

1 **A *MUC5B* gene polymorphism, rs35705950-T, confers protective effects in COVID-19 infection.**

2

3 Anurag Verma PhD^{1,2†}, Jessica Minnier PhD^{3,4,5†}, Jennifer E Huffman PhD^{6†}, Emily S Wan MD, MPH^{7,8†}, Lina Gao
4 PhD^{9,10}, Jacob Joseph MBBS, MD^{11,12}, Yuk-Lam Ho MPH¹³, Wen-Chih Wu MD, MPH^{14,15}, Kelly Cho PhD,
5 MPH^{16,17}, Bryan R Gorman PhD¹⁸, Nallakkandi Rajeevan PhD^{19,20}, Saiju Pyarajan^{21,22}, Helene Garcon MD⁶, James B
6 Meigs MD, MPH²³, Yan V Sun PhD^{24,25}, Peter D Reaven MD^{26,27}, John E McGeary PhD^{28,29}, Ayako Suzuki MD, PhD,
7 MSc^{30,31}, Joel Gelernter MD^{32,33}, Julie A Lynch PhD, RN^{34,35}, Jeffrey M Petersen MD^{36,37}, Seyedeh Maryam Zekavat
8 BS^{38,39}, Pradeep Natarajan MD, MMSc^{40,39,41}, Cecelia J Madison MBI⁴², Sharvari Dalal MD^{43,37}, Darshana N Jhala MD,
9 B Mus^{44,45}, Mehrdad Arjomandi MD⁴⁶, Elise Gatsby MPH⁴⁷, Kristine E Lynch MS PhD^{47,48}, Robert A Bonomo⁴⁹, Mat
10 Freiberg⁵⁰, Gita A Pathak PhD^{51,52}, Jin J Zhou PhD^{53,54}, Curtis J Donskey MD⁵⁵, Ravi K Madduri MS⁵⁶, Quinn S Wells
11 MD, PharmD⁵⁷, Rose DL Huang MS⁶, Renato Polimanti PhD^{51,52}, Kyong-Mi Chang MD⁵⁸, Katherine P Liao MD,
12 MPH⁵⁹, Philip S Tsao PhD⁶⁰, Peter W.F. Wilson MD^{25,61}, Adriana Hung⁶², Christopher J O'Donnell MD, MPH⁶³, John M
13 Gaziano⁶⁴, Richard L Hauger MD^{65,66}, Sudha K. Iyengar PhD^{67,68*} and Shih-Wen Luoh MD PhD^{69*}

14

15 ¹Corporal Michael Crescenz VA Medical Center, Philadelphia, PA, 19104, USA, ²Department of Medicine, Perelman
16 School of Medicine, University of Pennsylvania, Philadelphia, PA, 19104, USA, ³OHSU-PSU School of Public Health,
17 Oregon Health & Science University, Portland, OR, 97239, USA, ⁴Knight Cancer Institute, Biostatistics Shared
18 Resource, Oregon Health & Science University, Portland, OR, 97239, USA, ⁵VA Portland Health Care System, OR,
19 97239, USA, ⁶MAVERIC, VA Boston Healthcare System, Boston, MA, 02130, USA, ⁷Department of Medicine,
20 Pulmonary, Critical Care, Sleep, and Allergy Section, VA Boston Healthcare System, Boston, MA, 02115,
21 USA, ⁸Channing Division of Network Medicine, Brigham & Women's Hospital, Boston, MA, 02115, USA, ⁹Knight
22 Cancer Institute, Biostatistics Shared Resource, Oregon Health & Science University, ¹⁰VA Portland Health Care
23 System, ¹¹Department of Medicine, VA Boston Healthcare System, Boston, MA 02115, USA, ¹²Brigham & Women's
24 Hospital, ¹³MAVERIC, VA Boston Healthcare System, Boston, MA, 2130, USA, ¹⁴Department of Medicine,
25 Cardiology, Providence VA Healthcare System, Providence, RI, 02908, USA, ¹⁵Alpert Medical School & School of
26 Public Health, Brown University, Providence, RI, 02903, USA, ¹⁶MAVERIC, VA Boston Healthcare System, 150 S
27 Huntington Ave, Boston, MA, 02130, USA, ¹⁷Medicine, Aging, Brigham and Women's Hospital, Harvard Medical
28 School, Boston, MA, 02130, USA, ¹⁸VA Boston Healthcare System, 150 S Huntington Ave, Boston, MA, 02131,
29 USA, ¹⁹Yale Center for Medical Informatics, Yale School of Medicine, New Haven, CT, 06511', USA, ²⁰Clinical
30 Epidemiology Research Center (CERC), VA Connecticut Healthcare System, West Haven, CT, 06516', USA, ²¹VA
31 Boston Healthcare System, Boston, MA, 2115, ²²Department of Medicine, Harvard Medical School, ²³Medicine,
32 General Internal Medicine, Massachusetts General Hospital, 55 Fruit St, Boston, MA, 02115, USA, ²⁴Epidemiology,
33 Emory University School of Public Health, 1518 Clifton Rd. NE, Atlanta, GA, 30322, USA, ²⁵Atlanta VA Health Care
34 System, 1670 Clairmont Road, Decatur, GA, 30033, USA, ²⁶Department of Medicine, Phoenix VA Healthcare
35 System, Phoenix, AZ, 85012, USA, ²⁷Univ. of AZ, ²⁸Department of Psychiatry and Human Behavior, Providence VA
36 Medical Center, Providence, RI, 2908, USA, ²⁹Brown University Medical School, ³⁰Department of Medicine,
37 Gastroenterology, Durham VA Medical Center, 508 Fulton St, Durham, NC, 27705, USA, ³¹Department of Medicine,
38 Gastroenterology, Duke University, Durham, NC, 27710, USA, ³²Psychiatry, Human Genetics, Yale Univ. School of
39 Medicine, 950 Campbell Avenue, West Haven, CT, 06516, USA, ³³VA CT Healthcare Center, ³⁴VA Informatics &
40 Computing Infrastructure, VA Salt Lake City Health Care System, Salt Lake City, UT, USA, ³⁵Department of Medicine,
41 University of Utah School of Medicine, Salt Lake City, UT, ³⁶Pathology and Laboratory Medicine, Corporal Michael
42 Crescenz VA Medical Center, Philadelphia, PA, 19104, USA, ³⁷Perelman School of Medicine, University of
43 Pennsylvania, Philadelphia, PA, 19104, USA, ³⁸Computational Biology & Bioinformatics, Yale School of Medicine,
44 333 Cedar St, New Haven, CT, 6510, USA, ³⁹Program in Medical and Population Genetics, Cardiovascular Disease
45 Initiative, Broad Institute of Harvard and MIT, Cambridge, MA, 02142, USA, ⁴⁰Cardiovascular Research Center,
46 Massachusetts General Hospital, Boston, MA, 02114, USA, ⁴¹Department of Medicine, Harvard Medical School,
47 Boston, MA, 02115, USA, ⁴²Clinical Data Science Research Group, ORD, Portland VA Medical Center, Portland, OR,
48 97239, USA, ⁴³Pathology and Laboratory Medicine, Corporal Michael J Crescenz VA Medical Center, Philadelphia,
49 PA, 19104, USA, ⁴⁴Pathology and Laboratory Medicine, Corporal Michael Crescenz VA Medical Center, Philadelphia,
50 PA, 19104, ⁴⁵Perelman School of Medicine, University of Pennsylvania, 19104, ⁴⁶Medicine, Pulmonary and Critical
51 Care, San Francisco VA Healthcare System; University of California San Francisco, San Francisco, CA, 94121,

