

1 *Gene-gene interaction of AhR with and within the Wnt cascade affects suscepti-* 2 *bility to lung cancer*

3 Albert Rosenberger¹, Nils Muttray¹, Rayjean J Hung^{2,3}, David C Christiani⁴, Neil E Caporaso⁵, Geoffrey Liu^{6,7}, Stig E Bojesen^{8,9,10}, Loic Le Marchand¹¹,
4 Demetrios Albanes⁵, Melinda C Aldrich¹², Adonina Tardon¹³, Guillermo Fernández-Tardón¹³, Gad Rennert¹⁴, John K Field¹⁵, Michael P.A. Davies¹⁵,
5 Triantafillos Liloglou¹⁵, Lambertus A Kiemeny¹⁶, Philip Lazarus¹⁷, Bernadette Wendel¹, Aage Haugen¹⁸, Shanbeh Zienolddiny¹⁸, Stephen Lam¹⁹,
6 Matthew B Schabath²⁰, Angeline S Andrew²¹, Eric J Duell²², Susanne M Arnold²³, Gary E Goodman²⁴, Chu Chen²⁵, Jennifer A Doherty²⁶, Fiona Tay-
7 lor²⁷, Angela Cox²⁷, Penella J Woll²⁷, Angela Risch²⁸, Thomas R Muley^{29,30}, Mikael Johansson³¹, Paul Brennan³², Maria Teresa Landi⁵, Sanjay S Shete³³,
8 Christopher I Amos³⁴, Heike Bickeböller¹ on behalf of the INTEGRAL-ILCCO consortium

- 9 1. *Department of Genetic Epidemiology, University Medical Center, Georg-August-University Göttingen, Göttingen, Germany.*
- 10 2. *Lunenfeld-Tanenbaum Research Institute, Sinai Health System, University of Toronto, Toronto, Ontario, Canada.*
- 11 3. *Dalla Lana School of Public Health, University of Toronto, Toronto, Canada*
- 12 4. *Department of Environmental Health, Harvard T.H. Chan School of Public Health and Massachusetts General Hospital/Harvard Medical School,*
13 *Boston, Massachusetts, USA.*
- 14 5. *Division of Cancer Epidemiology and Genetics, National Cancer Institute, US National Institutes of Health, Bethesda, Maryland, USA.*
- 15 6. *Medical Oncology and Medical Biophysics, Princess Margaret Cancer Centre, Toronto, Ontario, Canada.*
- 16 7. *Medicine and Epidemiology, Dalla Lana School of Public Health, University of Toronto, Toronto, Ontario, Canada.*
- 17 8. *Department of Clinical Biochemistry, Herlev and Gentofte Hospital, Copenhagen University Hospital, Copenhagen, Denmark.*
- 18 9. *Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark.*
- 19 10. *Copenhagen General Population Study, Herlev and Gentofte Hospital, Copenhagen, Denmark.*
- 20 11. *Epidemiology Program, University of Hawaii Cancer Center, Honolulu, Hawaii, USA.*
- 21 12. *Department of Thoracic Surgery, Division of Epidemiology, Vanderbilt University Medical Center, Nashville, Tennessee, USA.*
- 22 13. *University of Oviedo, ISPA and CIBERESP, Faculty of Medicine, Oviedo, Spain.*
- 23 14. *Clalit National Cancer Control Center at Carmel Medical Center and Technion Faculty of Medicine, Haifa, Israel.*
- 24 15. *Roy Castle Lung Cancer Research Programme, The University of Liverpool, Department of Molecular and Clinical Cancer Medicine, Liverpool, UK.*
- 25 16. *Departments of Health Evidence and Urology, Radboud University Medical Center, Nijmegen, the Netherlands.*
- 26 17. *Department of Pharmaceutical Sciences, College of Pharmacy, Washington State University, Spokane, Washington, USA.*
- 27 18. *National Institute of Occupational Health, Oslo, Norway.*
- 28 19. *British Columbia Cancer Agency, Vancouver, British Columbia, Canada.*
- 29 20. *Department of Cancer Epidemiology, H. Lee Moffitt Cancer Center and Research Institute, Tampa, Florida, USA.*
- 30 21. *Department of Epidemiology, Geisel School of Medicine, Hanover, New Hampshire, USA.*
- 31 22. *Unit of Biomarkers and Susceptibility, Oncology Data Analytics Program, Catalan Institute of Oncology (ICO), Bellvitge Biomedical Research Insti-*
32 *tute (IDIBELL), Barcelona, Spain.*
- 33 23. *Markey Cancer Center, University of Kentucky, Lexington, Kentucky, USA.*
- 34 24. *Swedish Medical Group, Seattle, Washington, USA*
- 35 25. *Program in Epidemiology, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA.*
- 36 26. *Huntsman Cancer Institute, Department of Population Health Sciences, University of Utah, Salt Lake City, Utha, USA*
- 37 27. *Department of Oncology and Metabolism, University of Sheffield, Sheffield, UK.*
- 38 28. *University of Salzburg and Cancer Cluster Salzburg, Salzburg, Austria.*
- 39 29. *Translational Lung Research Center (TLRC) Heidelberg, Member of the German Center for Lung Research (DZL), Heidelberg, Germany.*
- 40 30. *Translational Research Unit, Thoraxklinik, University Hospital Heidelberg,, Heidelberg, Germany.*
- 41 31. *Department of Radiation Sciences, Umeå University, Umeå, Sweden.*
- 42 32. *International Agency for Research on Cancer, World Health Organization, Lyon, France.*
- 43 33. *Department of Biostatistics, Division of Basic Sciences, The University of Texas MD Anderson Cancer Center, Houston, Texas, USA*
- 44 34. *Dan L Duncan Comprehensive Cancer Center, Baylor College of Medicine, Houston, Texas, USA.*

