1　**Large-scale whole exome sequencing studies identify two genes, *CTSL* and *APOE*,**

2　**associated with lung cancer**

3

4　Jingxiong Xu[1], Wei Xu[2,3], Jiyeon Choi[4], Yonathan Brhane[1], David C. Christiani[5], Jui Kothari[6],

5　James McKay[7], John K. Field[8], Michael P.A. Davies[8], Geoffrey Liu[2,3], Christopher I. Amos[9,10],

6　Rayjean J. Hung[1,3], Laurent Briollais[1,3]*

7

8　1.Prosserman Centre for Population Health Research, Lunenfeld-Tanenbaum Research Institute,

9　Sinai Health, Toronto, Ontario~~ON~~, Canada.

10　2.Princess Margaret Cancer Center, University Health Network, Toronto, Ontario~~ON~~, Canada.

11　3.Dalla Lana School of Public Health, University of Toronto, Toronto, Ontario~~ON~~, Canada.

12　4.Division of Cancer Epidemiology and Genetics, National Cancer Institute, National Institutes of

13　Health, Bethesda, ~~MD~~Maryland, United States of America~~, USA.~~

14　5.T. H. Chan School of Public Health, Harvard University, Boston, ~~MA~~Massachusetts, ~~;~~

15　~~USA.~~United States of America

16　6.Department of Environmental Health, T. H. Chan School of Public Health, Harvard University,

17　Boston, ~~MA~~Massachusetts, ~~USA.~~United States of America

18　7.International Agency for Research on Cancer, Lyon, France.

19　8.Department of Molecular and Clinical Cancer Medicine, The University of Liverpool, Liverpool,

20　~~UK.~~United Kingdom

21　9.Dan L. Duncan Comprehensive Cancer Center, Department of Medicine, Baylor College of

22　Medicine, Houston,~~, TX~~ Texas, ~~USA~~United States of America

23    10.Institute for Clinical and Translational Research, Baylor College of Medicine, Houston,

24    ~~TX~~Texas, ~~, USA.~~United States of America

25    *~~Corresponding author~~

26    ~~Email:~~ laurent@lunenfeld.ca

27

28    **Abstract**

29    Common genetic variants associated with lung cancer have been well studied in the past decade.

30    However, only 12.3% heritability has been explained by these variants. In this study, we

31    investigate the contribution of rare variants (RVs) (minor allele frequency <0.01) to lung cancer

32    through two large whole exome sequencing case-control studies. We first performed gene-based

33    association tests using a novel Bayes Factor statistic in the International Lung Cancer Consortium,

34    the discovery study (European, 1042 cases vs. 881 controls). The top genes identified are further

35    assessed in the UK Biobank (European, 630 cases vs. 172 864 controls), the replication study.

36    After controlling for the false discovery rate, we found two genes, *CTSL* and *APOE,* significantly

37    associated with lung cancer in both studies. Single variant tests in UK Biobank identified 4 RVs

38    (3 missense variants) in *CTSL* and 2 RVs (1 missense variant) in *APOE* stongly associated with

39    lung cancer (OR between 2.0 and 139.0).  The role of these genetic variants in the regulation of

40    *CTSL* or *APOE* expression remains unclear. If such a role is established, this could have important

41    therapeutic implications for lung cancer patients.

42

43    Author summary

44    Lung cancer (LC) is the leading cause of cancer death accounting for 18% of all cancer deaths.

45    Previous studies have suggested genetic contribution to the disease. Common genetic variants

46    associated with LC have been well studied through large, collaborative, genome-wide association

47    studies (GWASs) in the past decade. However, they explained only about 12.3% of LC heritability.

48    It is therefore hypothesized that the unexplained variability might be partially due to rare variants

49    (RVs). In this study, we applied a novel gene-based test statistic based on a Bayes Factor approach,

50    to whole exome sequencing data from the International Lung Cancer consortium (ILCCO).

51  Independent replication of the top genes identified was performed using the UK Biobank data. We

52  found two genes, *CTSL* and *APOE,* significantly associated with LC in both studies. Within these

53  two genes, several RVs showed strong associations with lung cancer in the UK Biobank data.

54  These findings could suggest potential molecular mechanisms leading to lung cancer and more

55  importantly, possible therapeutic targets for personalized treatment.

## Introduction

Lung cancer (LC) is the most commonly diagnosed cancer in men and the third most commonly occurring cancer in women worldwide as estimated in 2018 [1], with an estimated 2.3 millions new cancers diagnosed annually. It is the leading cause of cancer death worldwide with 1.8 million annual deaths accounting for 18% of all cancer deaths [1]. Although reduction of tobacco consumption remains the most appropriate strategy to reduce LC burden, only 10%–15% of all smokers eventually develop LC [2-4]. In Asian countries, up to 30%–40% of lung cancer cases occur in never smokers [4], which suggests a possible role of genetic factors among others.

Common genetic variants associated with LC have been identified through large, collaborative, genome-wide association studies (GWASs), including susceptibility loci at *CHRNA3/5, TERT, HLA, BRCA2, CHEK2* [5,6]. Yet, they explained only about 12.3% of LC heritability reported in a recent GWAS[7]. It is therefore hypothesized that some of the unexplained variability might be due to rare variants (RVs) [8]. A recent study was able to identify 48 germline RVs with deleterious effects on LC in known candidate genes such as *BRCA*2 in a sample of 260 case patients with the disease and 318 controls [9]. More recently, Liu et al. [10] identified 25 deleterious RVs associated with LC susceptibility, including 13 reported in ClinVar. Of the five validated candidates, the authors identified two pathogenic variants in known LC susceptibility loci, *ATM* p.V2716A (Odds Ratio 19.55, 95%CI [5.04,75.6]) and *MPZL2* p.I24M frameshift deletion (Odds Ratio 3.88, 95%CI [1.71,8.8]); and three in novel LC susceptibility genes including *POMC*, *STAU2* and *MLNR*.

To improve the detection of RVs in sequencing studies, we recently proposed a gene-based test for case-control study designs using a Bayes Factors (BF) statistic [11], comparing the total RV counts between cases and controls. Informative priors can be included in this setting, making the BF also sensitive to allelic distribution differences at single variant sites between cases and

79  controls. To elucidate the inherited germline RVs associated with LC, we applied our novel BF

80  approach to whole exome sequencing (WES) data from the International Lung Cancer consortium

81  (ILCCO) [10], with the goal  to identify new genes associated with LC specifically focused on

82  RVs as well as potential causal variants within these genes. Independent replication of the most

83  promising genes and RVs was performed in the UK Biobank data [12].