52 USA, ⁴⁷VA Informatics and Computing Infrastructure (VINCI), VA Salt Lake City Healthcare System, Salt Lake City,
53 UT, ⁴⁸Internal Medicine, Epidemiology, University of Utah School of Medicine, Salt Lake City, UT, ⁴⁹Case Western
54 Reserve University, Cleveland, OH, ⁵⁰Vanderbilt University Medical Center, Nashville, TN, ⁵¹Department of
55 Psychiatry, Division of Human Genetics, Yale School of Medicine, New Haven, CT, 06511, USA, ⁵²VA Connecticut
56 Healthcare System, West Haven, CT, 06516, USA, ⁵³Medicine, University of California, Los Angeles, Los Angeles, CA,
57 90024, USA, ⁵⁴Epidemiology and Biostatistics, University of Arizona, AZ, 85724, USA, ⁵⁵Infectious Disease Section,
58 Louis Stokes Cleveland VA and Case Western Reserve University, Cleveland, OH, 44106, USA, ⁵⁶Data Science and
59 Learning, Argonne National Laboratory, 9700 S Cass Ave, Lemont, IL, 60439, ⁵⁷Departments of Medicine,
60 Biomedical Informatics, and Pharmacology, Vanderbilt University Medical Center, Nashville, TN, 37232,
61 USA, ⁵⁸Corporal Michael J Crescenz VA Medical Center, ⁵⁹Medicine, Rheumatology, VA Boston Healthcare System,
62 Boston, MA, 02130, ⁶⁰Precision Medicine, VA Palo Alto Health Care System, 3801 Miranda Avenue, Palo Alto, CA,
63 94304, USA, ⁶¹Emory University School of Medicine, Atlanta, GA, 30322, USA, ⁶²(No affiliation data
64 provided), ⁶³Medicine, Cardiology, VA Boston Healthcare System, 1400 VFW Parkway, Boston, MA, 2132, USA, ⁶⁴VA
65 Boston Healthcare System, Boston, MA, ⁶⁵Center of Excellence for Stress & Mental Health, VA San Diego
66 Healthcare System, San Diego, CA, 92161, USA, ⁶⁶Center for Behavioral Genetics of Aging, University of California
67 San Diego, La Jolla, CA, 92093, USA, ⁶⁷Case Western Reserve University, Cleveland, OH, 44106, USA, ⁶⁸Louis Stokes
68 Cleveland VA Medical Center, Cleveland, OH, 44106, USA, ⁶⁹VA Portland Health Care System, Portland, OR, 97239,
69 USA

70 **Joint Authorship**

71 [†]These authors contributed equally to this work

72 ^{*}These authors jointly supervised this work

73

74 Supported by MVP035 award from Million Veteran Program, Office of Research and Development,
75 Veterans Health Administration, and Veteran Affairs of the United States Government BX 004831
76 (P.W/K.C). This publication does not represent the views of the Department of Veteran Affairs of the
77 United States Government.

78

79 Author contributions: A.V, J.E.H, S.I, S.L, L.G and J.M and E.S.W. analyzed data. J.E.H., A.V., L.G., E.S.W.,
80 S.I., S.L. supervised data collection. A.V., J.E.H, L.G., J.M., E.S.W., S.I. and S.L. wrote the manuscript. A.V.,
81 J.E.H, L.G. J.M., E.S.W., S.I. and S.L. conceived the study design, supervised data collection and analysis,
82 and wrote and edited the manuscript. All authors revised and approved the final version to be
83 published.

84

85 Correspondence should be addressed to Shih-Wen Luoh (Shih-Wen.Luoh@va.gov), Sudha K Iyengar
86 (ski@case.edu)

87

88 Manuscript word count: **2,896**

89

90

91

92

93

94

95 **Abstract**

96 **Rationale:** A common *MUC5B* gene polymorphism, rs35705950-T, is associated with idiopathic
97 pulmonary fibrosis, but its role in the SARS-CoV-2 infection and disease severity is unclear.

98 **Objectives:** To assess whether rs35705950-T confers differential risk for clinical outcomes associated
99 with COVID-19 infection among participants in the Million Veteran Program (MVP) and COVID-19 Host
100 Genetics Initiative (HGI).

101 **Methods:** MVP participants were examined for an association between the incidence or severity of
102 COVID-19 and the presence of a *MUC5B* rs35705950-T allele. Comorbidities and clinical events were
103 extracted from the electronic health records (EHR). The analysis was performed within each ancestry
104 group in the MVP, adjusting for sex, age, age² and first twenty principal components followed by a
105 trans-ethnic meta-analysis. We then pursued replication and performed a meta-analysis with the trans-
106 ethnic summary statistics from the HGI. A phenome-wide association study (PheWAS) of the
107 rs35705950-T was conducted to explore associated pathophysiologic conditions.

108 **Measurements and Main Results:** A COVID-19 severity scale was modified from the World Health
109 Organization criteria, and phenotypes derived from the International Classification of Disease-9/10 were
110 extracted from EHR. Presence of rs35705950-T was associated with fewer hospitalizations
111 ($N_{\text{cases}}=25353$, $N_{\text{controls}}=631,024$; OR=0.86 [0.80-0.93], $p=7.4 \times 10^{-5}$) in trans-ethnic meta-analysis within
112 MVP and joint meta-analyses with the HGI ($N=1641311$; OR=0.89 [0.85-0.93], $p=1.9 \times 10^{-6}$). Moreover,
113 individuals of European Ancestry with at least one copy of rs35705950-T had fewer post-COVID-19
114 pneumonia events (OR=0.85 [0.76-0.96], $p=0.008$). PheWAS exclusively revealed pulmonary
115 involvement.

116 **Conclusions:** The *MUC5B* variant rs35705950-T is protective in COVID-19 infection.

117 **Keywords:** coronavirus disease 2019; severe acute respiratory syndrome coronavirus 2; idiopathic
118 pulmonary fibrosis; electronic health records; genetic association

119 Introduction

120 A respiratory disease caused by a novel coronavirus was first reported towards the end of 2019, now
121 known as SARS-CoV-2 (COVID-19). Despite massive public health measures and vaccination initiatives,
122 the COVID-19 pandemic remains a major global health threat. By September 2, 2021, the coronavirus
123 disease-2019 (COVID-19) pandemic had caused more than 219 million confirmed infections and more
124 than 4.5 million deaths worldwide(1).

125

126 Parenchymal fibrosis is a late complication of respiratory infections with COVID-19(2–4). Among chronic
127 lung diseases, idiopathic pulmonary fibrosis (IPF)(5), a disorder characterized by progressive pulmonary
128 scarring which is associated with a median survival of 2-3 years in the absence of lung
129 transplantation(6), shares several risk factors with those for severe COVID-19 disease, including
130 advanced age(7), cardiovascular disease, diabetes, and history of smoking(5). Thus, common
131 pathological processes may be shared between the fibrotic response towards COVID-19 infection and
132 those underlying IPF.

133

134 IPF likely develops from a multifaceted interaction between genetic and environmental factors, age-
135 related mechanisms, and epigenetic profibrotic reprogramming(8, 9). One of the most robust genetic
136 risk factors identified for IPF susceptibility is rs35705950-T, a common G to T transversion located
137 approximately 3 kb upstream of the mucin 5B, oligomeric mucus/gel-forming *MUC5B* gene (10, 11).

138 Laboratory evidence supports that rs35705950-T is: 1) a functional variant located within an enhancer
139 subject to epigenetic programming and 2) contributes to pathologic mis-expression in IPF (12).

140

141 Given the high minor allele frequency (MAF) of rs35705950-T (~20% among individuals of European
142 ancestry) and possible shared pathophysiological pathways between IPF and severe COVID-19 disease,

143 we examined the association between rs35705950-T and the clinical outcomes of COVID-19 infection in
144 the Million Veteran Program (MVP), a multi-ethnic cohort of over 650,000 U.S. Veterans with detailed
145 EHR and genotyping data(13). Following our primary analysis in the MVP, we validated our results with
146 a comparable analysis conducted in the Host Genetics Initiative (HGI), a global collaboration of over 160
147 genetic studies assembled to facilitate rapid discovery and dissemination of COVID-19 related science
148 (14).

149

150 **Methods**

151 **Data Sources**

152 Data from the MVP, a multi-ethnic genetic biobank sponsored by the United States Veterans Affairs
153 (VA), were analyzed (13). All protocols were approved by the VA Central Institutional Review Board and
154 all participants provided written informed consent. For detailed Materials and Methods, please see
155 **methods in the online data supplement.**

156

157 Demographic and pre-existing comorbidity data were collected from questionnaires and the VA EHR;
158 “pre-COVID” data was from the time of enrollment into the MVP to September 30, 2019. The cohort
159 demographics and a description of the clinical conditions for all tested patients in the two years
160 preceding the index dates are presented in a supplemental table (**Table E1**).

161

162 Genotyping was performed using a custom Thermo Fisher Axiom genotyping platform (MVP 1.0) which
163 included direct genotyping of rs35705950-T. Ancestry was defined using Harmonized ancestry, race, and
164 ethnicity (HARE) derived from self-report and genetic ancestry data(15).

165

166 *COVID-19 outcome definitions*

167

168 COVID-19 infection status from 02/2020 - 04/2021 was assessed by either self-report (if testing was
169 performed outside the VA) or by a positive polymerase chain reaction (PCR)-based test(16, 17). For
170 subgroup analyses of severity, only patients with confirmed PCR-based tests were examined. The *index*
171 *date* was defined as a COVID-19 diagnosis date, i.e., specimen date, or a self-reported date of diagnosis;
172 and for a hospitalized patient, the admission date up to 15 days prior to the COVID-19 case date.