45
46 **Corresponding author:** Albert Rosenberger, PhD
47 Universitätsmedizin Göttingen
48 Institut für Genetische Epidemiologie
49 Humboldtallee 32, 37073 Göttingen, Germany
50 Tel: +49 551 3914044;
51 Fax: +49 551 3914094;
52 Email: arosenb@gwdg.de
53 ORCID: 0000-0001-7848-1332

54 **Running title:** *AhR/Wnt genes in Lung Cancer*

55 **Funding** The National Institutes of Health (7U19CA203654-02/ 397 114564-5111078 Integrative Analysis of Lung
56 Cancer Etiology and Risk) supported this work. CARET is funded by the National Cancer Institute, National Institutes
57 of Health through grants U01 CA063673, UM1 CA167462, R01 CA 111703, RO1 CA 151989, U01 CA167462 and funds

58 from the Fred Hutchinson Cancer Research Center. The Boston Lung Cancer Study was funded by NCI grant
59 5U01CA209414. Other individual funding for participating studies and members of INTEGRAL-ILCCO are listed else-
60 where. [McKay JD, Hung RJ, Han Y, et al (2017) Nat Genet 49:1126–1132] The funders had no role in study design,
61 data collection and analysis, decision to publish, or preparation of the manuscript. We acknowledge support by the
62 Open Access Publication Funds of the Göttingen University.

63 **Conflict of interest** The authors declare no conflict of interest.

