1    **Additional Material**

2

3    **Fine-Scale Frequency of the *MUC5B* Promoter Variant Correlates with**
4    **Idiopathic Pulmonary Fibrosis Healthcare Burden**
5
6    Edmund Gilbert[1,2], Aoife Carolan[3,4], Mari Ozaki[3,4], Wan Lin Ng[3,5], Jisha Jasmin[3,6], David

7    A Schwartz[7], Michael T Henry[8], Gianpiero L. Cavalleri[1,2], Killian Hurley[3,4], Irish Pulmonary

8    Fibrosis Research Consortium*

9    **Methods**

10   **Genotype Data**

11   Using PLINK (1, 2)], we merged the Irish ancestry, genome-wide, SNP genotypes

12   previously reported from the Irish DNA Atlas (3)[REF] and the Trinity Student (4)

13   (dbGaP Study Accession: phs000789.v1.p1) datasets. New genotypes from the Irish

14   DNA Atlas were added to this dataset, using the same genotyping procedure as

15   previously described (3). Performing quality control of these genotypes, we removed

16   individuals with a missingness >5%, then SNP with a missingness >5%, minor allele

17   frequency <1%, and a p-value denoting deviation from Hardy-Weinberg expectations <

18   1e-9. To remove close relatedness from the dataset and thereby capture population-

19   level relatedness and diversity, we further removed one from each detected close family

20   pairs, using KING (5), prioritising participants from the Irish DNA Atlas due to their

21   added geographic data. This left a final dataset of 2,604 individuals (Irish DNA Atlas

22   n=381, Trinity Student n=2,223).

23   To cluster these 2,604 individuals, we phased autosomal genotypes using Beagle v5.4

24   (6) detected Identity-by-Descent segments with hap-ibd (7). For each pair of individuals,

25   we summed the total length of IBD segments with an individual length of > 3 cM and <

26   30 cM, and then constructed a network object from these summed lengths in python

27   and using the python implementation of the igraph package. We then performed

28   community detection using the Leiden algorithm (8) using the python implementation in

29   the leidenalg package, over two hierarchical levels of clustering. We detected the first

30   "top" level of clusters initially, and then further sub-divided each of the four top level

31   clusters into a second level of sub-clusters. Roughly, the top level corresponded to

32   Provincial differences, and the second level corresponded to County differences. We

33   disregarded small second level clusters that were too small or contained too few Irish

34   DNA Atlas individuals to annotate. This left a total of 2,465 individuals over 19 individual

35   second level clusters, each grouped into one of 4 top level clusters.

36   The rs35705950 variant was not directly genotyped in either Irish reference dataset,

37   therefore we sought to impute the genotype identity for each of our Irish Region

38   Reference samples. Using the Michigan Imputation Server (9) and the Haplotype

39   Reference Consortium panel (v2016), we imputed genotypes on chromosome 11.

40   rs35705950 imputation was of high quality ($r^2$ = 0.914), matching the relatively high

41   frequency of the variant overall. From these imputed genotypes, we estimated the allele

42   frequencies in each second-level cluster using the PLINK functions --freq and --within.

43   We then calculated confidence levels through a bootstrapping procedure. Briefly, for

44   each bootstrap iteration, we randomly shuffled cluster assignments to Irish references –

45   retaining original cluster sample size. Then we re-calculated the allele frequency for this

46   pseudo-cluster. From the allele frequencies over 100 bootstrap replicates, we estimated

47   the standard deviation and error for each cluster. Finally, to interpolate the allele

48   frequencies of the rs35705950-T allele, we utilised Kriging in R using the Krig() function

49   from the R package "fields". As an input dataset, for each Irish DNA Atlas sample with

50   geographic data, we assigned the rs35705950-T allele frequency based on the

51   frequency estimated for the cluster that they were assigned to.