173

174 Our analyses used harmonized definitions with the HGI to enable us to obtain larger sample sizes and
175 consistent results. In accordance with the HGI definitions, the three following analyses were performed:
176 (1) COVID Susceptibility: COVID-19 positive vs. population controls; (2) COVID Hospitalization-v1: COVID-
177 19 positive and hospitalized vs. population controls; (3) COVID-19 hospitalization-v2: COVID-19 positive
178 and hospitalized vs. COVID-19 positive but not hospitalized.

179

180 Our other analyses focused on data present in MVP only and addressed the outcome severity and post-
181 index events. For these sets of analyses, we only focused on patients who received their PCR-based
182 COVID-19 testing within VA systems. COVID-19 severity scale was derived from the WHO COVID-19
183 Disease Progression Scale(18) as mild, moderate (hospitalization), severe (Intensive Care Unit-level
184 care), or death within 30 days of PCR-confirmed COVID-19 infection. All data and variables were
185 assessed centrally by the MVP data core's Shared Data Repository (SDR).

186

187 *Post-index analytic constructs and study design*

188 The ICD codes used to pull the pneumonia events within 60 days post-index (pneumonia60d) are
189 presented in **Table E2**. Pre-index conditions were derived using natural language processing (NLP)-
190 boosted unstructured notes, ICD and Current Procedural Terminology (CPT) codes, and medications are

191 taken 2 years prior. Post-index conditions, including pneumonia, were derived using ICD and CPT codes,
192 and medications 60 days after the index date. Association with post-index pneumonia events
193 (pneumonia60d) were performed among patients who received confirmatory COVID-19 PCR testing at
194 VA sites.

195

196 **Statistical analysis**

197 Firth logistic regression(19, 20) as implemented in the R (v3.6.1) package “brglm2” (version 0.7.1)(21)
198 was used to examine the association between COVID-19 outcomes and rs35705950-T (additive model)
199 separately by ancestry, with adjustment for age, age², sex, and ethnicity-specific principal components.
200 Trans-ethnic meta-analyses were performed using random-effects models in “metafor” (version 2.4-
201 0)(22). Interactions between COVID-19 infection status and rs35705950-T on the outcome of post-index
202 pneumonia at 60 days were assessed using a multiplicative interaction term followed by stratified
203 analyses by COVID-19 infection status, with additional covariate adjustment of pre-index pneumonia.

204

205 **Phenome-wide and Laboratory-wide association studies (PheWAS and LabWAS)**

206 Associations between rs35705950-T allele and pre-existing comorbid conditions and laboratory values
207 were examined using preclinical data prior to the COVID-19 era (Sept 2019). Individuals with ≥ 2
208 Phecodes(23) were defined as cases. Phecodes with < 200 cases within each ancestry group were
209 excluded, resulting in 1618 (EUR), 1289 (AFR), 994 (HIS) Phecodes. LabWAS was conducted for 69 clinical
210 tests; for individuals with repeated measures, the median of the individuals’ EHR record was used.
211 Logistic/Firth regression and linear regression were used for Phecodes and laboratory measurements,
212 respectively. Bonferroni-adjusted thresholds for significance (by ancestry) were: EUR = 3.09×10^{-05}
213 ($0.05/1618$), AFR = 3.8×10^{-05} ($0.05/1289$), HIS = 5.03×10^{-05} ($0.05/994$). Analyses were performed using
214 PLINK2(24) (Additional details in supplemental methods).

215 **Meta-analysis with HGI**

216 Data from Release 5 (01/18/2021) of the COVID-19 Host Genetics Initiative (HGI) were utilized for
217 replication via an inverse-variance weighted meta-analysis using plink2a(24) and
218 GWAMA(25)(Additional details in supplemental methods).

219

220 **Results**

221 **Elucidation of the shared genetics with the MUC5B rs35705950-T by PheWAS and LabWAS**

222 In order to understand the pathophysiology associated with the *MUC5B* rs35705950-T allele, and more
223 specifically how the presence of the *MUC5B* rs35705950-T allele(s) might impact the susceptibility and
224 severity of COVID-19, we performed PheWAS and LabWAS to search for the *MUC5B* rs35705950-T allele
225 associated conditions prior to COVID-19 infection. The sample sizes for MVP participants used for
226 PheWAS and COVID-19 association studies, as well as HGI participants examined in this study, are shown
227 in **Table 1 (Figure E1)**. The results of the PheWAS are shown in **Figure 1** and **Table E3**.

228

229 In the PheWAS analysis between this *MUC5B* variant and 1605 phenotypes (cases > 200) from
230 participants of European ancestry, we found significant associations ($P_{\text{bonferroni}} < 2.5 \times 10^{-6}$) with 12
231 respiratory conditions. Consistent with the previous finding in IPF, rs35705950-T was associated with
232 increased risk of Idiopathic fibrosing alveolitis (phecode = 504.1; OR = 2.85 [2.65 - 3.05], $P = 8.90 \times 10^{-186}$),
233 other alveolar and parietoalveolar pneumonopathy (phecode = 504; OR = 2.64 [2.50 - 2.78], $P = 7.07$
234 $\times 10^{-289}$), and postinflammatory pulmonary fibrosis (phecode = 502; OR = 2.34 [2.23 - 2.45], $P = 8.90 \times 10^{-186}$).
235 Additionally, we also observed significant associations with respiratory failure (Phecode: 509),
236 ventilatory dependence (Phecode 509.8), lung transplant (Phecode: 510.2) and pneumonia (Phecode:
237 480) (**Figure 1, Table E3**). Notably, we evaluated Phecodes for influenza infection (481) in our PheWAS
238 analysis and did not observe an association with *MUC5B* rs35705950-T ($p < 0.05$; the power to detect a

239 difference was >95% as there were 4728 cases of influenza in EUR).

240

241 We identified, as in EUR, a significant association of this *MUC5B* variant with an increased risk of three

242 pulmonary conditions in African ancestry participants: idiopathic alveolitis (Phecode: 504.1), other

243 alveolar and parietoalveolar pneumonopathy (Phecode:504), and post-inflammatory fibrosis (Phecode:

244 502) (**Figure 1, Table E3**). Two of these associations, other alveolar and parietoalveolar pneumonopathy

245 (Phecode:504) and post-inflammatory fibrosis (Phecode: 502), were also seen in HIS ancestry, suggesting

246 shared etiology.

247

248 We performed a Laboratory-wide association study of the *MUC5B* rs35705950-T with median values of

249 clinical laboratory tests measured prior to the COVID-19 pandemic. We only included quantitative traits

250 with 1000 or more individuals. Among EUR participants, we evaluated 63 lab measurements and 10 had

251 a significant association with the rs35705950-T. Increased level of neutrophils (absolute count) had the

252 most significant association (beta= 0.05, $p=6.24 \times 10^{-23}$). This specific association has not been previously

253 reported. Other significant associations with increased levels were white blood cell counts, neutrophil

254 fraction, estimated glomerular filtration rate (eGFR), eosinophils (absolute count), monocytes (absolute

255 count), and platelets (**Figure 2, Table E4**). The variant had an association with reduced levels of albumin,

256 lymphocyte fraction, and creatinine (**Figure 2, Table E4**). There was no significant association with lab

257 measurements in AFR or HIS, but among HIS monocytes (absolute count) were significant (beta =0.0078,

258 $p 1.66 \times 10^{-04}$) in the same direction as in EUR.

259

260 **Association of the *MUC5B* rs35705950-T allele with the COVID-19 infection or hospitalization in the**

261 **MVP and meta-analysis with HGI**

262

263 We tested for association between *MUC5B* rs35705950-T with three COVID-19 phenotype definitions 1)
264 COVID-19 positive as cases vs all the other participants in the MVP as controls 2) COVID-19 positive that
265 required hospitalization for treatment vs all the other participants in the MVP as controls 3) COVID-19
266 positive that required hospitalization for treatment vs COVID-19 positives that didn't require
267 hospitalization as controls. First, we performed the analysis in three major ancestries separately
268 (European, African, and Hispanic). Then, we meta-analyzed the summary statistics with the COVID-19
269 HGI (Freeze 5) using an inverse-variance weighted method (GWAMA)(25). Among the three COVID-19
270 phenotypes, the most significant association of rs35705950-T allele carriers was with fewer
271 hospitalization events (OR = 0.89 [0.85-0.93], $p=1.88 \times 10^{-6}$, **Figure 3** and **Table 2**).

