., Zuin, R.,
43 Fabbri, L.M. and Saetta, M. (2006) Transforming growth factor-beta type II receptor in
44 pulmonary arteries of patients with very severe COPD. *Eur. Respir. J.*, **28**, 556-562.

45 51 Cho, M.H., McDonald, M.L., Zhou, X., Mattheisen, M., Castaldi, P.J., Hersh, C.P., Demeo,
46 D.L., Sylvia, J.S., Ziniti, J., Laird, N.M. *et al.* (2014) Risk loci for chronic obstructive pulmonary
47 disease: a genome-wide association study and meta-analysis. *Lancet Respir. Med.*, **2**, 214-225.

1 52 De Winter, J.P., De Vries, C.J., Van Achterberg, T.A., Ameerun, R.F., Feijen, A., Sugino, H.,
2 De Waele, P., Huylebroeck, D., Verschueren, K. and Van Den Eijden-Van Raaij, A.J. (1996)
3 Truncated activin type II receptors inhibit bioactivity by the formation of heteromeric
4 complexes with activin type I. receptors. *Exp. Cell Res.*, **224**, 323-334.

5 53 McDonald, M.L., Mattheisen, M., Cho, M.H., Liu, Y.Y., Harshfield, B., Hersh, C.P., Bakke, P.,
6 Gulsvik, A., Lange, C., Beaty, T.H. *et al.* (2014) Beyond GWAS in COPD: probing the landscape
7 between gene-set associations, genome-wide associations and protein-protein interaction
8 networks. *Hum. Hered.*, **78**, 131-139.

9 54 Boueiz, A., Chase, R., Lamb, A., Lee, S., Naing, Z.Z.C., Cho, M.H., Parker, M.M., Hersh, C.P.,
10 Crapo, J.D., Stergachis, A.B. *et al.* (2017) Integrative genomics analysis identifies ACVR1B as a
11 candidate causal gene of emphysema distribution in non-alpha 1-antitrypsin deficient
12 smokers. *bioRxiv*, in press.

13 55 Tilley, A.E., Walters, M.S., Shaykhiev, R. and Crystal, R.G. (2015) Cilia dysfunction in
14 lung disease. *Annu. Rev. Physiol.*, **77**, 379-406.

15 56 Wain, L.V., Sayers, I., Soler Artigas, M., Portelli, M.A., Zeggini, E., Obeidat, M., Sin, D.D.,
16 Bosse, Y., Nickle, D., Brandsma, C.A. *et al.* (2014) Whole exome re-sequencing implicates
17 CCDC38 and cilia structure and function in resistance to smoking related airflow obstruction.
18 *PLoS Genet.*, **10**, e1004314.

19 57 Kenna, K.P., van Doormaal, P.T., Dekker, A.M., Ticozzi, N., Kenna, B.J., Diekstra, F.P., van
20 Rheenen, W., van Eijk, K.R., Jones, A.R., Keagle, P. *et al.* (2016) NEK1 variants confer
21 susceptibility to amyotrophic lateral sclerosis. *Nat. Genet.*, **48**, 1037-1042.

22 58 Regan, E.A., Hokanson, J.E., Murphy, J.R., Make, B., Lynch, D.A., Beaty, T.H., Curran-
23 Everett, D., Silverman, E.K. and Crapo, J.D. (2010) Genetic epidemiology of COPD (COPDGene)
24 study design. *COPD*, **7**, 32-43.

25 59 Silverman, E.K., Chapman, H.A., Drazen, J.M., Weiss, S.T., Rosner, B., Campbell, E.J.,
26 O'Donnell, W.J., Reilly, J.J., Ginns, L., Mentzer, S. *et al.* (1998) Genetic epidemiology of severe,
27 early-onset chronic obstructive pulmonary disease. Risk to relatives for airflow obstruction
28 and chronic bronchitis. *Am. J. Respir. Crit. Care Med.*, **157**, 1770-1778.

29 60 Cingolani, P., Platts, A., Wang le, L., Coon, M., Nguyen, T., Wang, L., Land, S.J., Lu, X. and
30 Ruden, D.M. (2012) A program for annotating and predicting the effects of single nucleotide
31 polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2;
32 iso-3. *Fly (Austin)*, **6**, 80-92.

33 61 Chen, H., Wang, C., Conomos, M.P., Stilp, A.M., Li, Z., Sofer, T., Szpiro, A.A., Chen, W.,
34 Brehm, J.M., Celedon, J.C. *et al.* (2016) Control for Population Structure and Relatedness for
35 Binary Traits in Genetic Association Studies via Logistic Mixed Models. *Am. J. Hum. Genet.*, **98**,
36 653-666.