52   **Discharge rates for IPF**

53   We estimated the country wide health care use burden for IPF by extracting the

54   discharge rates associated with the diagnosis of IPF (J841, ICD-10-AM/ACHI/ACS Eight

55    Edition https://www.ihacpa.gov.au/resources/icd-10-amachiacs-eighth-edition) using

56    data from the Hospital Inpatient Enquiry (HIPE) supplied by The Healthcare Pricing

57    Office, Ireland. To avoid changes in hospital admission and discharge patterns during

58    the COVID-19 pandemic we examined hospital discharges from 2015 to 2019. We

59    associated search of these discharges with the home county of the patient to estimate

60    county-based burden. To account for population size differences in Irish regions, we

61    also extracted county population sizes from the Irish 2016 census recorded by the Irish

62    Central Statistics Office (CSO), selecting the population size of individuals self-

63    identifying as "White Irish" to best proxy European-Irish ancestry. From these county-

64    based discharge rates and population sizes we then estimated the proportion that each

65    county contributes to the total number of discharges, and separately proportional

66    population size. From these, we then estimated the county discharge rate weighted by

67    population size by dividing the county's discharge proportion by the county's population

68    size proportion. Then, to estimate each county's average allele frequency, we assigned

69    each Irish DNA Atlas individual to the Irish county that their ancestral geographic

70    position resides in, as well as recording the rs35705950-T allele frequency of that

71    individual's cluster. Then for each county, we took the mean cluster rs35705950-T allele

72    frequency using all Atlas participants assigned to that county by geographic position.

73    Thereby controlling for counties with a mixture of genetic clusters, also accounting for

74    uneven cluster representation in that county.

75    To control for the effects of age and sex on IPF risk when regressing county IPF burden

76    with risk allele frequency, for each county in the Republic of Ireland we extracted the

77    male-female ratio and average age as collected by the CSO in the 2016 census.

78  **References:**

79  1. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, Maller J, Sklar

80     P, de Bakker PI, Daly MJ, Sham PC. PLINK: a tool set for whole-genome

81     association and population-based linkage analyses. *Am J Hum Genet* 2007; 81:

82     559-575.

83  2. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation

84     PLINK: rising to the challenge of larger and richer datasets. *Gigascience* 2015; 4:

85     7.

86  3. Gilbert E, O'Reilly S, Merrigan M, McGettigan D, Molloy AM, Brody LC, Bodmer W,

87     Hutnik K, Ennis S, Lawson DJ, Wilson JF, Cavalleri GL. The Irish DNA Atlas:

88     Revealing Fine-Scale Population Structure and History within Ireland. *Sci Rep*

89     2017; 7: 17199.

90  4. Desch KC, Ozel AB, Siemieniak D, Kalish Y, Shavit JA, Thornburg CD, Sharathkumar

91     AA, McHugh CP, Laurie CC, Crenshaw A, Mirel DB, Kim Y, Cropp CD, Molloy AM,

92     Kirke PN, Bailey-Wilson JE, Wilson AF, Mills JL, Scott JM, Brody LC, Li JZ,

93     Ginsburg D. Linkage analysis identifies a locus for plasma von Willebrand factor

94     undetected by genome-wide association. *Proc Natl Acad Sci U S A* 2013; 110:

95     588-593.

96  5. Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, Chen WM. Robust

97     relationship inference in genome-wide association studies. *Bioinformatics* 2010;

98     26: 2867-2873.

99  6. Browning BL, Tian X, Zhou Y, Browning SR. Fast two-stage phasing of large-scale

100    sequence data. *Am J Hum Genet* 2021; 108: 1880-1890.

101    7. Zhou Y, Browning SR, Browning BL. A Fast and Simple Method for Detecting Identity-

102        by-Descent Segments in Large-Scale Data. *Am J Hum Genet* 2020; 106: 426-437.

103    8. Traag VA, Waltman L, van Eck NJ. From Louvain to Leiden: guaranteeing well-

104        connected communities. *Sci Rep* 2019; 9: 5233.

105    9. Das S, Forer L, Schonherr S, Sidore C, Locke AE, Kwong A, Vrieze SI, Chew EY, Levy

106        S, McGue M, Schlessinger D, Stambolian D, Loh PR, Iacono WG, Swaroop A,

107        Scott LJ, Cucca F, Kronenberg F, Boehnke M, Abecasis GR, Fuchsberger C. Next-

108        generation genotype imputation service and methods. *Nat Genet* 2016; 48: 1284-

109        1287.

110