272

273 **Association of the *MUC5B* rs35705950-T allele with fewer pneumonia events within 60 days of COVID-**
274 **19 infection in the MVP**

275 In 9,216 COVID-19 infected MVP patients, the adjusted odds ratio for post-index pneumonia was 14.8%
276 less with each additional *MUC5B* rs35705950-T allele (OR = 0.852 [0.757-0.958], $p=0.008$). In COVID-19
277 negative patients, the adjusted odds for post-index pneumonia was 7.8% higher with each additional
278 *MUC5B* rs35705950-T allele (OR=1.078 [1.001-1.162], $p=0.048$). This differential effect of an additional
279 *MUC5B* rs35705950-T allele on post-index pneumonia in COVID-19 positive vs. COVID-19 negative
280 patients was statistically significant (p -value for interaction 0.0009) in EUR (**Table 3**, **Table E5**).

281

282 **Association of the *MUC5B* rs35705950-T allele with severe outcomes of COVID-19 infection in the**
283 **MVP**

284 Presence of a *MUC5B* rs35705950-T allele was not associated with severe outcomes of COVID-19
285 infection in the MVP. The *MUC5B* rs35705950-T allele was not associated with severe outcomes with
286 mortality (OR = 1.01 [0.58-1.20], $p=0.72$) nor mortality alone (OR = 0.91 [0.72-1.16], $p=0.25$) in EUR

287 ancestry individuals(**Table E6**).

288

289 **Discussion**

290 The data herein establishes that the “T” allele of rs35705950-T in *MUC5B*, which has been associated
291 with an *increased* risk for the development of IPF, confers a *decreased* risk of hospitalization and
292 pneumonia following COVID-19 infection among MVP participants of European ancestry. The protective
293 effect of the rs35705950-T, in addition to being counterintuitive, is in stark contrast to the increased risk
294 of severe COVID-19 disease observed for other well-established causal variants or IPF, including variants
295 located in the *TERC*, *DEPTOR*, and *FAM13A*(26).

296

297 The protein product of *MUC5B* is a major gel-forming mucin in the lung that plays a key role in
298 mucociliary clearance (MCC) and host defense(27). *MUC5B* protein is secreted from proximal
299 submucosal glands and distal airway secretory cells(28–30). Mucus traps inhaled particles, including
300 bacteria, and transporting them out of the airways by ciliary and cough-driven forces. Mucin also helps
301 remove endogenous debris including dying epithelial cells and leukocytes. *MUC5AC* and *MUC5B* are two
302 major secreted forms of mucins in the lung.

303

304 The rs35705950-T is located within an enhancer region of *MUC5B*; the “T” allele demonstrates gain-of-
305 function and is associated with enhanced expression of the *MUC5B* transcript in lung tissue from
306 unaffected subjects and patients with IPF(31). In patients with IPF, excess *MUC5B* protein is especially
307 observed in epithelial cells in the respiratory bronchiole and honeycomb cyst(29, 30, 32), regions of the
308 lung involved in lung fibrosis.

309

310 Mouse models found that *Muc5b* is required for mucociliary clearance, for controlling bacterial

311 infections in the airways and middle ear, and for maintaining immune homeostasis in mouse lungs(33).
312 *Muc5b* deficiency caused materials to accumulate in the upper and lower airways. This defect led to
313 chronic infection by multiple bacterial species, including *Staphylococcus aureus*, and to inflammation
314 that failed to resolve normally. Apoptotic macrophages accumulated, phagocytosis was impaired, and
315 interleukin-23 (IL-23) production was reduced in *Muc5b*(-/-) mice. By contrast, in transgenic mice that
316 overexpress *Muc5b*, macrophage functions improved (33). *Muc5B* over-expressing transgenic mice have
317 been shown to be more susceptible to the fibroproliferative effects of bleomycin (34), consistent with a
318 role in IPF. Paradoxically, while the “T” allele of rs35705950-T increases susceptibility towards the
319 development of IPF, the same allele has also been associated with *decreased* mortality among IPF
320 patients(35).

321

322 Our analyses demonstrating a significant interaction between COVID-19 infection and the prospective
323 development of pneumonia suggest a possible mechanism by which the protective effect of
324 rs35705950-T is mediated. Whether enhanced pulmonary macrophage function or quantitative or
325 qualitative changes in mucous production resulting from the minor allele of rs35705950-T are
326 responsible for the observed protective effect should be explored in future work. Of note, the *MUC5B*
327 rs35705950-T allele did not decrease the risk of pneumonia in COVID-19 tested negative participants
328 (**Table 3**), suggesting that the protective effect may be specific to COVID-19 related pneumonia. More
329 studies in the future are needed to further investigate this phenomenon.

330

331 No extrapulmonary association was noted on PheWAS analysis suggesting a very circumscribed
332 molecular and clinical effect of this promoter variant. This supports the notion that the effect of
333 rs35705950-T on COVID-19 infection is mediated in pulmonary tissues. The *Muc5b* over-expression in
334 the distal airway may specifically or non-specifically affect the SARS-CoV-2 viral infection in the lung,

335 leading to decreased incidence of pneumonia and hospitalization in the infected individuals.

336

337 The human *MUC5B* rs35705950-T allele does not appear to be sufficient to cause pulmonary fibrosis.

338 Although ~20% of the non-Hispanic white populations have a copy of the *MUC5B* rs35705950-T

339 allele(31, 33), IPF is a rare disease with a population prevalence of less than 0.1% (36). Additional

340 genetic and/or environmental insults are likely needed in the development of IPF in humans. Since the

341 overwhelming number of individuals with the *MUC5B* rs35705950-T allele will not know their *MUC5B*

342 status, it is unlikely that the reason for our observation is due to a change in health behaviors of

343 participants that carry this variant.

344

345 The *MUC5B* rs35705950-T allele was associated with elevated neutrophil counts. This could be due in

346 part to the association of this allele with an increased incidence of pneumonia. It is worth noting that

347 neutrophils are a major source of alpha-defensin and elevated alpha-defensin levels were seen in the

348 serum of IPF patients; the levels of alpha-defensin in the serum correlated with the lung function decline

349 in the IPF patients(37, 38).

350

351 Longer follow-up of SARS-CoV-2 infected individuals with the *MUC5B* rs35705950-T allele is needed. One

352 would need to be cautious regarding the longer-term outcome of COVID-19 in the *MUC5B* rs35705950-T

353 allele positive individuals as a fibrotic response has been reported in the survivors of severe COVID-19.

354 This is of particular importance if the manipulation of *MUC5B* expression is considered in the

355 prevention/treatment of COVID-19.

356

357 The *MUC5B* rs35705950-T allele variant resides within an enhancer subject to lineage- and disease-

358 dependent epigenetic remodeling. It was postulated that this G to T transversion in the *MUC5B*

359 rs35705950-T allele might lead to the removal of a binding site for the GCF transcription repressor(12,
360 39). A potential avenue for chromatin-based therapies in which *MUC5B* enhancer chromatin
361 architecture serves as a target to block the *MUC5B* mis-expression was proposed(12, 39). Additional
362 small molecule and signaling inhibitors targeting IPF are being studied as well(40). These strategies are
363 generally aiming at reducing fibrosis or the effects associated with *MUC5B* over-expression. How these
364 strategies or alternatives can be utilized to treat/prevent COVID-19 remains to be studied.

365

366 In conclusion, we show in this study a common *MUC5B* promoter variant leading to *MUC5B* over-
367 expression is associated with fewer hospitalizations and pneumonia events after SARS-CoV-2 infection.
368 Our study provides a strong rationale to stratify patient populations based on common and disease-
369 related genetic polymorphism in order to better understand the mechanisms and their clinical
370 implications in COVID-19. How the *MUC5B* rs35705950-T allele association may shed light on the
371 pathogenesis and/or management of COVID-19 remains to be fully examined.

372

373 **Strengths & Limitations**

374

375 MVP is a large genomic medicine database with diverse ethnicity and geography. MVP participants are
376 predominantly males but it represents a large multi-ethnic, prospective cohort available. Successful
377 replication in the HGI and meta-analysis is a strength as well as our ability to investigate specific clinical
378 events post-index. PheWAS was designed as a broad screen to test for potentially clinically relevant
379 associations between genes and phenotypes and helped in the understanding of potential disease
380 mechanisms but has limited power to detect associations among uncommon conditions, especially
381 when further stratified by genetic ancestry.

382

383 **Acknowledgments**

384 This research is based on data from the Million Veteran Program, Office of Research and Development,
385 Veterans Health Administration, and was supported by award MVP035. This publication does not
386 represent the views of the Department of Veteran Affairs of the United States Government.