64

65 Abstract

66 Background: Aberrant *Wnt* signalling, regulating cell development and stemness, influences the development of
67 many cancer types. The Aryl hydrocarbon receptor (*AhR*) mediates tumorigenesis of environmental pollutants. Com-
68 plex interaction patterns of genes assigned to *AhR/Wnt*-signalling were recently associated with lung cancer suscep-
69 tibility. Aim: To assess the association and predictive ability of *AhR/Wnt*-genes with lung cancer in cases and controls
70 of European descent. Methods: Odds ratios (OR) were estimated for genomic variants assigned to the *Wnt* agonist
71 and the antagonistic genes *DKK2*, *DKK3*, *DKK4*, *FRZB*, *SFRP4* and *Axin2*. Logistic regression models with variable se-
72 lection were trained, validated and tested to predict lung cancer, at which other previously identified SNPs that have
73 been robustly associated with lung cancer risk could also enter the model. Further, decision trees were created to
74 investigate variant x variant interaction. All analyses were performed for overall lung cancer and for subgroups. Re-
75 sults: No genome-wide significant association of *AhR/Wnt*-genes with overall lung cancer was observed, but within
76 the subgroups of ever smokers (e.g. maker rs2722278 *SFRP4*; OR=1.20; 95%-CI: 1.13-1.27; p=5.6 10⁻¹⁰) and never
77 smokers (e.g. maker rs1133683 *Axin2*; OR=1.27; 95%-CI: 1.19-1.35; p=1.0 10⁻¹²). Although predictability is poor,
78 *AhR/Wnt-variants* are unexpectedly overrepresented in optimized prediction scores for overall lung cancer and for
79 small cell lung cancer. Remarkably, the score for never-smokers contained solely two *AhR/Wnt-variants*. The optimal
80 decision tree for never smokers consists of 7 *AhR/Wnt-variants* and only two lung cancer variants Conclusions: The
81 role of variants belonging to *Wnt/AhR*-pathways in lung cancer susceptibility may be underrated in main-effects
82 association analysis. Complex interaction patterns in individuals of European descent have moderate predictive ca-
83 pacity for lung cancer or subgroups thereof, especially in never smokers.

84 **Key words**: susceptibility, association, gene-gene integration, prediction, polygenic risk score, decision trees, never
85 smoker, small cell lung cancer

86 Background

87 Lung cancer (LC) is the most common cancer worldwide since 1985. It is the leading cause of cancer related death
88 around the world. [1] It was estimated for 2020, that globally 2.2 million new LC-cases were diagnosed, which are
89 11.4% of all new cancer cases. In the same year 1.8 million LC-cases died, which are 18% of all cancer related
90 deaths. [2] The lifetime risk of developing a clinical manifest lung cancer (from birth to age 74) is higher in men
91 (3.78%) than in women (1.77%).

92 The *Wnt* signalling pathway is a multi-regulator of e.g. cell proliferation, differentiation, genetic stability, and much
93 more. It is crucial in the development of embryos and in the dynamic balance of adult tissues, so also that of the
94 lung. With respect to LC, changes of the *Wnt* signalling pathway have been observed for *Wnt* ligands, frizzled,
95 TCF/LEF (T cell factor/lymphoid enhancer factor)-dependent transcription, and *Wnt* inhibitor silencing. [3]

96 Genome-wide association studies (GWAS) have identified dozens of susceptibility loci throughout the genome that
97 are associated with the susceptibility to lung cancer or one of its histological subtypes. [4–11] Genes related to *Wnt*
98 signalling, one of the key pathway regulating cell development and stemness, were not detected as being associated
99 to LC susceptibility in individuals of European descent so far, unlike *TERT* (5p15.33) that was one of the first for which
100 a robust association was observed. [12] Aberrant *Wnt* signalling is often observed in expression profiles of many
101 cancers, but to date no association of *Wnt/Ahr* genes with susceptibility to cancer of any type has been ob-
102 served. [13–15] Administration of RNAi against *Wnt* was shown to reduce tumour burden in lung adenocarcinoma
103 (adenoLC). [16] In non-small cell lung cancer (NSCLC), overexpressed *miR-582-3p* maintains stemness features by
104 negatively targeting the regulators of *Wnt* signalling *Axin2*, *DKK3* and *SRP1* for degradation, thereby increasing β -
105 catenin mediated *Wnt* activity. [17] *TERT* expression was found to be directly enhanced by binding of β -catenin to
106 its promoter region and thereby links telomerase activity to *Wnt* signalling. [13] This is inasmuch important, as *TERT*
107 is one of the first and most robust susceptibility genes for LC identified by GWAS. [18, 19] The tight regulatory ma-
108 chinery of the *Wnt* pathway has several major antagonists, such as Secreted Frizzled related protein (*sFRP*), Dickkopf
109 5 (*DKK*) protein and *Axin2* protein. [20] Evidence also exists for a crosstalk between *Ahr* and *Wnt* signalling. [21]