84

85  **Methods**

86  *Ethics Statement*

87  All participants provided written informed consent, and the study was reviewed and approved by

88  institutional ethic committee of each study site including HSPH-MGH, University Health Network

89  and Mount Sinai Hospital in Toronto (Toronto), University of Liverpool in UK (Liverpool) and

90  IARC.

91

92  *Study population for gene-based and RV discovery*

93  Case patients with LC and matched healthy individuals were identified from four independent case

94  series that form the ILCCO consortium, including Harvard University School of Public

95  Health/Massachusetts General Hospital (HSPH-MGH), University Health Network and Mount

96  Sinai Hospital in Toronto (Toronto), University of Liverpool in UK (Liverpool) and the

97  International Agency for Research on Cancer (IARC). The original data includes 2047 samples,

98  of which 44 are HapMap controls and 68 were flagged by the Center for Inherited Disease

99  Research (CIDR) as duplicates, related individuals or quality control outliers. Whole exome

100 sequencing was performed for selected LC cases and frequency-matched unaffected controls, to

101 identify novel common and rare genetic variants associated with LC risk. To enrich the relevance

102    of genetics in the cases, LC patients were preferentially selected from those with a family history

103    of LC among first-degree relative or early-onset (<60 years). About the same number of controls

104    were selected, frequency-matched by age and sex with the cases. To adjust for population

105    stratification, principal components (PCs) were derived from the genome-wide data from the

106    ILCCO. The analysis was restricted to those with European ancestry. The representation of the top

107    3 PCs (S1 Fig) identified one outlier participant with possible non-European ancestry, and was

108    removed from the analysis. We further removed 10 individuals with genotype missing rate >10%

109    and one individual was flagged with very low heterozygosity rate (> 6 standard deviations below

110    the mean heterozygosity). After the filtering steps, a total of 1923 subjects remained in the study

111    and were included in the analyses. ~~All participants provided written informed consent, and the~~

112    ~~study was reviewed and approved by institutional ethic committee of each study site including~~

113    ~~HSPH MGH, University Health Network and Mount Sinai Hospital in Toronto (Toronto),~~

114    ~~University of Liverpool in UK (Liverpool) and IARC.~~

115

116    *Study population for gene-based and RV replication*

117    We used UK Biobank WES data as the validation set [13,14]. Among the total number of 200,643

118    samples, our analysis includes all LC patients after excluding those diagnosed at most 5 years

119    before any other primary cancers and controls with no cancer diagnosis history. We also removed

120    at random one individual from each pair of individuals closer than $3^{rd}$ degree relatives (kinship

121    coefficient > 0.0884), and subjects who self-reported a non-white ethnic background. After the

122    filtering, 173,494 individuals remained in the study.

123

124    *Germline Sequencing/QC*

125 *ILCCO:* The sequencing of whole exomes and additional targeted regions of DNA samples from

126 all 4 different sites was performed at the CIDR. Targeted regions were selected based on previous

127 associations with LC or with histological LC subtypes from GWASs on common variants [5,6].

128 After initial quality control (QC) analysis by CIDR [10], the mean on-target coverage was 52X

129 and more than 97% of targeted bases had a depth greater than 10X. Further QC analysis was

130 performed including the following steps: i) Exclusion of variants with QUAL<100 indicating a

131 low probability that there is a variant at a site  or mean GQ<50 indicating low probabilities that

132 genotype calls were correct across individuals at a site so that Ts/Tv ratio is greater than 2 (S2 Fig);

133 ii) Exclusion of singleton variants (variant with occurrence of only 1 minor allele) when minor

134 allele has GQ<50 or depth <20; iii) Exclusion of non-biallelic variants and variants on the sex

135 chromosome; iv) Exclusion of variants with p-value of Hardy-Weinberg equilibrium test <1e-7 in

136 the control samples ; v) Set individual genotype as missing if GQ<30 or depth<10; vi) Exclusion

137 of variants with minor allele frequency (MAF)>1% (MAF was estimated using study population).

138 The MAF distribution of the remaining RVs is given in Table 1.

139 **Table 1. MAF distribution of genetic variants in the discovery study (ILCCO)**

| MAF | 0 | (0,0.01) | [0.01,0.05] | [0.05,0.5) | Total |
|---|---|---|---|---|---|
| #(Rare Variants) | 136485 | 1022101 | 60288 | 129789 | 1348663 |
| Proportion (%) | 10.12 | 75.79 | 4.47 | 9.62 | 100 |

140
141

142 *UK Biobank*: We performed the following QC steps for all genes selected in the discovery set: i)

143 exclude variants that are not bi-allelic and those with QUAL<10; ii) filter out variants with mean

144 GQ<30 as well as singleton variants with depth <20 or GQ<40; iii) set genotype missing if

145 depth<10 or GQ<20, and exclude variants with missing genotype rate >10%; iv) exclude variants

146 with MAF>1% (MAF estimated using study population).

147     In both the discovery and replication studies, for our gene-based analyses, we considered -/+ 1k

148     bp up- and down-stream sites of each gene (including non-exonic RVs) for the analysis.

149

150     *Gene-based analysis*

151     To increase the power of discovering genes associated with LC, we applied a gene-based approach

152     based on a Bayes Factor (BF) statistic that we recently developed, to both the discovery and

153     replication studies [11]. It was designed specifically to test the association between a set of RVs

154     located in the same region or in a gene and a disease outcome in the context of case-control designs.

155     An advantage of our BF approach over existing methods is the possibility to introduce an

156     "informative" prior to gain power to detect gene-based associations, where this prior is sensitive

157     to allelic differences between cases and controls for a particular gene (S1 Text). Compared to the

158     commonly-used SKAT gene-based test [15], our BF approach is more sensitive to an excess of

159     small p-values from single RV tests within each gene while SKAT has better power to detect genes

160     exhibiting systematic allelic differences between cases and controls across all RVs. This difference

161     was discussed in details in [11] and illustrated on two genes that showed large discrepancy in

162     overall ranking when applying these two approaches [11]. In this study, we applied two versions

163     of the BF test statistic, $BF_{KS}$ and $BF_{SKAT}$, where either a Kolmogorov-Smirnov (KS) or SKAT p-

164     value is used as informative prior. This gave us higher chance to detect genes that may have

165     different underlying RV allelic distribution differences between cases and controls. The respective

166     advantage of each approach is described in details in the S1 Text. In this paper, we mainly focused

167     on $BF_{KS}$ and used $BF_{SKAT}$ as a secondary analysis.

168     To assess the sensitivity of the association tests on confounding variables, we conducted sensitivity

169     analyses on the genome-wide significant genes and adjusted our analyses for age, sex, smoking

170  and the top 5 PCs used to control for population stratification. Both the BF and the prior

171  components (KS or SKAT p-value) were adjusted. The extention of $BF_{KS}$ and $BF_{SKAT}$ incorporating

172  covariates is described in S1 Text.

173

174  *Single RV-based analysis*

175  For the two genes that passed a gene-based replication genome-wide significance level (see below),

176  i.e., *APOE* and *CTSL*, we performed single RV tests only with UK Biobank since this study has

177  larger coverage of RVs. We used the Firth's bias-reduced logistic regression to deal with sparse

178  allelic counts [16]. Analyses were adjusted for age, sex, smoking status (ever vs. never smoking)

179  and the top five PCs. RVs that pass a FDR adjusted q value [17] of 0.01 were selected.