37 62 Jiang, D. and McPeck, M.S. (2014) Robust rare variant association testing for
38 quantitative traits in samples with related individuals. *Genet. Epidemiol.*, **38**, 10-20.

39 63 Shihab, H.A., Gough, J., Cooper, D.N., Stenson, P.D., Barker, G.L., Edwards, K.J., Day, I.N.
40 and Gaunt, T.R. (2013) Predicting the functional, molecular, and phenotypic consequences of
41 amino acid substitutions using hidden Markov models. *Hum. Mutat.*, **34**, 57-65.

42 64 Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A.,
43 Paulovich, A., Pomeroy, S.L., Golub, T.R., Lander, E.S. *et al.* (2005) Gene set enrichment
44 analysis: a knowledge-based approach for interpreting genome-wide expression profiles.
45 *Proc. Natl. Acad. Sci. U. S. A.*, **102**, 15545-15550.

46 65 Morrow, J., Qiu, W., D.L., D., Houston, I., Pinto Plata, V.M., B.R., C., Marchetti, N., Criner,
47 G.J., Bueno, R. and Washko, G.R. (2015) Network Analysis of Gene Expression in Severe COPD
48 Lung Tissue Samples (abstract). *Am. J. Respir. Crit. Care Med.*, **191**.

1 66 Blake, J.A., Eppig, J.T., Kadin, J.A., Richardson, J.E., Smith, C.L., Bult, C.J. and the Mouse
2 Genome Database, G. (2017) Mouse Genome Database (MGD)-2017: community knowledge
3 resource for the laboratory mouse. *Nucleic Acids Res.*, **45**, D723-D729.

4 **Legends to Figures**

5 **Figure 1. Network of the two sets of top genes with p-value < 0.01 in the case-control and family-based**
6 **analyses focusing on rare, deleterious variants.** The nodes in red are the top genes identified in the
7 COPDGene case-control analysis; the nodes in pink are the first neighbors of the red nodes. The nodes in
8 blue are the top genes identified in the BEOCOPD-ICGN family-based analysis; the nodes in light blue are
9 the first neighbors of the blue nodes. Genes that are in both sets are colored in purple. Edges connecting
10 genes to the largest hub *TULP3* are colored in red. These genes form one large well-connected component.
11 Larger sized nodes indicate hubs (circle) and genes reported to be associated with COPD or lung function
12 (diamond). Hubs include *TULP3*, *VNN2*, *ACVR2B*, *KCNA5*, and *GRB2*, which are the top genes with the
13 most number of degrees in this network. 14 genes out of 513 are near GWAS loci for COPD or lung
14 function (*CHRM3*, *DNLZ*, *EFEMP1*, *EFEMP2*, *EGFL8*, *GANAB*, *GNG3*, *PARN*, *PIP4K2B*, *NOTCH4*,
15 *RUVBL1*, *SEC16A*, *TARS*, *TEKT5*, *THSD4*). The zoomed-in panel on the left shows the genes in the
16 intraciliary transport particle A pathway (GO:0030991). The zoomed-in panel on the right shows the genes
17 in the transforming growth factor beta-receptor binding pathway (GO:0005160) and *ACVR2B*.

18 **Figure 2. A flow chart of the study design.** COPDGene (pink) samples were sequenced in two batches
19 (Baylor and ESP, see Methods). The family-based studies (blue) included two cohorts. Forty-nine pedigrees
20 of the Boston Early-Onset study samples were sequenced and analyzed previously (26); we combined these
21 data with another subset of these BEOCOPD and additional samples from the ICGN study. All of these
22 sequenced subjects from BEOCOPD and ICGN were called together, forming the BEOCOPD-ICGN dataset
23 (blue). We applied the family-based GESE test to the most severe cases and resistant controls in the
24 BEOCOPD-ICGN dataset. We also performed single-variant, gene-based and pathway-based association
25 tests in COPDGene and the BEOCOPD-ICGN samples. A final network analysis was conducted to look at
26 the topological relationship between the top results from the two datasets.