387

388 We are grateful to our Veterans for their contribution to MVP. Full acknowledgments for the VA Million
389 Veteran Program COVID-19 Science Initiative can be found in the supplemental methods.

390

391 **Conflict of Interest**

392 CJO is an employee of Novartis Institute for Biomedical Research. PN reports grant support from Amgen,
393 Apple, AstraZeneca, Boston Scientific, and Novartis, personal fees from Apple, AstraZeneca, Blackstone
394 Life Sciences, Genentech, and Novartis, and spousal employment at Vertex, all unrelated to the present
395 work.

396

397 **References**

- 398 1. Website. at <(https://coronavirus.jhu.edu/map.html)>.
- 399 2. Zhao Y-M, Shang Y-M, Song W-B, Li Q-Q, Xie H, Xu Q-F, Jia J-L, Li L-M, Mao H-L, Zhou X-M, Luo H,
400 Gao Y-F, Xu A-G. Follow-up study of the pulmonary function and related physiological
401 characteristics of COVID-19 survivors three months after recovery. *EClinicalMedicine* 2020;
- 402 3. Yan X, Huang H, Wang C, Jin Z, Zhang Z, He J, Yin S, Fan M, Huang J, Chen F, Zeng Y, Han X, Zhu Y.
403 Follow-up study of pulmonary function among COVID-19 survivors 1 year after recovery. *J Infect*
404 2021;doi:10.1016/j.jinf.2021.05.034.
- 405 4. McGroder CF, Zhang D, Choudhury MA, Salvatore MM, D'Souza BM, Hoffman EA, Wei Y, Baldwin

- 406 MR, Garcia CK. Pulmonary fibrosis 4 months after COVID-19 is associated with severity of illness
407 and blood leucocyte telomere length. *Thorax* 2021;doi:10.1136/thoraxjnl-2021-217031.
- 408 5. King CS, Nathan SD. Idiopathic pulmonary fibrosis: effects and optimal management of
409 comorbidities. *Lancet Respir Med* 2017;5:72–84.
- 410 6. George PM, Patterson CM, Reed AK, Thillai M. Lung transplantation for idiopathic pulmonary
411 fibrosis. *Lancet Respir Med* 2019;7:271–282.
- 412 7. Ley B, Collard HR, King TE Jr. Clinical course and prediction of survival in idiopathic pulmonary
413 fibrosis. *Am J Respir Crit Care Med* 2011;183:431–440.
- 414 8. Selman M, Pardo A. Revealing the pathogenic and aging-related mechanisms of the enigmatic
415 idiopathic pulmonary fibrosis. an integral model. *Am J Respir Crit Care Med* 2014;189:1161–1172.
- 416 9. Selman M, Pardo A. The leading role of epithelial cells in the pathogenesis of idiopathic
417 pulmonary fibrosis. *Cell Signal* 2020;66:109482.
- 418 10. Zhang Y, Noth I, Gibson KF, Ma S-F, Richards TJ, Bon JM, Lindell KO, Branch RA, Nicolae D,
419 Sciruba F, Garcia AN, Kaminski N. A Variant In The Promoter Of MUC5B Is Associated With
420 Idiopathic Pulmonary Fibrosis And Not Chronic Obstructive Pulmonary Disease. *B102 INTERSTITIAL*
421 *LUNG DISEASE: NOVEL MANAGEMENT AND OUTCOME STRATEGIES* 2011;doi:10.1164/ajrccm-
422 conference.2011.183.1_meetingabstracts.a6395.
- 423 11. Moore C, Blumhagen RZ, Yang IV, Walts A, Powers J, Walker T, Bishop M, Russell P, Vestal B,
424 Cardwell J, Markin CR, Mathai SK, Schwarz MI, Steele MP, Lee J, Brown KK, Loyd JE, Crapo JD,
425 Silverman EK, Cho MH, James JA, Guthridge JM, Cogan JD, Kropski JA, Swigris JJ, Bair C, Kim DS, Ji
426 W, Kim H, *et al.* Resequencing Study Confirms That Host Defense and Cell Senescence Gene
427 Variants Contribute to the Risk of Idiopathic Pulmonary Fibrosis. *Am J Respir Crit Care Med*
428 2019;200:199–208.
- 429 12. Gally F, Sasse SK, Kurche JS, Gruca MA, Cardwell JH, Okamoto T, Chu HW, Hou X, Poirion OB,

430 Buchanan J, Preissl S, Ren B, Colgan SP, Dowell RD, Yang IV, Schwartz DA, Gerber AN. The MUC5B-
431 associated variant rs35705950 resides within an enhancer subject to lineage- and disease-
432 dependent epigenetic remodeling. *JCI Insight* 2021;6.:

433 13. Gaziano JM, Concato J, Brophy M, Fiore L, Pyarajan S, Breeling J, Whitbourne S, Deen J, Shannon
434 C, Humphries D, Guarino P, Aslan M, Anderson D, LaFleur R, Hammond T, Schaa K, Moser J, Huang
435 G, Muralidhar S, Przygodzki R, O’Leary TJ. Million Veteran Program: A mega-biobank to study
436 genetic influences on health and disease. *Journal of Clinical Epidemiology* 2016;

437 14. COVID-19 Host Genetics Initiative. Mapping the human genetic architecture of COVID-19.
438 *Nature* 2021;doi:10.1038/s41586-021-03767-x.

439 15. Fang H, Hui Q, Lynch J, Honerlaw J, Assimes TL, Huang J, Vujkovic M, Damrauer SM, Pyarajan S,
440 Gaziano JM, DuVall SL, O’Donnell CJ, Cho K, Chang K-M, Wilson PWF, Tsao PS, VA Million Veteran
441 Program, Sun YV, Tang H. Harmonizing Genetic Ancestry and Self-identified Race/Ethnicity in
442 Genome-wide Association Studies. *Am J Hum Genet* 2019;105:763–772.

443 16. Ioannou GN, Locke E, Green P, Berry K, O’Hare AM, Shah JA, Crothers K, Eastment MC, Dominitz
444 JA, Fan VS. Risk Factors for Hospitalization, Mechanical Ventilation, or Death Among 10 131 US
445 Veterans With SARS-CoV-2 Infection. *JAMA Netw Open* 2020;3:e2022310.

446 17. Obeidat M, Frank AR, Icardi MS, Klutts JS. VA-Wide, Multicenter Verification Study of the
447 Cepheid Xpert SARS-CoV-2 Assay. *Acad Pathol* 2021;8:23742895211011911.

448 18. WHO Working Group on the Clinical Characterisation and Management of COVID-19 infection. A
449 minimal common outcome measure set for COVID-19 clinical research. *Lancet Infect Dis*
450 2020;20:e192–e197.

451 19. Firth D. Bias reduction of maximum likelihood estimates. *Biometrika* 1993;

452 20. Kosmidis I, Firth D. Jeffreys-prior penalty, finiteness and shrinkage in binomial-response
453 generalized linear models. *Biometrika* 2021;

- 454 21.Kosmidis I. brglm2: Bias reduction in generalized linear models. *R package version 0.1* 2017;5.:
- 455 22.Viechtbauer W. Conducting Meta-Analyses in R with the metafor Package. *Journal of Statistical*
- 456 *Software* 2010;
- 457 23.Denny JC, Ritchie MD, Basford MA, Pulley JM, Bastarache L, Brown-Gentry K, Wang D, Masys DR,
- 458 Roden DM, Crawford DC. PheWAS: demonstrating the feasibility of a phenome-wide scan to
- 459 discover gene-disease associations. *Bioinformatics* 2010;26:1205–1210.
- 460 24.Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: rising to
- 461 the challenge of larger and richer datasets. *GigaScience* 2015;
- 462 25.Mägi R, Morris AP. GWAMA: software for genome-wide association meta-analysis. *BMC*
- 463 *Bioinformatics* 2010;11:288.
- 464 26.Fadista J, Kraven LM, Karjalainen J, Andrews SJ, Geller F, COVID-19 Host Genetics Initiative,
- 465 Baillie JK, Wain LV, Jenkins RG, Feenstra B. Shared genetic etiology between idiopathic pulmonary
- 466 fibrosis and COVID-19 severity. *EBioMedicine* 2021;65:103277.
- 467 27.Evans CM, Fingerlin TE, Schwarz MI, Lynch D, Kurche J, Warg L, Yang IV, Schwartz DA. Idiopathic
- 468 Pulmonary Fibrosis: A Genetic Disease That Involves Mucociliary Dysfunction of the Peripheral
- 469 Airways. *Physiol Rev* 2016;96:1567–1591.
- 470 28.Rose MC, Voynow JA. Respiratory tract mucin genes and mucin glycoproteins in health and
- 471 disease. *Physiol Rev* 2006;86:245–278.
- 472 29.Seibold MA, Smith RW, Urbanek C, Groshong SD, Cosgrove GP, Brown KK, Schwarz MI, Schwartz
- 473 DA, Reynolds SD. The idiopathic pulmonary fibrosis honeycomb cyst contains a mucociliary
- 474 pseudostratified epithelium. *PLoS One* 2013;8:e58658.
- 475 30.Nakano Y, Yang IV, Walts AD, Watson AM, Helling BA, Fletcher AA, Lara AR, Schwarz MI, Evans
- 476 CM, Schwartz DA. MUC5B Promoter Variant rs35705950 Affects MUC5B Expression in the Distal
- 477 Airways in Idiopathic Pulmonary Fibrosis. *American Journal of Respiratory and Critical Care*