180

181  *Significance threshold for gene-based replication analysis*

182  We denote $P_d$ the $P$ value for selecting genes in the discovery cohort (ILCCO) and $P_r$ the $P$ value

183  for selecting a gene in the replication cohort (UK Biobank). We set $\gamma$ as the significance threshold

184  for selecting genes in the discovery cohort and which will be followed-up for replication in UK

185  biobank and $\lambda$ the significance level in the replication cohort. To control the gene-based family-

186  wise error rate (FWER) $\alpha$, we can determine $\gamma$ and $\lambda$ such that,

187
$$FWER_{(P_d \leq \gamma, P_r \leq \lambda)} = Pr(V \geq 1) \leq \alpha,$$

188  where V is number of genes declared achieved signficiance levels in both discovery and validation

189  studies, $P_d \leq \gamma$ and $P_r \leq \lambda$, where $\gamma$ and $\lambda$ were determined through permutation analysis, as

190  follows. First, we repeated analyses of ILCCO (discovery set) and UK Biobank (validation set)

191  studies 100 times, where each time the phenotype of individuals was permuted. Second, we

192  determined the two thresholds such that among 100 replicates, the number of identified significant

193 genes is less or equal to $100 \times \alpha = 100 \times 0.05 = 5$, for a genome-wide control of FWER$\leq$ 5%.

194 We found the following thresholds, $\gamma = 5 \times 10^{-4}$ and $\lambda = 0.05$ in the discovery and validation study,

195 respectively, when using $BF_{KS}$ as the test statistic (i.e., our main statistic). Therefore, in our

196 application analysis, the set of genes that passed a significance threshold of $\gamma = 5 \times 10^{-4}$ in the

197 discovery (ILCCO) cohort and $\lambda = 0.05$ in the replication (UK Biobank) cohort were declared

198 associated with the disease and replicated.

199

200 **Results**

201 *Characteristics of patients in the discovery and replication studies*

202 Our discovery study (ILCCO) includes 1042 lung cancer cases and 881 controls (HSPH-MGH,

203 426 cases and 270 controls; Toronto, 259 cases and 258 controls; Liverpool, 64 cases and 69

204 controls; IARC, 293 cases and 284 controls). The replication study (UK Biobank) includes a total

205 of 630 cases and 172,864 controls. In the discovery study, the distributions of sex and age are

206 comparable between cases and controls. However, in the replication study, there is an excess of

207 males in cases compared to controls (52.7% vs. 45.2%, $P=1.9 \times 10^{-4}$) and cases are older age at

208 enrollment compared to controls (mean=62.0 vs. 56.7 years, $P<2.2 \times 10^{-16}$) (Table 2). As expected,

209 there is a higher proportion of never smokers in controls compared to cases (35.2% vs. 11.8%

210 $P<2.2 \times 10^{-16}$ and 54.6% vs. 14.8% $P<2.2 \times 10^{-16}$ in the discovery and replication study, respectively).

211
212 **Table 2. Basic demographic characteristics in the discovery and validation studies**
213

| | Discovery (ILCCO) | | | Replication (UK Biobank) | | |
|---|---|---|---|---|---|---|
| | controls n=881 | cases n=1042 | p-value | controls n=172864 | cases n=630 | p-value |
| Sex, No. (%) | | | NS | | | 1.9E-04 |
| M | 513 (58.2) | 613 (58.8) | | 78163 (45.2) | 332 (52.7) | |
| F | 368 (41.8) | 429 (41.2) | | 94701 (54.8) | 298 (47.3) | |

| | | | | | | |
|---|---|---|---|---|---|---|
| Age, mean (SD) | 60.8 (11.8) | 62.2 (12.3) | NS | 56.7 (8.0) | 62.0 (5.8) | <2.2E-16 |
| Smoking, No. (%) | | | <2.2E-16 | | | <2.2E-16 |
| Never | 310 (35.2) | 123 (11.8) | | 94378 (54.6) | 93 (14.8) | |
| Former | 375 (42.6) | 421 (40.4) | | 61770 (35.7) | 319 (50.6) | |
| Current | 193 (21.9) | 492 (47.2) | | 16119 (9.3) | 214 (34.0) | |
| Missing | 3 (0.3) | 6 (0.6) | | 597 (0.3) | 4 (0.6) | |

214    NS: not significant

215

216    *Gene-Based analysis*

217    In the discovery study, a total of 13,872 genes with at least 20 bi-allelic RVs were analyzed based

218    on the QC pipeline described. The QQ plots corresponding to $2\log(BF_{KS})$ and $2\log(BF_{SKAT})$

219    statistics are presented in Fig 1 and confirm that they are both asymptotically distributed as $\chi^2(3)$.

220    Using a significance level of $\gamma = 5 \times 10^{-4}$ in the discovery cohort (see Methods section), a total of

221    17 genes based on $BF_{KS}$ and 14 genes using $BF_{SKAT}$ (Tables 3-4) were selected for replication. The

222    2 top genes are *CTSL* ($P=4.9 \times 10^{-5}$) and *TBX4* ($P=6.5 \times 10^{-5}$) with $BF_{KS}$, *VAV2* ($P=1.9 \times 10^{-5}$) and

223    *DENND4B* ($P=4.3 \times 10^{-5}$) with $BF_{SKAT}$. Four genes are found by both test statistics including *CTSL*,

224    *TBX4*, *C8orf44*, and *DGKB*. Using a significance level of $\lambda = 0.05$ (see Methods section) in the

225    replication study, we were able to replicate only one gene, *CTSL* ($P=2.7 \times 10^{-3}$), when using the

226    $BF_{KS}$ test and the two genes *APOE* ($P=1.9 \times 10^{-3}$) and *CTSL* ($P=6.9 \times 10^{-6}$) based on the $BF_{SKAT}$ test

227    (Tables 3-4). For each gene identified in the discovery set, we calculated an overall p-value in

228    Ttables 3-4 by combining p-values from the discovery and validation sets using Fisher's method

229    [18].

230

231    *Sensitivity analysis*

232 We found that the association signal for *CTSL* did not change much after adjustment for

233 confounders using $BF_{KS}$ (unadjusted: discovery p-value=4.87E-05, validation p-value=2.75E-03;

234 adjusted: discovery p-value=2.84E-05, validation p-value=3.88E-03) (S1 Table) and $BF_{SKAT}$

235 (unadjusted: discovery p-value=4.30E-04, validation p-value=1.31E-05; adjusted: discovery p-

236 value=1.32E-03, validation p-value=4.33E-05) (S2 Table). While the adjusted association using

237 $BF_{SKAT}$ on *APOE* (discovery p-value=2.12E-03, validation p-value=8.24E-03) (S2 Table) was not

238 as significant as the unadjusted $BF_{SKAT}$ (discovery p-value=2.56E-04, validation p-value=4.01E-

239 03). Of note, in this analysis, 9 out of 1923 individuals were removed from ILCCO study due to

240 the missing smoking status and 761 out of 173,494 individuals were removed from UK Biobank

241 study due to the missing values of smoking and/or PCs.