110 Aryl hydrocarbon receptor (7p21.1; *Ahr*) is a ligand induced transcription factor, which is translocated into the nu-
111 cleus. It is known to mediate the toxicity and tumorigenesis of a variety of environmental pollutants, including for

112 NSCLC. *AhR* upregulates the enzyme CYP1A1 when cells are exposed to carcinogenic metabolites, such as some pol-
113 ycyclic aromatic hydrocarbons (PAHs) found in cigarette smoke. The CYP1A1 coding gene is discussed as a suscepti-
114 bility gene for LC. *AhR* is a major determinant in the process of smoking driven LC. [22–24] The complexity of both
115 the *AhR* signalling pathway and the *Wnt* signalling cascade is reflected by interaction effects of genomic variants
116 within genes, which control their function. [25] Recently, the association of the *Wnt*-genes *DKK4* (8p11.21), *DKK3*
117 (11p15.3), *DKK2* (4q25), *FRZB* (2q32.1, also known as *sFRP3*), *SFRP4* (7p14.1), *Axin2* (17q24.1) and a potential inter-
118 action with *AhR* was investigated with respect to the susceptibility to LC in a sample of 600 subjects from North
119 India. [25, 26] A notable association with LC, e.g. for the *SFRP4* variant rs1802073 (OR=3.19; 95%-CI 1.81-5.63), was
120 reported. Classification And Regression Tree (CART) analysis revealed an interaction of *DKK2* and *SFRP4* polymor-
121 phisms to be the best (off all investigated) predictors for LC; especially within smokers. They also reported to have
122 identified several high-risk subgroups in smokers, e.g. characterised by *DKK2* (rs17037102 / rs419558) and *Axin2*
123 (rs9915936). A similar picture was observed in a sample of 270 subjects from Istanbul, Turkey. [27] A two-way inter-
124 action between *DKK3* (rs3206824) and *SFRP4* (rs1802074) was found to be predictive of LC.

125 We aimed to assess a possible association of *AhR* pathway and *Wnt* signalling cascade with LC within the large-scale
126 series of cases and controls of European descent hold by the International Lung Cancer Consortium (ILCCO) / Inte-
127 grative analysis of Lung Cancer Etiology and Risk (INTEGRAL). To do this, we also evaluated the contribution of these
128 genes to genetic prediction of LC as a complement to known LC-related markers.

129 Methods

130 The work presented has been reviewed and approved by the ILCCO Steering Committee.

131 Cases and Controls

132 Phenotype and genotype data of 58,181 entries of the data repository of ILCCO were extracted. Details of the repos-
133 itory is described previously. [4, 28] QC control samples, individuals without information on smoking status or age,
134 and samples of poor genotyping quality or sex discrepancies, were excluded. To avoid population stratification, this
135 analysis is focused on European-ancestry population (defined as more than 95% probability of being of European
136 descent). 14,068 incident LC-cases and 12,390 cancer-free controls of European descent remained for analysis. Those

137 genotyped with other genome-wide array in addition to OncoArray were separated to form an independent valida-
138 tion set (2nd validation set) of size (n=4,359, including 2,360 LC-cases and 1,999 controls).

139 Selected Markers

140 For this investigation we extracted the genotypes of 113 genomic variants (markers) assigned to 58 genes, previously
141 associated with the risk for LC in European decent people or one of its histological subtypes through a wide variety
142 of approaches [4–11] or proxies thereof (called *LC-marker*), and 296 markers assigned to 7 genes involved in *Wnt*
143 signalling and listed in Bahl et al. [25, 26] and Yilmaz et al. [27] (called *AhR/Wnt-marker*). Thus, we focused this anal-
144 ysis to genes previously investigated with respect to LC. Fifty of these 409 markers were eliminated before analysis
145 due to a MAF<1% (minor allele frequency), or departure from HWE (Hardy–Weinberg equilibrium) in genotypes
146 (unaffected $p<10^{-7}$, affected $p<10^{-12}$), or low imputation accuracy (info<0.8). Seventy-eight of the remaining *LC-mark-*
147 *ers* were genotyped with the OncoArray (44 thereof are proxy SNPs identified using LDlink [29]) and 32 had to be
148 imputed. Two-hundred twenty-one of the remaining *AhR/Wnt-markers* were genotyped and 28 have been imputed.
149 A list of these markers extracted from ILCCO OncoArray repository is given in the appendix.