478 *Medicine* 2016;

479 31.Seibold MA, Wise AL, Speer MC, Steele MP, Brown KK, Loyd JE, Fingerlin TE, Zhang W,
480 Gudmundsson G, Groshong SD, Evans CM, Garantziotis S, Adler KB, Dickey BF, du Bois RM, Yang IV,
481 Herron A, Kervitsky D, Talbert JL, Markin C, Park J, Crews AL, Slifer SH, Auerbach S, Roy MG, Lin J,
482 Hennessy CE, Schwarz MI, Schwartz DA. A common MUC5B promoter polymorphism and
483 pulmonary fibrosis. *N Engl J Med* 2011;364:1503–1512.

484 32.Conti C, Montero-Fernandez A, Borg E, Osadolor T, Viola P, De Lauretis A, Stock CJ, Bonifazi M,
485 Bonini M, Caramori G, Lindahl G, Blasi FB, Nicholson AG, Wells AU, Sestini P, Renzoni E. Mucins
486 MUC5B and MUC5AC in Distal Airways and Honeycomb Spaces: Comparison among Idiopathic
487 Pulmonary Fibrosis/Usual Interstitial Pneumonia, Fibrotic Nonspecific Interstitial Pneumonitis, and
488 Control Lungs. *Am J Respir Crit Care Med* 2016;193:462–464.

489 33.Roy MG, Livraghi-Butrico A, Fletcher AA, McElwee MM, Evans SE, Boerner RM, Alexander SN,
490 Bellinghausen LK, Song AS, Petrova YM, Tuvim MJ, Adachi R, Romo I, Bordt AS, Bowden MG, Sisson
491 JH, Woodruff PG, Thornton DJ, Rousseau K, De la Garza MM, Moghaddam SJ, Karmouty-Quintana
492 H, Blackburn MR, Drouin SM, Davis CW, Terrell KA, Grubb BR, O’Neal WK, Flores SC, *et al.* Muc5b is
493 required for airway defence. *Nature* 2014;505:412–416.

494 34.Hancock LA, Hennessy CE, Solomon GM, Dobrinskikh E, Estrella A, Hara N, Hill DB, Kissner WJ,
495 Markovetz MR, Grove Villalon DE, Voss ME, Tearney GJ, Carroll KS, Shi Y, Schwarz MI, Thelin WR,
496 Rowe SM, Yang IV, Evans CM, Schwartz DA. Muc5b overexpression causes mucociliary dysfunction
497 and enhances lung fibrosis in mice. *Nat Commun* 2018;9:5363.

498 35.Peljto AL, Zhang Y, Fingerlin TE, Ma S-F, Garcia JGN, Richards TJ, Silveira LJ, Lindell KO, Steele
499 MP, Loyd JE, Gibson KF, Seibold MA, Brown KK, Talbert JL, Markin C, Kossen K, Seiwert SD, Murphy
500 E, Noth I, Schwarz MI, Kaminski N, Schwartz DA. Association between the MUC5B promoter
501 polymorphism and survival in patients with idiopathic pulmonary fibrosis. *JAMA* 2013;309:2232–

- 502 2239.
- 503 36.Raghu G, Remy-Jardin M, Myers JL, Richeldi L, Ryerson CJ, Lederer DJ, Behr J, Cottin V, Danoff
- 504 SK, Morell F, Flaherty KR, Wells A, Martinez FJ, Azuma A, Bice TJ, Bouros D, Brown KK, Collard HR,
- 505 Duggal A, Galvin L, Inoue Y, Jenkins RG, Johkoh T, Kazerooni EA, Kitaichi M, Knight SL, Mansour G,
- 506 Nicholson AG, Pipavath SNJ, *et al.* Diagnosis of Idiopathic Pulmonary Fibrosis. An Official
- 507 ATS/ERS/JRS/ALAT Clinical Practice Guideline. *Am J Respir Crit Care Med* 2018;198:e44–e68.
- 508 37.Mukae H, Iiboshi H, Nakazato M, Hiratsuka T, Tokojima M, Abe K, Ashitani J, Kadota J, Matsukura
- 509 S, Kohno S. Raised plasma concentrations of alpha-defensins in patients with idiopathic pulmonary
- 510 fibrosis. *Thorax* 2002;57:623–628.
- 511 38.Konishi K, Gibson KF, Lindell KO, Richards TJ, Zhang Y, Dhir R, Bisceglia M, Gilbert S, Yousem SA,
- 512 Song JW, Kim DS, Kaminski N. Gene expression profiles of acute exacerbations of idiopathic
- 513 pulmonary fibrosis. *Am J Respir Crit Care Med* 2009;180:167–175.
- 514 39.Korfei M, Stelmaszek D, MacKenzie B, Skwarna S, Chillappagari S, Bach AC, Ruppert C, Saito S,
- 515 Mahavadi P, Klepetko W, Fink L, Seeger W, Lasky JA, Pullamsetti SS, Krämer OH, Guenther A.
- 516 Comparison of the antifibrotic effects of the pan-histone deacetylase-inhibitor panobinostat versus
- 517 the IPF-drug pirfenidone in fibroblasts from patients with idiopathic pulmonary fibrosis. *PLoS One*
- 518 2018;13:e0207915.
- 519 40.Montesi SB, Fisher JH, Martinez FJ, Selman M, Pardo A, Johannson KA. Update in Interstitial Lung
- 520 Disease 2019. *Am J Respir Crit Care Med* 2020;202:500–507.

521

522

523

524

525

526 **Figures and Tables**

527 **Table 1.** Demographics for COVID-19 tested positive and all MVP participants examined in this study.

Characteristics	Million Veteran Program	COVID-19 Positive
	Number (%)	Number (%)
Total Patients	658,582	13,841
Male	592516 (90)	12,320 (89)
Genetic Ancestry		
European	464961 (70)	8011 (58)
African	123120 (19)	3749 (27)
Hispanic	52183 (8)	1903 (14)
Asian	8329 (1)	178 (1)
Other	9989 (2)	0
Muc5B rs35705950		
0 copy	10604 (1.6)	353 (25)
1 copy	2161 (0.03)	75 (0.05)
2 copies		
Comorbidities		
Obesity (phecode = 278)	283197 (43)	8905 (64)
Hypertension (phecode = 401.1)	451998 (69)	10617 (77)
Type 2 Diabetes (phecode = 250.2)	227575 (34)	10491 (76)
Coronary Artery Disease (phecode = 411.4)	152136 (23)	4182 (30)
Chronic Kidney Disease (phecode = 585.2)	10046 (15)	533 (38)
Outcomes		
Hospitalized	-	4491 (32)
Severe	-	657 (47)
Deceased	-	644 (46)

528

529

530

531

532

533

534

535

536 **Table 2.** Association of rs35705950 in *MUC5B* with (i) COVID-19 Positive vs Population Controls, (ii)
 537 COVID-19 Positive, Hospitalized vs Population Controls, and (iii) COVID-19 Positive, Hospitalized vs
 538 COVID-19 Positive, not Hospitalized. Odds ratio (OR) and 95% confidence interval (95% CI) is reported
 539 for the minor (T) allele, and results are shown for VA Million Veteran Program (MVP) African Americans
 540 (AFR), European Americans (EUR), Hispanic/Latino Americans (HIS), and trans-ethnic meta-analysis (ALL),
 541 the COVID-19 Host Genetics Initiative (HGI) trans-ethnic release 5 meta-analysis excluding MVP and
 542 23&Me, and the meta-analysis of MVP and HGI (META).