242

243 **Table 3. Results of gene-based analyses using $BF_{KS}$ test[a] in the discovery and replication**
244 **studies**

| Rank | Genes | Chr | #(Sites) | Discovery (ILCCO) | | Replication (UK Biobank) | | Combined P |
|------|-------|-----|----------|-------|-------|-------|-------|-----------|
| | | | | KS P[b] | $BF_{KS}$ P[c] | KS P[b] | $BF_{KS}$ P[c] | Fisher's method |
| 1 | *CTSL* | 9 | 25 | 1.32E-03 | **4.87E-05** | 8.43E-01 | **2.75E-03** | **2.26E-06** |
| 2 | *TBX4* | 17 | 37 | 1.48E-03 | 6.49E-05 | 9.67E-01 | 9.96E-01 | 6.88E-04 |
| 3 | *RASL10B* | 17 | 53 | 4.05E-04 | 6.75E-05 | 1.00E+00 | 9.81E-01 | 7.03E-04 |
| 4 | *MUC3A* | 7 | 94 | 1.30E-04 | 7.33E-05 | 5.95E-01 | 6.42E-01 | 5.16E-04 |
| 5 | *AMN* | 14 | 22 | 1.68E-04 | 8.08E-05 | 9.07E-01 | 9.71E-01 | 8.20E-04 |
| 6 | *KRTAP19-4* | 21 | 21 | 3.38E-05 | 1.27E-04 | 8.74E-02 | 8.76E-02 | 1.38E-04 |
| 7 | *KRTAP19-5*[d] | 21 | 20 | 3.38E-05 | 1.28E-04 | NA | NA | NA |
| 8 | *CPB2* | 13 | 25 | 1.74E-03 | 1.46E-04 | 1.00E+00 | 6.71E-01 | 1.01E-03 |
| 9 | *C8orf44* | 8 | 38 | 1.11E-02 | 2.17E-04 | 6.74E-01 | 1.82E-01 | 4.39E-04 |
| 10 | *ZW10* | 11 | 48 | 6.49E-04 | 2.23E-04 | 8.79E-01 | 7.19E-01 | 1.56E-03 |
| 11 | *INHA* | 2 | 68 | 1.05E-04 | 2.51E-04 | 1.00E+00 | 9.04E-01 | 2.13E-03 |
| 12 | *DGKB* | 7 | 79 | 3.33E-02 | 3.27E-04 | 9.73E-01 | 9.34E-01 | 2.77E-03 |
| 13 | *FBXO6* | 1 | 55 | 2.09E-03 | 3.34E-04 | 3.47E-01 | 3.51E-01 | 1.18E-03 |
| 14 | *PHF12* | 17 | 82 | 2.00E-03 | 3.57E-04 | 1.00E+00 | 8.35E-01 | 2.71E-03 |
| 15 | *LEMD3* | 12 | 46 | 8.70E-04 | 3.58E-04 | 1.00E+00 | 9.98E-01 | 3.20E-03 |
| 16 | *OR5AC2* | 3 | 70 | 1.09E-04 | 3.85E-04 | 1.00E+00 | 1.00E+00 | 3.41E-03 |
| 17 | *FGF8* | 10 | 38 | 9.89E-02 | 4.52E-04 | 9.93E-01 | 7.29E-01 | 2.97E-03 |

**Table 4. Results of gene-based analyses using $BF_{SKAT}$[a] in the discovery and replication studies**

| Rank | Genes | Chr | #(Sites) | Discovery (ILCCO) | | Replication (UK Biobank) | | Combined P |
|------|-------|-----|----------|-------------------|---|--------------------------|---|------------|
| | | | | SKAT P[b] | $BF_{SKAT}$ P[c] | SKAT P[b] | $BF_{SKAT}$ P[c] | Fisher's method |
| 1 | *VAV2* | 9 | 121 | 3.09E-04 | 1.95E-05 | 6.72E-01 | 5.72E-01 | 1.39E-05 |
| 2 | *DENND4B* | 1 | 69 | 2.21E-05 | 4.31E-05 | 9.96E-01 | 6.35E-01 | 3.15E-04 |
| 3 | *TBX4* | 17 | 37 | 1.95E-03 | 8.41E-05 | 8.21E-01 | 9.41E-01 | 8.27E-04 |
| 4 | *RHBDL3* | 17 | 27 | 9.09E-03 | 1.06E-04 | 1.63E-01 | 2.91E-01 | 3.51E-04 |
| 5 | *C8orf44* | 8 | 38 | 5.89E-03 | 1.19E-04 | 9.97E-01 | 2.52E-01 | 3.43E-04 |
| 6 | *CCT8* | 21 | 46 | 2.43E-02 | 2.41E-04 | 9.87E-01 | 9.99E-01 | 2.25E-03 |
| 7 | *SIGLEC11* | 19 | 24 | 3.10E-03 | 2.46E-04 | 7.23E-01 | 5.81E-01 | 1.41E-03 |
| 8 | *APOE* | 19 | 25 | 2.65E-04 | **2.56E-04** | 6.10E-03 | **4.01E-03** | **1.52E-05** |
| 9 | *POMK* | 8 | 33 | 3.00E-02 | 3.27E-04 | 9.54E-01 | 7.50E-01 | 2.29E-03 |
| 10 | *DGKB* | 7 | 79 | 4.34E-02 | 4.20E-04 | 3.79E-01 | 5.10E-01 | 2.02E-03 |
| 11 | *CTSL* | 9 | 25 | 1.29E-02 | **4.30E-04** | 3.08E-03 | **1.31E-05** | **1.13E-07** |
| 12 | *CPB2* | 13 | 25 | 5.55E-03 | 4.42E-04 | 2.98E-01 | 2.65E-01 | 1.18E-03 |
| 13 | *ITGB6* | 2 | 61 | 3.23E-02 | 4.93E-04 | 9.83E-01 | 9.40E-01 | 4.02E-03 |
| 14 | *VCPIP1* | 8 | 39 | 1.73E-02 | 4.94E-04 | 8.73E-01 | 7.00E-01 | 3.10E-03 |

*Single RV-based analysis*

In UK Biobank, a total of 155 bi-allelic RVs for *CTSL* and 174 for *APOE* were included in the

analysis. In *CTSL*, 4 RVs were found associated with LC at an FDR q-value of 0.01, including

variant at positions 87728433 (rs771328780), 87729621 (rs778002071), 87730426 (rs777251059)

and 87727608 (rs112682750) on chromosome 9 (Table 5), where the last 3 were missense variants.

In *APOE*, 2 RVs passed this significance level, including variant at position 44907893 (rs number

not available) and 44906640 (rs1568615382) on chromosome 19. Most of the variants found to be

264    associated with LC risk are very rare (MAF$<10^{-4}$ in controls), except one missense variant in *CTSL*,

265    rs112682750,  has a MAF of $7.7\times10^{-3}$.