27

1 **Tables**

2 **Table 1. Baseline characteristics.**

Datasets	COPDGene		ICGN-BEOCOPD						
	ESP		Baylor		Severe Cases ¹	Moderate Cases	Resistant controls ¹	Other Controls	Other
	Case	Control	Case	Control					
N	192	188	293	316	853	431	101	118	51 ²
# Females	92	103	117	146	412	199	53	68	30
# Males	100	85	176	170	441	232	48	50	21
Age, yr	58.2 (5.1)	69.5 (5.6)	68 (6.4)	61.9 (5.6)	56.1 (12.2)	59.4 (12.1)	55.5 (11.3)	37.0 (20.4)	53.9 (24.4)
Pack-years	45.0 (26.2)	45.0 (23.5)	51.0 (29.0)	50.6 (19.1)	43.8 (31.0)	37.4 (25.6)	30.8 (21.9)	0.0 (7.6)	24.0 (39.9)
FEV ₁ % predicted	30.0 (15.9)	98.2 (12.8)	30.2 (15.8)	92.7 (14.3)	30.0 (17.6)	65.5 (14.1)	98.0 (15.3)	96.3 (14.6)	77.4 (10.6)
FEV1/FVC	0.33 (0.10)	0.78 (0.07)	0.35 (0.12)	0.76 (0.07)	0.33 (0.14)	0.56 (0.13)	0.76 (0.06)	0.81 (0.10)	0.71 (0.12)

3 N – Number of subjects

4 Median (IQR) is presented for age, pack-years, FEV₁ % predicted, and FEV₁/FVC ratio for each dataset.

5 ¹ Only severe cases (GOLD level III and IV) and resistant controls (see text) were included in the GESE test
6 of the family-based data. All subjects were included in the association analysis of the family-based data.

7 ² 51 subjects in the ICGN-BEOCOPD data had lung function values not consistent with either case or
8 control status.

9

10

11 **Table 2. Results of the GESE analysis on the BEOCOPD-ICGN dataset.**

GENE	GESE p-value	Number of segregating families
<i>TBC1D10A</i> *	1.1E-06	2
<i>RFPL1</i> *	1.6E-06	4
<i>DHODH</i> *#	6.9E-05	2
<i>CYP4F12</i> *	1.0E-04	4
<i>ANAPC7</i> *	1.5E-04	1
<i>RGS5</i> *#	1.5E-04	2
<i>CD101</i> *#	1.7E-04	5
<i>KCNMB4</i> *#	1.8E-04	1
<i>ARMC12</i> *	2.1E-04	4
<i>VPS41</i> *#	3.9E-04	5

1 Variants included are loss-of-function variants with MAF < 0.1%. The third column shows the number of
 2 families each gene is segregating in (present in all the cases and not in the controls). Genes marked with *
 3 show expression in the lung (defined as at least 50% of samples with FPKM > 0.5 in the Lung Genomics
 4 Resource Consortium RNA-seq samples). Genes marked with # are differentially expressed by FEV1%
 5 predicted in lung tissue (65).

6

7

8 **Table 3. Nominally significant gene-based results in COPDGene for 329 candidate genes after**
 9 **conditioning on the lead GWAS SNP.**

Gene	(conditioned on) lead GWAS SNP	#SNV	SKATO (unadjusted)	SKATO (conditional)
<i>SEC16A</i>	rs10870202	38	9.52E-04	1.05E-03
<i>CDC7</i>	rs1192404	5	5.87E-03	6.78E-03
<i>CCDC38</i>	rs12820313	5	7.20E-03	7.66E-03
<i>CYP2A6</i>	rs12459249	11	6.87E-03	1.15E-02
<i>TRIP11</i>	rs7155279	24	1.53E-02	2.47E-02
<i>CNGBI</i>	rs12447804	26	3.27E-02	2.89E-02
<i>PBLD</i>	rs7095607	3	6.99E-02	3.43E-02
<i>RRP15</i>	rs10429950	2	3.75E-02	4.07E-02
<i>TNXB</i>	rs2070600	63	5.95E-02	4.39E-02
<i>CYFIP2</i>	rs10515750, rs1990950	4	4.13E-02	4.46E-02
<i>EGFL8</i>	rs2070600	9	8.11E-02	4.60E-02
<i>CHRNA5</i>	rs17486278	5	1.83E-02	1.22E-01
<i>AGER</i>	rs2070600	12	4.89E-03	1.36E-01

10 The SKAT-O tests included functional (MODERATE effect defined by SnpEff) and rare (MAF < 5%)
 11 variants in the COPDGene study.