Analysis	Population	N Case	N Control	Total N	EAF	OR (95% CI)	P
Positive vs Population Control	MVP (AFR)	6,411	114,781	121,192	0.02	0.99 [0.87, 1.12]	0.826
	MVP (EUR)	15,814	443,428	459,242	0.11	0.96 [0.92, 0.99]	0.019
	MVP (HIS)	3,128	47,462	50,590	0.07	0.95 [0.85, 1.05]	0.275
	MVP (ALL)	25,353	605,671	631,024	0.09	0.96 [0.93, 1.00]	0.060
	HGI (ALL)	25,652	1,282,972	1,308,624	0.11	0.98 [0.95, 1.01]	0.134
	META	51,005	1,888,643	1,939,648	0.10	0.97 [0.95, 0.99]	4.57E-03
Hospitalized vs Population Control	MVP (AFR)	1,739	119,453	121,192	0.02	0.83 [0.64, 1.07]	0.147
	MVP (EUR)	3,325	455,917	459,242	0.11	0.87 [0.80, 0.94]	5.43E-04
	MVP (HIS)	657	49,933	50,590	0.07	0.86 [0.68, 1.07]	0.182
	MVP (ALL)	5,721	625,303	631,024	0.09	0.86 [0.80, 0.93]	7.35E-05
	HGI (ALL)	9,086	1,001,201	1,010,287	0.11	0.91 [0.85, 0.97]	4.12E-03
	META	14,807	1,626,504	1,641,311	0.10	0.89 [0.85, 0.93]	1.88E-06
Hospitalized vs Not Hospitalized	MVP (AFR)	1,739	4,672	6,411	0.02	0.80 [0.59, 1.08]	0.141
	MVP (EUR)	3,325	12,489	15,814	0.11	0.89 [0.81, 0.97]	0.012
	MVP (HIS)	657	2,471	3,128	0.07	0.88 [0.68, 1.14]	0.319
	MVP (ALL)	5,721	19,632	25,353	0.08	0.88 [0.81, 0.96]	2.64E-03
	HGI (ALL)	4,420	11,093	15,513	0.16	0.97 [0.88, 1.08]	0.575
	META	10,141	30,725	40,866	0.11	0.91 [0.86, 0.98]	7.20E-03

543

544

545 **Table 3.** Fewer pneumonia events developed within 60 days post COVID-19 infection for MVP EUR
 546 individuals with the presence of a *MUC5B* rs35705950 allele. Odds ratios are estimated from Firth
 547 logistic regression adjusting for pre-index pneumonia, age, age², and PC1-20, including an interaction
 548 between additive *MUC5B* rs35705950 allele and COVID-19 infection.

549
 550

	COVID-19 negative	COVID-19 positive	COVID-19 & <i>MUC5B</i> p-value for interaction
	OR (95% CI) of a <i>MUC5B</i> allele		p=0.0009
	1.08 (1.00, 1.16) p=0.048	0.89 (0.76, 0.96) p=0.008	
<i>MUC5B</i>=0	<i>MUC5B</i>=1	<i>MUC5B</i>=2	
OR (95% CI) of COVID-19 positive status			
10.00 (9.35, 10.70) p<0.0001	7.91 (6.97, 8.97) p<0.0001	6.26 (4.83, 8.08) p<0.0001	

551
 552
 553

554 **Figure 1. Phenome-Wide Association Study (PheWAS) of *MUC5B* rs35705950 allele in the Million**
 555 **Veteran Program.** A PheWAS plot shows associations of rs35705950 and phenotypes derived from the
 556 electronic health records data prior to COVID-19 in MVP participants from A) European ancestry B)
 557 African ancestry and C) Hispanic ancestry. The phenotypes are shown on the x-axis and organized by
 558 disease categories. The p-value (-log₁₀) of each association is shown on the y-axis the direction of the
 559 triangle represents the direction of effect of the associations - with the upward triangle as increased risk
 560 and the downward triangle as reduced risk. The red line indicates the significance threshold based on
 561 the Bonferroni correction. The forest plot of Bonferroni significant associations are shown within the
 562 right top corner of each PheWAS plot. The Bonferroni threshold for each ancestry group is shown in the
 563 forest plot.

564

565 **Figure 2. Laboratory-Wide Association Study (PheWAS) of *MUC5B* rs35705950 allele in the Million**
 566 **Veteran Program.** A LabWAS plot shows associations of rs35705950 and median values of laboratory
 567 measures extracted from electronic health records data prior to COVID-19 in MVP participants. The
 568 bottom panel shows the -log₁₀ (p-value) on the y-axis and laboratory test descriptions on the x-axis.
 569 Triangles points up have increasing effects and points down have decreasing effects. The colors
 570 represent the different ancestry groups. The top panel shows beta from the regression model for each
 571 laboratory measure. The significant results are highlighted in the color corresponding to ancestry groups
 572 and other results are plotted in grey.