**Table 5. Results of single RV-based association analysis in the genes *CTSL* and *APOE* using UK Biobank data**

| Gene (Variant, position) | ClinVar Significance [3719] | Overall (N=173,494) | | Cases (N=630) | | Controls (N=172,864) | | Association | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | MAF | #Carriers | MAF | #Carriers | MAF | #Carriers | P value[a] | FDR q-value[b] | Odds Ratio[a] (95% CI) |
| *CTSL* (rs771328780, chr9: 87,728,433) | Unknown | 4.0E-5 | 14 | 1.6E-03 | 2 | 3.5E-05 | 12 | 6.7E-5 | 7.1E-4 | 83.9 (18.2-387.2) |
| *CTSL* (rs778002071, chr9: 87,729,621) | Missense | 2.0E-05 | 7 | 7.9E-4 | 1 | 1.7E-05 | 6 | 8.0E-4 | 4.3E-3 | 139.0 (20.7-933.7) |
| *CTSL* (rs777251059, chr9: 87,730,426) | Missense | 1.4E-05 | 5 | 7.9E-04 | 1 | 1.2E-05 | 4 | 3.9E-3 | 9.8E-3 | 54.8 (7.8-382.6) |
| *CTSL* (rs112682750, chr9: 87,727,608) | Missense | 7.8E-03 | 2694 | 1.5E-02 | 19 | 7.7E-03 | 2675 | 7.8E-03 | 0.01 | 2.0 (1.3,3.1) |
| *APOE* (chr19: 44,907,893) | Unknown | 1.2E-05 | 4 | 7.9E-04 | 1 | 8.7E-06 | 3 | 2.8E-4 | 5.5E-3 | 276.3 (38.5-1985.3) |
| *APOE* (rs1568615382 chr19: 44,906,640) | Missense | 3.2E-05 | 11 | 7.9E-04 | 1 | 2.9E-05 | 10 | 1.4E-4 | 0.01 | 90.0 (15.5-523.9) |

[a]Based on the Firth biased-corrected logistic regression [165]

[b]Only RVs with a q-value ≤ 0.01 were selected.

All the 6 RVs are associated with increased LC risk as indicated by an odds-ratio>1 in UK Biobank. One of the 6 RVs was present in ILCCO, rs112682750 in *CTSL*, but it did not show association with LC after adjustment for age, sex, smoking and PCs (*P*=0.19).

*Genomic region analysis of rs112682750 in CTSL*

Using cancer cell lines from the USCS genome browser, a genomic analysis of the region around rs112682750 indicates that this variant is located within a promoter/enhancer region of *CTSL* in lung related cells (S3 Fig). This suggests that rs112682750 might affect the transcription of *CTSL*.

*Annotation of Single RVs in CTSL and APOE*

We searched functional annotation for the 6 associated RVs identified from *CTSL and APOE* using Ensembl Variant Effect Predictor (VEP) [~~19~~20], Combined Annotation Dependent Depletion (CADD) [~~20~~21,~~21~~22] and Functional Annotation of Variants – Online Resource (FAVOR) [~~22~~23]. The search results indicated that rs778002071 (*CTSL*) was categorized as deleterious nonsynonymous variant, according to all three annotation resources, and the rest 5 RVs were predicted to be tolerated (benign) by at least one resource (Table 6).

**Table 6 Functional annotation of rare variants in the genes *CTSL* and *APOE***

| SNP | Allele | Amino acids | Codons | PolyPhen Category[a] | Val[b] | SIFT Category[c] | Val[d] | FAVOR aPC-Protein-Function[e] Category | PHRED | Percentile | CADD[f] PHRED |
|---|---|---|---|---|---|---|---|---|---|---|---|
| rs771328780 (*CTSL*, 87,728,433) | G | - | - | - | - | - | - | intronic | 2.97 | - | 3.90 |
| rs778002071 (*CTSL*, 87,729,621) | A | G/S | Ggc/Agc | possibly damaging | 0.861 | deleterious | 0.02 | exonic, nonsynonymous | 28.03 | 0.16 | 26.10 |
| rs777251059 (*CTSL*, 87,730,426) | C | G/A | gGt/gCt | benign | 0.059 | tolerated | 0.33 | - | - | - | 21.60 |
| rs112682750 (*CTSL*, 87,727,608) | C | N/T | aAt/aCt | benign | 0.001 | tolerated | 0.99 | exonic, nonsynonymous | 22.17 | 0.61 | 15.00 |
| - (*APOE*, 44,907,893) | A | Q | caG/caA | - | - | - | - | - | - | - | 3.97 |
| rs1568615382 (*APOE*, 44,906,640) | G | A/T | Gct/Act | Possibly damaging | 0.536 | tolerated | 0.09 | - | - | - | 22.9 |

a. PolyPhen category of change [~~389~~19].
b. PolyPhen score: It predicts the functional significance of an allele replacement from its individual features. Range: [0, 1] (default: 0) [~~389~~19].
c. SIFT category of change [~~3940~~24].
d. SIFT score, ranges from 0.0 (deleterious) to 1.0 (tolerated). Range: [0, 1] (default: 1) [~~3940~~24].
e. Protein function annotation PC: the first PC of the standardized scores of "SIFTval, PolyPhenVal, Grantham, Polyphen2_HDIV_score, Polyphen2_HVAR_score, MutationTaster_score, MutationAssessor_score" in PHRED scale. Range: [2.974, 86.238] [~~22~~23].
f. The CADD score in PHRED scale (integrative score). A higher CADD score indicates more deleterious. Range: [0.001, 84] [~~2021~~,~~2122~~].

**Discussion**

By focusing on rare variants using whole exome sequencing data, we identified two new genes, *CTSL* and *APOE*, associated with LC in the ILCCO study, that were replicated in the UK Biobank study. In *CTSL*, 3 missense RVs and 1 RV with unknown significance were discovered as associated with LC in the UK Biobank study. In *APOE*, 1 missense variant and 1 with unknown significance were discovered.

The Cathepsin L gene (*CTSL),* is a ubiquitously expressed lysosomal endopeptidase that is primarily involved in terminal degradation of intracellular and endocytosed proteins [2125]. *CTSL* has recently gained attentions for its roles in SARS-CoV2 entry to host cell by cleaving receptor-bound viral spike protein, which results in further activation and infection[2426,2527]. While potential functional connection between viral infection and lung cancer susceptibility remains to be established, *CTSL* also has roles relevant in tumorigenesis and progression. *CTSL* upregulation has been reported in a wide range of human malignancies including ovarian, breast, prostate, lung, gastric, pancreatic and colon cancers [2628]. Importantly, evidence indicates that *CTSL* expression may be linked to cancer grade and stage. In LC patients, higher *CTSL* activity has been reported compared to non-malignant tissue as well as association between tumor grade and upregulated serum levels [2729]. The role of *CTSL* in promoting tumor progression and metastatic aggressiveness has also been suggested [2830]. Significant interest in the development of *CTSL* intervention strategies has also emerged. For example, *CTSL* downregulation through RNA interference in different tumor models (including glioma, osteosarcoma, myeloma and melanoma) resulted in consistent inhibition of tumorigenicity and invasiveness of neoplastic cells [29-3231-34]. The identification of patients who might benefit from anti-CTSL therapy remains an important clinical question. The identification of new RVs that correlate with LC risk in our study could

19

therefore help identify these patients. Although the impacts of these variants to CTSL levels or activity in early vs. late stages of lung tumorigenesis need to be established, potential regulatory function of the most common variant we identified in *CTSL*, rs112682750, for instance, could be hypothesized.