150 Association analysis

151 We first performed association analysis for each marker separately using the program PLINK. [30, 31] Crude (model
152 1) and adjusted odds ratios (ORs) were estimated along with 95%-confidence intervals within log-additive models.
153 Sex, age and smoking status and the first 3 principal components (PCs) to adjust for population stratification (model
154 2); and in addition the 6 most significantly associated *LC-markers* (rs55781567, 15q25.1 *CHRNA5*; rs11780471,
155 8p21.2 *CHRNA2*; rs7705526, 5p15.33 *TERT*; rs56113850, 19q13.2 *CYP2A6*; rs71658797, 1p31.1 *AK5*; rs11571833,
156 13q13.1 *BRCA2*) (model 3) were included in adjusted models. ORs were estimated for overall LC, small cell LC (SCLC),
157 squamous cell LC (SqCLC), adenocarcinoma LC (adenoLC), ever smokers, never smokers and individuals aged
158 ≤ 55 years (early onset LC) as subgroups. We generated QQ-plots for the *AhR/Wnt-markers* and estimated the ge-
159 nomic inflation factor λ . To account for multiple testing, genome-wide statistical significance was considered to cor-
160 respond to a p-value of 10^{-7} or lower, suggestive significance to a p-value between 10^{-5} and 10^{-7} and nominal signif-
161 icance to a p-value between 0.05 and 10^{-5} .

162 Logistic Regression - Predicting models with model selection

163 We fitted logistic regression models with variable selection to find appropriate polygenic risk scores (PRS) in order
164 to predict the disease (LC) status (affected or unaffected). Any *AhR/Wnt-marker* or the *LC-marker* could be included
165 in the model without preference. To avoid multi-collinearity we removed one of two SNPs in LD to another ($R^2 > 0.8$,
166 pruning). The remaining entered the models as potential predictors. We performed forward selection until the
167 Bayesian information criterion (BIC, most stringent selection), the *Akaike* information criterion (AIC, less stringent
168 selection, contains in general more predictors) or the sample size corrected AIC (AICC) indicate a best solution (and
169 10 more selection steps). The resulting PRSs are called BIC-, AIC- and AICC-scores. Note, that for the purpose of
170 model building, the AIC-selection is asymptotically equivalent to cross-validation (CV). [32, 33] To avoid overfitting,
171 we assigned individuals to a training or a validation set (to build a score) and a testing set (to examine the score
172 performance) with a 1/3 probability each. For comparison, we also generated a BIC^{LC}-score with at least one marker,
173 only allowing *LC-markers* to enter the model building. To compare the importance for LC prediction of the sets g of
174 *LC-makers* and *AhR/Wnt-markers*, respectively, we contrasted the importance-values defined as $I_g = \sum_{m \in g} |\beta_m| \cdot$
175 MAF_m for each score (MAF_m the minor allele frequency and β_m the logistic regression coefficient of marker m). The
176 superiority of the AIC-scores over the BIC^{LC}-score and the BIC-score was tested applying the nonparametric test of
177 DeLong, DeLong, and Clarke-Pearson 1988 (1-sided) on AUCs of ROC (area under the receiver operation character-
178 istic curve). [34] In addition, a corresponding precision-recall plot was created for the SCLC.

179 Decision trees

180 Decision trees were created to examine marker x marker interaction with respect to the LC prediction. Any *AhR/Wnt-*
181 *marker* or the *LC-marker* could be included in a tree without preference. This was accomplished in the entire sample
182 and in all subgroups defined above. The R packages *rpart* and *DescTools* were used. [35, 36] To avoid trees being
183 formed by spurious epistasis we removed one of two SNPs in LD to another ($R^2 > 0.8$, pruning). Since overfitting is a
184 point of concern when building decision trees, the complexity parameter was first optimized applying 10-fold cross-
185 validation, grading the performance on the validation set by Somers' D (concordance of true and predicted LC-sta-
186 tus). The ability of the optimal trees to predict the LC-status was then tested within the independent sample of 4,359
187 cases and controls. True positive (TP) and true negative (TN) rates are given.