12

13

14

15

1 **Table 4. Selected set of likely pathogenic variants annotated by ClinVar in *SERPINA1* and telomere-related genes.**

GENE	SNP	COPDGene		BECOPD-ICGN		MAF	IMPACT	CLNSIG	Disease association
		Case	Cont	Case	Cont				
<i>RTEL1</i>	20:62324513:T:C	0	1	.	.	6.11E-05	missense_variant	5	Telomeropathy
<i>SERPINA1</i>	rs28929474	29	12	60	8	1.83E-02	missense_variant	5 5	Alpha-1 antitrypsin deficiency
<i>SERPINA1</i>	rs17580	53	29	148	32	3.04E-02	missense_variant	5	Alpha-1 antitrypsin deficiency
<i>SERPINA1</i>	rs28929470	4	6	5	1	4.95E-03	missense_variant	5	Alpha-1 antitrypsin deficiency
<i>SERPINA1</i>	rs28931570	2	4	8	2	1.62E-03	missense_variant	4 5	Alpha-1 antitrypsin deficiency
<i>SERPINA1</i>	rs121912714	.	.	3	0	7.04E-04	missense_variant	4	Alpha-1 antitrypsin deficiency
<i>TERT</i>	rs61748181	40	33	92	16	4.97E-02	missense_variant	5 2	Telomeropathy
<i>TERT</i>	5:1278865	1	0	.	.	7.49E-05	missense_variant	5	Telomeropathy
<i>TERT</i>	5:1280427	1	0	.	.	.	missense_variant	.	Telomeropathy
<i>TERT</i>	rs35719940	27	22	.	.	2.11E-02	missense_variant	5 2 2 5 5 5 3 2	Telomeropathy
<i>TERT</i>	rs34094720	4	7	.	.	1.53E-02	missense_variant	2	Telomeropathy
<i>TERT</i>	rs141425941	1	1	.	.	2.68E-04	missense_variant	5	Telomeropathy
<i>TINF2</i>	rs142777869	2	1	1	0	7.25E-04	missense_variant	5	Telomeropathy

2

3 Case, Cont - number of alternative alleles carried by the cases and controls in each dataset. Note that in the family-based data, there are approximately 6

4 times more cases than controls. MAF - minor allele frequency using Non-Finnish Europeans from ExAC. IMPACT - functional impact of each variant

5 annotated by SnpEff. CLNSIG and Disease association are annotations from ClinVar; 2 = benign, 3 = likely benign, 4 = likely pathogenic, 5= pathogenic.

6

7

1

2 **Table 5. Homozygous variants in COPD-related Mendelian genes.**

GENE	SNP	COPDGene		BECOPD-ICGN		MAF	IMPACT	CLNSIG	Disease association
		Case	Cont	Case	Cont				
<i>ATP6V0A2</i>	12:124206896*	0	0	1	0	.	Splice_acceptor_variant	.	Cutis laxa
<i>CFTR</i>	rs1800076	0	1	3	0	2.48E-02	missense_variant	2 2 5	Cystic fibrosis
<i>NPC2</i>	rs140130028	.	.	2	0	0.00551	splice_donor_variant	5	Niemann-Pick disease type C2
<i>SERPINA1</i>	rs17580	1	1	4	1	3.04E-02	missense_variant	5	Apha-1-antitrypsin deficiency
<i>TERT</i>	rs61748181	5	0	2	0	4.97E-02	missense_variant	5 2	Telomeropathy
<i>TERT</i>	rs35719940	1	0	.	.	2.11E-02	missense_variant	5 2 2	Telomeropathy

3

4 Variants with homozygous genotypes in 29 Mendelian genes and were annotated with significance 4 and above by ClinVar, or have potential deleterious
5 effects (MAF < 0.1% and predicted to be deleterious by FATHMM, SIFT, and CADD (>15). Case, Cont - number of alternative alleles carried by the cases
6 and controls in each dataset. Note that in the family-based data, there are approximately 6 times more cases than controls. MAF - minor allele frequency
7 using Non-Finnish Europeans from ExAC. IMPACT - functional impact of each variant annotated by SnpEff. CLNSIG and Disease association are
8 annotations from ClinVar; 2 = benign, 3 = likely benign, 4 = likely pathogenic, 5 = pathogenic.

9

10

11 Abbreviations

12 COPD: Chronic obstructive pulmonary disease

13 GWAS: Genome-wide association study

14 MAF: Minor allele frequency

15 GESE: Gene-based segregation test

16 FEV₁: Forced expiratory volume in 1 second

17 GOLD: Global Initiative for Chronic Obstructive Lung Disease

18 FPKM: Fragments Per Kilobase of transcript per Million mapped reads

19 PPI: Protein-protein interaction network