The apolipoprotein E gene (*APOE*) codes for a protein associated with lipid particles, that mainly functions in lipoprotein-mediated lipid transport between organs via the plasma and interstitial fluids. *APOE* is also associated with atherosclerogenesis, which itself has been involved in tumor development. *APOE* has been shown to act as a growth factor that can influence carcinogenesis [3335]. In patients with LC, the levels of *APOE* gene expression were significantly higher in cancer tissue than in adjacent non-cancer tissue [3436]. Serum *APOE* has also been associated with lymph node metastasis in lung adenocarcinoma patients [3537]. It was also reported that high expression of *APOE* promotes cancer cell proliferation and migration and contributes to an aggressive clinical course in patients with lung adenocarcinoma [3638]. *APOE* has also raised interest for therapeutic interventions. For instance, *APOE* was involved in the inhibition of melanoma metastasis and angiogenesis by stimulating the immune response to tumor cells [3739]. Identification of genetic variants that could regulate *APOE* expression could therefore have important therapeutic implications. Of note, *APOE* was only detected with one version of our BF approach (i.e., $BF_{SKAT}$) and further validation of this gene is warranted.

The strengths of our study include the large sample sizes available for discovery and replication of the gene-based analyses and the use of UK Biobank data for RV discoveries. Our statistical approach for gene discovery, the Bayes Factor statistic, has also been shown to have increased power compared to competing approaches such as SKAT and the Burden test [11]. Another significant  advantage is its sensitivity to detect single RV associations through the definition of

informative priors. Under our statistical framework, the discovery of RVs can therefore be thought as a two-step approach where the first step is a gene-based analysis and the second step, an RV association test within the set of significantly associated genes.

Our study contrasts with Liu et al.'s analysis of the ILLCO data [10] in several aspects. They performed single RV analyses focusing only on suspected deleterious variants. In a second step, they performed gene-based tests using only genes that included RVs that were significantly associated with LC after controlling for multiple comparisons from a Burden test. In comparison, we tested all the genes in the discovery cohort and did not make any assumption regarding the possible functional effect of the RVs.

The discovery of RVs in the context of sequencing studies remains a field of intensive research. The limitations of this study include the need for further validation and characterization of the two genes and RVs identified, in particular to correlate them with disease progression outcomes and LC subtypes. Also, the benefit for therapeutic interventions may be considered as it could lead to a more personalized treatment of LC patients targeting specific gene/pathway mechanisms such as the immune response system.

## References

1. Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A. Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. CA Cancer J Clin. 2018 Nov;68(6):394-424.
2. Mattson ME, Pollack ES, Cullen JW. What are the odds that smoking will kill you? Am J Public Health. 1987 Apr;77(4):425-31. doi: 10.2105/ajph.77.4.425. Erratum in: Am J Public Health 1987 Jul;77(7):818.
3. Scagliotti GV, Longo M, Novello S. Nonsmall cell lung cancer in never smokers. Curr Opin Oncol. 2009 Mar;21(2):99-104.
4. Lee YJ, Kim JH, Kim SK, Ha SJ, Mok TS, Mitsudomi T, Cho BC. Lung cancer in never smokers: change of a mindset in the molecular era. Lung Cancer. 2011 Apr;72(1):9-15.
5. Timofeeva MN, Hung RJ, Rafnar T, Christiani DC, Field JK, Bickeböller H, Risch A, McKay JD, Wang Y, Dai J, Gaborieau V, McLaughlin J, Brenner D, Narod SA, Caporaso NE, Albanes D, Thun M, Eisen T, Wichmann HE, Rosenberger A, Han Y, Chen W, Zhu D, Spitz M, Wu X, Pande M, Zhao Y, Zaridze D, Szeszenia-Dabrowska N, Lissowska J, Rudnai P, Fabianova E, Mates D, Bencko V, Foretova L, Janout V, Krokan HE, Gabrielsen ME, Skorpen F, Vatten L, Njølstad I, Chen C, Goodman G, Lathrop M, Benhamou S, Vooder T, Välk K, Nelis M, Metspalu A, Raji O, Chen Y, Gosney J, Liloglou T, Muley T, Dienemann H, Thorleifsson G, Shen H, Stefansson K, Brennan P, Amos CI, Houlston R, Landi MT; Transdisciplinary Research in Cancer of the Lung (ILCCO) Research Team. Influence of common genetic variation on lung cancer risk: meta-analysis of 14 900 cases and 29 485 controls. Hum Mol Genet. 2012 Nov 15;21(22):4980-95.
6. Wang Y, McKay JD, Rafnar T, Wang Z, Timofeeva MN, Broderick P, Zong X, Laplana M, Wei Y, Han Y, Lloyd A, Delahaye-Sourdeix M, Chubb D, Gaborieau V, Wheeler W, Chatterjee N, Thorleifsson G, Sulem P, Liu G, Kaaks R, Henrion M, Kinnersley B, Vallée M, LeCalvez-Kelm F, Stevens VL, Gapstur SM, Chen WV, Zaridze D, Szeszenia-Dabrowska N, Lissowska J, Rudnai P, Fabianova E, Mates D, Bencko V, Foretova L, Janout V, Krokan HE, Gabrielsen ME, Skorpen F, Vatten L, Njølstad I, Chen C, Goodman G, Benhamou S, Vooder T, Välk K, Nelis M, Metspalu A, Lener M, Lubiński J, Johansson M, Vineis P, Agudo A, Clavel-Chapelon F, Bueno-de-Mesquita HB, Trichopoulos D, Khaw KT, Johansson M, Weiderpass E, Tjønneland A, Riboli E, Lathrop M, Scelo G, Albanes D, Caporaso NE, Ye Y, Gu J, Wu X, Spitz MR, Dienemann H, Rosenberger A, Su L, Matakidou A, Eisen T, Stefansson K, Risch A, Chanock SJ, Christiani DC, Hung RJ, Brennan P, Landi MT, Houlston RS, Amos CI. Rare variants of large effect in BRCA2 and CHEK2 affect risk of lung cancer. Nat Genet. 2014 Jul;46(7):736-41. doi: 10.1038/ng.3002. Epub 2014 Jun 1. Erratum in: Nat Genet. 2017 Mar 30;49(4):651.
7. McKay JD, Hung RJ, Han Y, Zong X, Carreras-Torres R, Christiani DC, Caporaso NE, Johansson M, Xiao X, Li Y, Byun J, Dunning A, Pooley KA, Qian DC, Ji X, Liu G, Timofeeva MN, Bojesen SE, Wu X, Le Marchand L, Albanes D, Bickeböller H, Aldrich MC, Bush WS, Tardon A, Rennert G, Teare MD, Field JK, Kiemeney LA, Lazarus P, Haugen A, Lam S, Schabath MB, Andrew AS, Shen H, Hong YC, Yuan JM, Bertazzi PA, Pesatori AC, Ye Y, Diao N, Su L, Zhang R, Brhane Y, Leighl N, Johansen JS, Mellemgaard A, Saliba W, Haiman CA, Wilkens LR, Fernandez-Somoano A, Fernandez-Tardon G, van der Heijden HFM, Kim JH, Dai J, Hu Z, Davies MPA, Marcus MW, Brunnström H, Manjer J, Melander O, Muller DC, Overvad K, Trichopoulou A, Tumino R, Doherty JA, Barnett MP, Chen C, Goodman GE,