188 All statistical analyses were performed with SAS® 9.4, PLINK 1.90 and 2.0 or R 4.0.2.

189 Gene Expression

190 We extracted information on gene expression from the *Human Protein Atlas* [37, 38] and *LungGENS* [39, 40].

191 Results

192 Sample description

193 The analysed sample consists of 14,068 LC-cases and 12,390 controls with median age of 63. Sixty-three percent
194 were male, 52% of cases and 28% of controls were current smokers. The most frequent histological subtype is ade-
195 nocarcinoma (38%), followed by squamous cell carcinoma (SqCLC) (26%) and small cell lung cancer (SCLC) (10%). The
196 proportion of never-smokers was largest within the subgroup of adenocarcinoma cases (14%), but almost the same
197 between those cases aged ≤ 55 years (10%) and aged > 55 years (9%). Details on smoking status and histological
198 subtypes are presented in Table 1.

199 Table 1 Smoking by LC status and subgroups

200 Association analysis

201 We first performed association analysis for each *Wnt/AhR-marker* separately. The p-values for an association of
202 *AhR/Wnt-markers* with LC range from 0.005 (rs12115174; 8p11.21 *DKK4*; OR=0.9211) to 1 (model 2; adjusted for
203 sex, age, smoking status and population stratification); with a negligible genomic inflation ($\lambda=1.02$). A nominally sig-
204 nificant association ($10^{-5} < p \leq 0.05$) was observed for only 8 of the 249 markers (~3%). The corresponding point
205 estimates of OR range from 0.88 (rs1053070054; 8p11.21 *DKK4*; p=0.007) to 1.12 (rs74596148; 7p14.1 *SFRP4*;
206 p=0.25). A QQ-plot indicates that achieved p-values almost perfectly agree with the expectation of no associated
207 marker (see Figure 1). P-values and OR are in moderate agreement between the models (e.g. model 2 to model 3;
208 additionally adjusted by *LC-markers*: Kendall's $\rho_{p}=0.75$, $\rho_{OR}=0.78$).

209 Figure 1: Association of *AhR/Wnt-marker*

210 Subgroup analysis: When dividing the cases according to histological subtypes (SCLC; SqCLC and adenoLC) the ob-
211 servation of no detectable association for *WNT/AhR-markers* remains. Merely the number of nominally significant
212 association ($10^{-5} < p \leq 0.05$) increases to 12 (5%) or 21 (8%) of the 249 markers for SqCLC and SCLC, respectively,

213 hence close to the expected type 1 error. (Additional file 1: S-Table 2). When dividing the cases and controls accord-
214 ing to their smoking behaviour (ever and never smokers), genome-wide significance ($p \leq 10^{-7}$) was achieved for 7
215 and 8 markers, respectively. Another 12 and 3 markers, respectively, were found suggestively significant
216 ($10^{-7} < p \leq 10^{-5}$) (see Additional file 1: S-Figure 1) for ever and never smokers. Those markers found associated among
217 ever smokers have mainly been directly genotyped and are assigned to *SFRP4* and *DKK4*. E.g. for marker rs2722278
218 we estimated an OR=1.20 (95%-CI: 1.13-1.27), yielding a p-value of $5.6 \cdot 10^{-10}$. Those markers found associated among
219 never smokers have mainly been imputed and are mostly assigned to *Axin2*, but also to *AHR*, *FRZB* and *DKK2*. Marker
220 rs17037102, assigned to *DKK2*, was the only one found associated with LC by Bahl et al. and in this analysis (see
221 Table 2 and Additional file 1: S-Table 3). Interestingly, the ORs of these markers estimated by model 3 (additionally
222 adjusted for selected *LC-marker*) differ from that estimated by model 2. They are closer to one and no more signifi-
223 cant. E.g. for rs1133683 (*Axin2*) we observe an OR=1.27 (95%-CI: 1.19-1.35, $p=1 \times 10^{-12}$) fitting model 2, but OR=0.95
224 (95%-CI: 0.86-1.06, $p=0.3586$) fitting model 3.