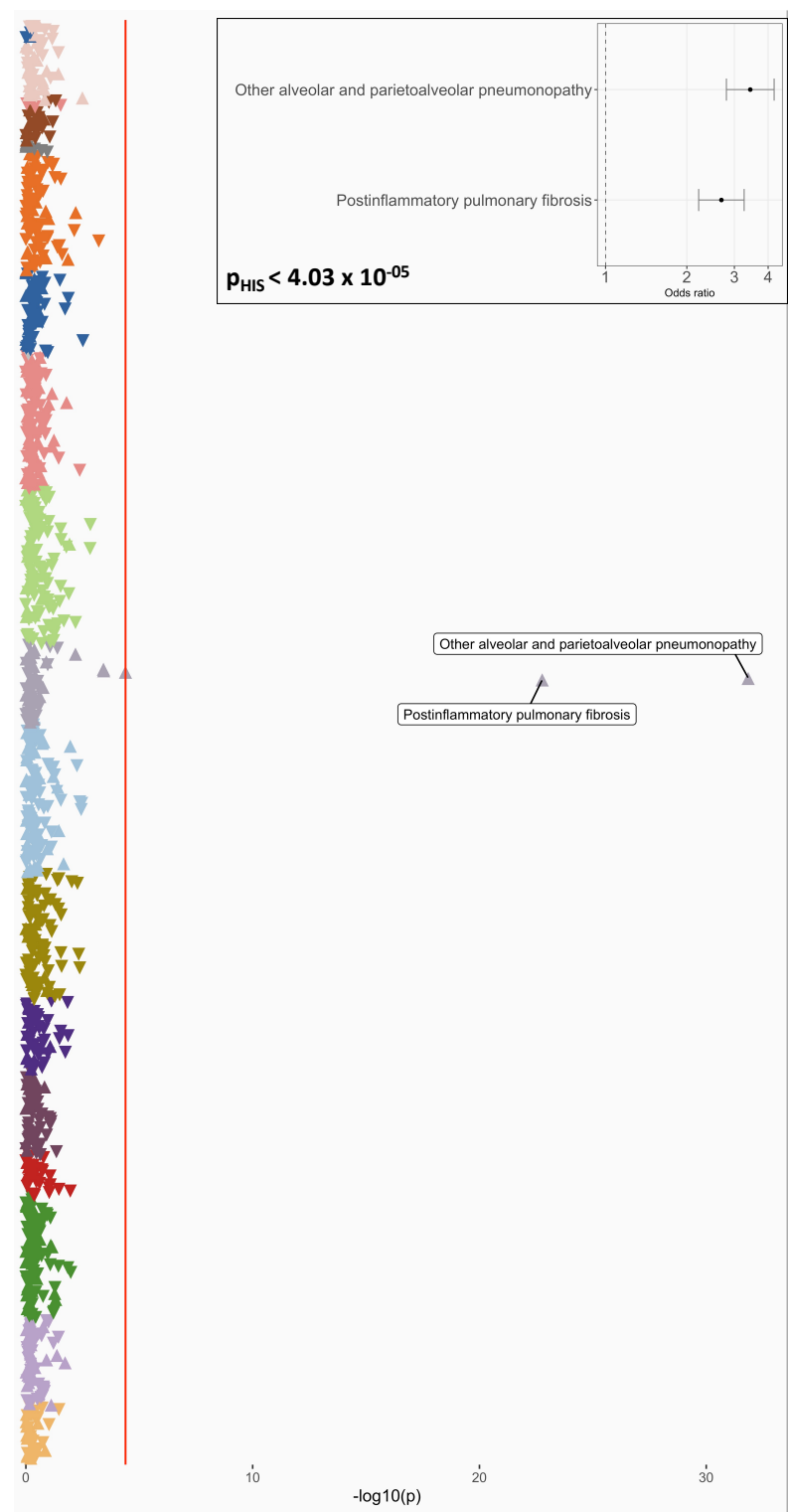
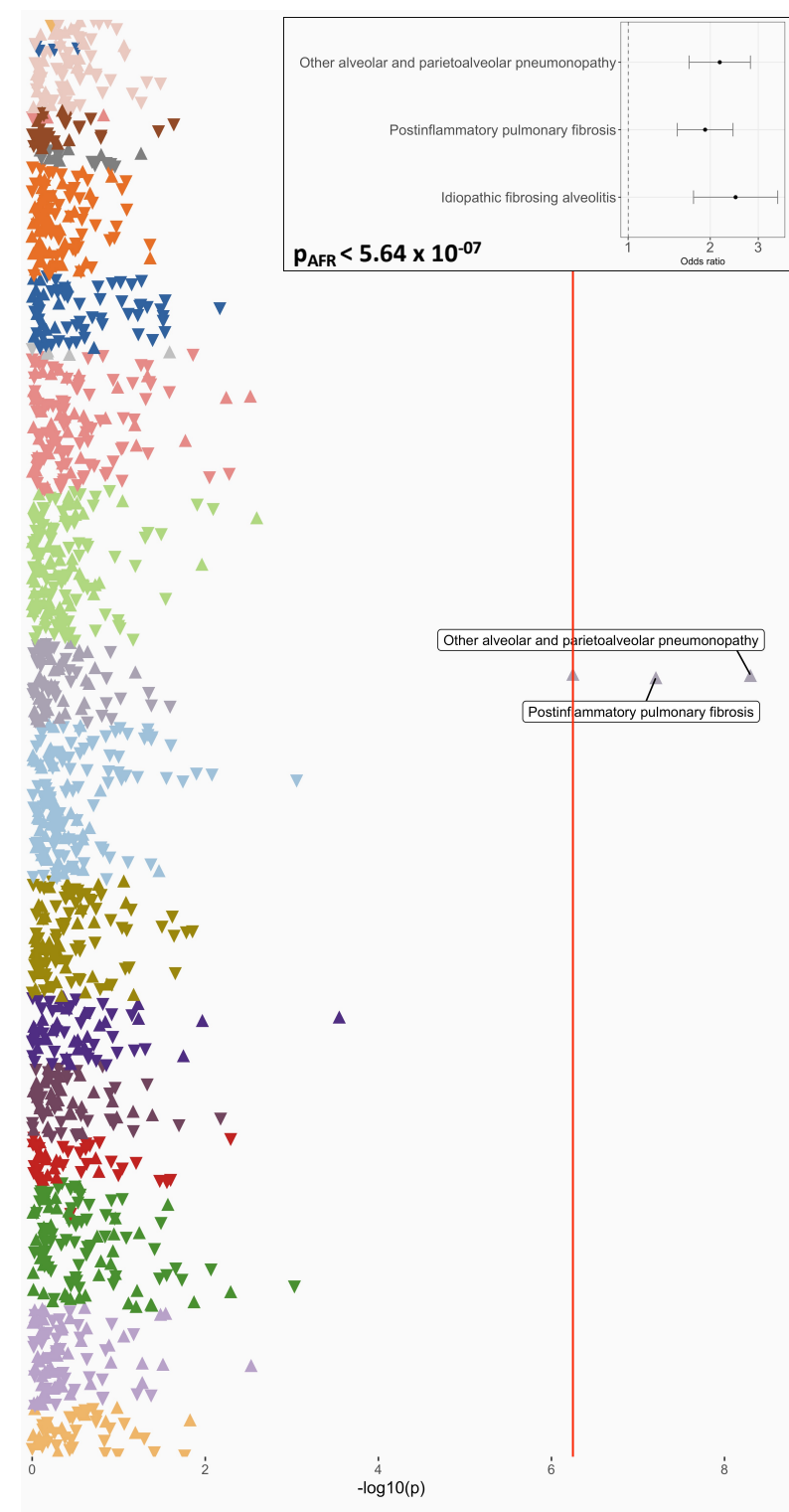
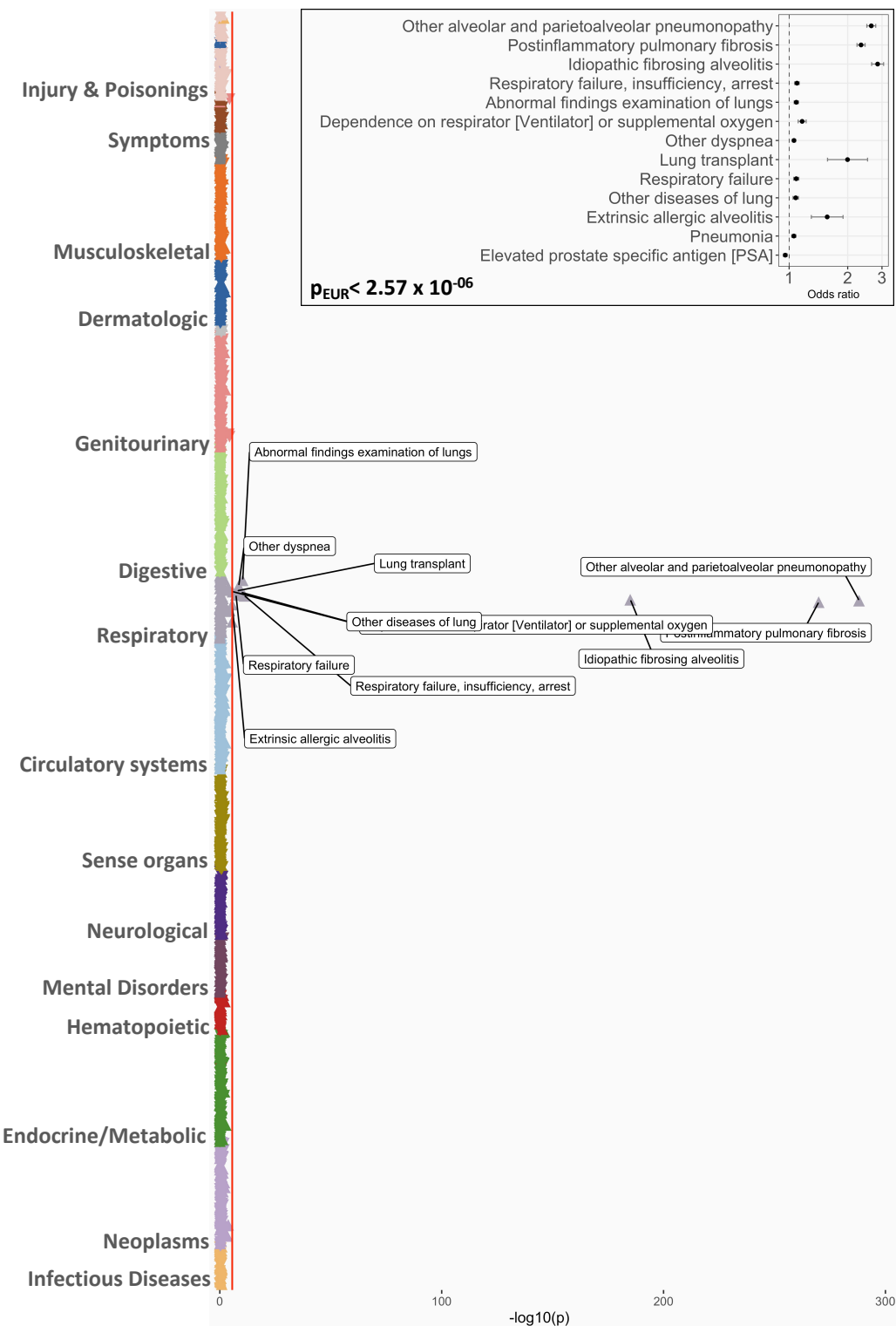
573

574 **Figure 3.** Forest plot association of rs35705950 in *MUC5B* with (i) COVID-19 Positive vs Population
 575 Controls, (ii) COVID-19 Positive, Hospitalized vs Population Controls, and (iii) COVID-19 Positive,

576 Hospitalized vs COVID-19 Positive, not Hospitalized. Odds ratio (OR) and 95% confidence interval (95%
577 CI) is reported for the minor (T) allele, and results are shown for VA Million Veteran Program (MVP)
578 African Americans (AFR), European Americans (EUR), Hispanic/Latino Americans (HIS), and trans-ethnic
579 meta-analysis (ALL), the COVID-19 Host Genetics Initiative (HGI) trans-ethnic release 5 meta-analysis
580 excluding MVP and 23&Me, and the meta-analysis of MVP and HGI (META).

581

582



Direction of effect

▲ Risk

▼ Protective

A

B

C