Cox A, Taylor F, Woll P, Brüske I, Wichmann HE, Manz J, Muley TR, Risch A, Rosenberger A, Grankvist K, Johansson M, Shepherd FA, Tsao MS, Arnold SM, Haura EB, Bolca C, Holcatova I, Janout V, Kontic M, Lissowska J, Mukeria A, Ognjanovic S, Orlowski TM, Scelo G, Swiatkowska B, Zaridze D, Bakke P, Skaug V, Zienolddiny S, Duell EJ, Butler LM, Koh WP, Gao YT, Houlston RS, McLaughlin J, Stevens VL, Joubert P, Lamontagne M, Nickle DC, Obeidat M, Timens W, Zhu B, Song L, Kachuri L, Artigas MS, Tobin MD, Wain LV; SpiroMeta Consortium, Rafnar T, Thorgeirsson TE, Reginsson GW, Stefansson K, Hancock DB, Bierut LJ, Spitz MR, Gaddis NC, Lutz SM, Gu F, Johnson EO, Kamal A, Pikielny C, Zhu D, Lindströem S, Jiang X, Tyndale RF, Chenevix-Trench G, Beesley J, Bossé Y, Chanock S, Brennan P, Landi MT, Amos CI. Large-scale association analysis identifies new lung cancer susceptibility loci and heterogeneity in genetic susceptibility across histological subtypes. Nat Genet. 2017 Jul;49(7):1126-1132

8.  Manolio, T., Collins, F., Cox, N., Goldstein, D., Hindorff, L., Hunter, D.J., et. al. (2009) Finding the missing heritability of complex diseases. Nature, 461, 747–753.

9.  Liu, Y., Lusk, C.M., Cho, M.H., Silverman, E.K., Qiao, D., Zhang, R. et al. (2018) Rare variants in known susceptibility loci and their contribution to risk of lung cancer. Journal of Thoracic Oncology, 13, 1483–1495.

10. Liu Y, Xia J, McKay J, Tsavachidis S, Xiao X, Spitz MR, Cheng C, Byun J, Hong W, Li Y, Zhu D, Song Z, Rosenberg SM, Scheurer ME, Kheradmand F, Pikielny CW, Lusk CM, Schwartz AG, Wistuba II, Cho MH, Silverman EK, Bailey-Wilson J, Pinney SM, Anderson M, Kupert E, Gaba C, Mandal D, You M, de Andrade M, Yang P, Liloglou T, Davies MPA, Lissowska J, Swiatkowska B, Zaridze D, Mukeria A, Janout V, Holcatova I, Mates D, Stojsic J, Scelo G, Brennan P, Liu G, Field JK, Hung RJ, Christiani DC, Amos CI. Rare deleterious germline variants and risk of lung cancer. NPJ Precis Oncol. 2021 Feb 16;5(1):12.

11. Xu J, Xu W, Briollais L. A Bayes factor approach with informative prior for rare genetic variant analysis from next generation sequencing data. Biometrics. 2021 Mar;77(1):316-328.

12. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, Motyer A, Vukcevic D, Delaneau O, O'Connell J, Cortes A, Welsh S, Young A, Effingham M, McVean G, Leslie S, Allen N, Donnelly P, Marchini J. The UK Biobank resource with deep phenotyping and genomic data. Nature. 2018 Oct;562(7726):203-209.

13. Backman, J.D., Li, A.H., Marcketta, A. et al. Exome sequencing and analysis of 454,787 UK Biobank participants. Nature 599, 628–634 (2021).

14. Szustakowski, J.D., Balasubramanian, S., Kvikstad, E. et al. Advancing human genetics research and drug discovery through exome sequencing of the UK Biobank. Nat Genet 53, 942–948 (2021).

15. Wu MC, Lee S, Cai T, Li Y, Boehnke M, Lin X. Rare-variant association testing for sequencing data with the sequence kernel association test. Am J Hum Genet. 2011 Jul 15;89(1):82-93. doi: 10.1016/j.ajhg.2011.05.029. Epub 2011 Jul 7. PMID: 21737059; PMCID: PMC3135811.

16. Firth D (1993). Bias reduction of maximum likelihood estimates. Biometrika 80, 27-38. Heinze G, Schemper M (2002). A solution to the problem of separation in logistic regression. Statistics in Medicine 21: 2409-2419.

17. Storey, John D. (2002). "A direct approach to false discovery rates". Journal of the Royal Statistical Society, Series B (Statistical Methodology). **64** (3): 479–498.

18. Cinar, O. & Viechtbauer, W. (2022). The poolr package for combining independent and dependent p values. Journal of Statistical Software, 101(1), 1‑42. https://doi.org/10.18637/jss.v101.i01

19. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, Kondrashov AS, Sunyaev SR. A method and server for predicting damaging missense mutations. Nat Methods. 2010 Apr;7(4):248-9. doi: 10.1038/nmeth0410-248. PMID: 20354512; PMCID: PMC2855889.

19.20. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, Flicek P, Cunningham F. The Ensembl Variant Effect Predictor. Genome Biology. 2016 Jun;17(1):122. doi: 10.1186/s13059-016-0974-4

20.21. Kircher M, Witten DM, Jain P, O'Roak BJ, Cooper GM, Shendure J. A general framework for estimating the relative pathogenicity of human genetic variants. Nat Genet. 2014 Feb 2. doi: 10.1038/ng.2892. PubMed PMID: 24487276.

21.22. Rentzsch P, Witten D, Cooper GM, Shendure J, Kircher M. CADD: predicting the deleteriousness of variants throughout the human genome. Nucleic Acids Res. 2018 Oct 29. doi: 10.1093/nar/gky1016. PubMed PMID: 30371827.