225 Table 2 Significantly associated *AhR/Wnt-markers* within never and ever smokers

226 Logistic Regression - Predicting models with model selection

227 We further fit logistic regression models with variable selection to evaluate the contribution of *AhR/Wnt-markers* to
228 a polygenic risk scores (PRS), but without postulating the usefulness of the score as such. Eight *LC-markers* from only
229 eight *LC-genes* (*CYP2A6*, *CHRNA5*, *TERT*, *AMICA1*, *CHRNA3*, *COPS2*, *HCG4* and *CHRNA2*) were selected for the BIC-
230 score (most stringent selection) to predict overall LC. Hence, the BIC-score and the BIC^{LC}-score are identical. In con-
231 trast, the AIC-score (for overall LC identical to the AICC-score) includes 20 *LC-markers* and remarkable 17 *AhR/Wnt-*
232 *markers*, with *LC-markers* being more important than the *AhR/Wnt-markers* (importance ratio 0.56: 0.34) (see Figure
233 2, Additional file 1: S-Figure 3 and S-Table 4). The ability to distinguish cases and controls from susceptibility genes
234 only was, as expected, poor for each of the scores (see Additional file 1: S-Table 5). In the training set the perfor-
235 mance of the AIC/AICC-score (AUC=0.607) exceeded those of the BIC/BIC^{LC}-score (AUC=0.582) significantly
236 ($p < 0.001$). Within the test set (AUCs: 0.577 and 0.576) and the 2nd validation set (AUCs: 0.553 and 0.548), the higher
237 complexity with additional *AhR/Wnt-markers* did not improve discriminability for overall LC ($p=0.87$ and $p=0.35$).

238 Similar score composition and performance was observed for most subgroups. The BIC-scores in the subgroups
239 adenoLC (involved marker LC:*AhR/Wnt*=6:--), SCLC (3:--) and smokers (7:--) contained *LC-markers* only, whereas
240 *AhR/Wnt-markers* are included even under this stringent variable selection in the subgroups SqCLC (5:1) and Early
241 onset LC (2:2). However, between 14 and 31 *AhR/Wnt-markers* entered these subgroup's AIC-scores. For these sub-
242 groups, the importance of the LC-markers for the AIC-score is higher than that of the included *Ahr/Wnt-markers*.

243 Figure 2: Comparison of score composition

244 Most important, we observed a significantly higher predictive accuracy (larger AUCs) of the *AhR/Wnt-markers* en-
245 riched AIC-scores compared to BIC^{LC}-score in the subgroup of **SCLC** patients ($p=0.019$; $AUC_{AIC}=0.577$ $AUC_{BIC}=0.546$)
246 within the test set (see Additional file 1: S-Figure 4). For this subgroup, the selected *AhR/Wnt-markers* contribute to
247 the AIC-score more than twice as much as the *LC-markers* (importance ration 0.60: 1.49). The precision-recall plot
248 of Figure 3 indicates that a positive SCLC prediction based on the AIC-score can be trusted more than that based on
249 *LC-markers* alone (BIC^{LC}-score). In the 2nd validation set the score-specific AUCs were similar but no more significantly
250 different ($p=0.08$; $AUC_{AIC}=0.564$ vs. $AUC_{BIC}=0.531$). The AIC-score of this SCLC-subgroup is composed of 12 *LC-mark-*
251 *ers* (assigned to *CHRNA5*, *HCG4*, *DNAJB4* (4x each), *CYP2A6*, *CHRNA3*, *CHRNA2*, *AMICA1*, *KCNJ4*, *AS1*, *BRCA2*, *EGFL8*
252 and *WNK1* (2x each)) and 27 *AhR/Wnt-markers* (assigned to all *AhR/Wnt*-genes except *DKK3*). However, only one LC
253 patient in the test set ($n=434$) and one in the 2nd validation set ($n=164$) was recognized as a patient at a threshold of
254 50% case probability.