22.23. Li X, Li Z, Zhou H, Gaynor SM, Liu Y, Chen H, Sun R, Dey R, Arnett DK, Aslibekyan S, Ballantyne CM, Bielak LF, Blangero J, Boerwinkle E, Bowden DW, Broome JG, Conomos MP, Correa A, Cupples LA, Curran JE, Freedman BI, Guo X, Hindy G, Irvin MR, Kardia SLR, Kathiresan S, Khan AT, Kooperberg CL, Laurie CC, Liu XS, Mahaney MC, Manichaikul AW, Martin LW, Mathias RA, McGarvey ST, Mitchell BD, Montasser ME, Moore JE, Morrison AC, O'Connell JR, Palmer ND, Pampana A, Peralta JM, Peyser PA, Psaty BM, Redline S, Rice KM, Rich SS, Smith JA, Tiwari HK, Tsai MY, Vasan RS, Wang FF, Weeks DE, Weng Z, Wilson JG, Yanek LR, NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium, TOPMed Lipids Working Group, Neale BM, Sunyaev SR, Abecasis GR, Rotter JI, Willer CJ, Peloso GM, Natarajan P, and Lin X. Dynamic incorporation of multiple in silico functional annotations empowers rare variant association analysis of large whole-genome sequencing studies at scale. Nature Genetics 2020; 52(9): 969-983. PMID: 32839606. DOI: 10.1038/s41588-020-0676-4.

Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, Kondrashov AS, Sunyaev SR. Nat Methods 7(4):248-249 (2010).

24. Ng PC, Henikoff S. SIFT: Predicting amino acid changes that affect protein function. Nucleic Acids Res. 2003 Jul 1;31(13):3812-4. doi: 10.1093/nar/gkg509. PMID: 12824425; PMCID: PMC168916.

23.25. Dennemärker J, Lohmüller T, Müller S, Aguilar SV, Tobin DJ, Peters C, Reinheckel T. Impaired turnover of autophagolysosomes in cathepsin L deficiency. Biol Chem. 2010 Aug;391(8):913-22.

24.26. Zhao, MM., Yang, WL., Yang, FY. et al. Cathepsin L plays a key role in SARS-CoV-2 infection in humans and humanized mice and is a promising target for new drug development. Sig Transduct Target Ther 6, 134 (2021).

25.27. Dong, Q., Li, Q., Duan, L. et al. Expressions and significances of CTSL, the target of COVID-19 on GBM. J Cancer Res Clin Oncol 148, 599–608 (2022). https://doi.org/10.1007/s00432-021-03843-9

26.28. Chauhan SS, Goldstein LJ, Gottesman MM. Expression of cathepsin L in human tumors. Cancer Res. 1991 Mar 1;51(5):1478-81.

27.29. Chen Q, Fei J, Wu L, Jiang Z, Wu Y, Zheng Y, Lu G. Detection of cathepsin B, cathepsin L, cystatin C, urokinase plasminogen activator and urokinase plasminogen activator receptor in the sera of lung cancer patients. Oncol Lett. 2011 Jul;2(4):693-699.

28.30.  Sudhan DR, Siemann DW. Cathepsin L targeting in cancer treatment. Pharmacol Ther. 2015 Nov;155:105-16.

29.31.  Kirschke H, Eerola R, Hopsu-Havu VK, Brömme D, Vuorio E. Antisense RNA inhibition of cathepsin L expression reduces tumorigenicity of malignant cells. Eur J Cancer. 2000 Apr;36(6):787-95.

30.32.  Krueger S, Kellner U, Buehling F, Roessner A. Cathepsin L antisense oligonucleotides in a human osteosarcoma cell line: effects on the invasive phenotype. Cancer Gene Ther. 2001 Jul;8(7):522-8.

31.33.  Levicar N, Dewey RA, Daley E, Bates TE, Davies D, Kos J, Pilkington GJ, Lah TT. Selective suppression of cathepsin L by antisense cDNA impairs human brain tumor cell invasion in vitro and promotes apoptosis. Cancer Gene Ther. 2003 Feb;10(2):141-51.

32.34.  Yang Z, Cox JL. Cathepsin L increases invasion and migration of B16 melanoma. Cancer Cell Int. 2007 May 8;7:8.

33.35.  Chen YC, Pohl G, Wang TL, Morin PJ, Risberg B, Kristensen GB, Yu A, Davidson B, Shih IeM. Apolipoprotein E is required for cell proliferation and survival in ovarian cancer. Cancer Res. 2005 Jan 1;65(1):331-7.

34.36.  Trost Z, Marc J, Sok M, Cerne D. Increased apolipoprotein E gene expression and protein concentration in lung cancer tissue do not contribute to the clinical assessment of non-small cell lung cancer patients. Arch Med Res. 2008 Oct;39(7):663-7.

35.37.  Luo J, Song J, Feng P, Wang Y, Long W, Liu M, Li L. Elevated serum apolipoprotein E is associated with metastasis and poor prognosis of non-small cell lung cancer. Tumour Biol. 2016 Aug;37(8):10715-21.

36.38.  Su WP, Chen YT, Lai WW, Lin CC, Yan JJ, Su WC. Apolipoprotein E expression promotes lung adenocarcinoma proliferation and migration and as a potential survival marker in lung cancer. Lung Cancer. 2011 Jan;71(1):28-33.

37.39.  Pencheva N, Tran H, Buss C, Huh D, Drobnjak M, Busam K, Tavazoie SF. Convergent multi-miRNA targeting of ApoE drives LRP1/LRP8-dependent melanoma metastasis and angiogenesis. Cell. 2012 Nov 21;151(5):1068-82.

38. Landrum MJ, Kattman BL. ClinVar at five years: Delivering on the promise. Hum Mutat. 2018 Nov;39(11):1623-1630.

39.1.   Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, Kondrashov AS, Sunyaev SR. Nat Methods 7(4):248-249 (2010).

40.1.   Ng PC, Henikoff S. SIFT: Predicting amino acid changes that affect protein function. Nucleic Acids Res. 2003 Jul 1;31(13):3812-4. doi: 10.1093/nar/gkg509. PMID: 12824425; PMCID: PMC168916.

**Figures**

**Fig 1. QQ plot of ILCCO WES study.**

The departure of the right tail from the 45 degree line represents the association signals from the study. (A) illustrates results using BF with KS prior. Under the null hypothesis (no association between genes and phenotype), $2\log BF_{ks} \sim \chi^2(3)$. (B) shows results using BF with SKAT prior. Similarly, $2\log BF_{SKAT} \sim \chi^2(3)$ under the null hypothesis.

**Supporting Information files**

S1 Text. Method Supplement.

S1 Table. Results of gene-based analysis using adjusted $BF_{KS}$ test in the discovery and replication.

S2 Table. Results of gene-based analysis using adjusted $BF_{SKAT}$ test in the discovery and replication.

S1 Figure. Population Structure shown in top 3 principal components.

S2 Figure. Relationship between QUAL and mean GQ vs. Ts/Tv ratio.

S3 Figure. Genetic region of rs112682750 (pos: 87727608, build 38) within *CTSL* gene.