255 Figure 3: ROC and precision-recall-plot: SCLC

256 Interestingly the BIC-score for **never smokers** was built by only two *AhR/Wnt-markers* (assigned to *Axin2* and *SFRP4*)
257 but not a single LC-marker. Further, the *LC-markers* are the minority in the composite of the AIC-score (15:23). They
258 also contribute less to the AIC-score than the *AhR/Wnt-markers* (importance ratio of 0.96 : 1.46). The median pre-
259 dicted case probability, in the test set (24.8%) and 2nd validation set (25.6%), exceeds that of controls by 1%- to 2%-
260 points. However, AUC differed neither in the test set ($p=0.13$; $AUC_{AIC}=0.540$ $AUC_{BIC}=0.514$) nor in the 2nd validation
261 set ($p=0.36$; $AUC_{AIC}=0.535$ $AUC_{BIC}=0.526$) significantly. Nevertheless, this observation highlights the value of the
262 *AhR/Wnt-markers* in the subgroup of never smokers.

263 Decision trees

264 Finally, we generated decision trees to evaluate the contribution of *AhR/Wnt-markers* to LC prediction that allow for
265 a complex interaction structure, but without postulating the usefulness of the trees as such. The decision tree for
266 overall LC (whole sample) consists off solely a single decision node (rs55781567 assigned to *CHRNA5*), achieving a
267 Somers' concordance index $D=0.0565$ in the 2nd validation set (see Additional file 1: S-Table 6 and S-Figure 2). A
268 single-node decision-tree was also found optimal for participants aged ≤ 55 years (split: rs1051730 assigned to
269 *CHRNA3*), achieving a Somers' concordance index $D=0.096$. These two, unsophisticated trees are characterised by
270 balanced TP- (about 62%) and TN-rates (about 44%).

271 The decision trees for ever smokers, SCLC and SqCLC were more complex achieving Somers' concordance indexes D
272 of 0.007, -0.0005 and 0.0126, respectively. The trees for SCLC and SqCLC are characterised by an extreme TP-rate
273 $< 5\%$ and TN-rate $> 99\%$; the tree for Ever Smokers by a TP-rate $> 99\%$ and TN-rate $< 5\%$. Remarkably, a marker as-
274 signed to *CHRNA5* was always chosen as the first and most important split for the trees for ever smokers, for SCC
275 and SqCLC. However, markers assigned to *AhR/Wnt-genes* (smoker: *DKK2*; SCLC: *FRZB*; SqCLC; *DKK2* and *DKK3*) ap-
276 pear at lower-level decision-nodes (Additional file 1: S-Figure 5, 6, 7 and 8). With the same program settings, no
277 decision tree could be created for adenocarcinoma.

278 Most notable is the optimal decision tree for the 5,242 **never smokers** (75% LC-cases, 25% controls), the only one
279 that does not contain a marker belonging to the *CHRN* (*Cholinergic receptors nicotinic subunits*) gene group (see
280 Figure 4). The tree is built from only two *LC-markers* but 7 *AhR/Wnt-markers*, achieving a Somers' concordance index
281 $D=-0.002$. One can make out three branches of this tree. Branch I covers two thirds of individuals ($n=754$, 66% of
282 1141 in the 2nd validation set): All of these are graded as "unaffected" based on only the two *LC-markers*: first deci-
283 sion node (rs885518 assigned to *MTAP*) and second decision node (rs7705526 assigned to *TERT* that links telomerase
284 activity to *Wnt signalling*). For branch II an additional node (rs17214897 assigned to *DKK2*) is taken into account,
285 covering a further tenth (9.9%) of never smokers. In this branch, very few subjects of the training set (1.7% within
286 branch II eq. 0.17% of all never smokers) are graded "affected". However, one in four individuals of the 2nd validation
287 set belonging to both branches, I and II, is truly "affected" but has not been detected (TP-rate=0%, TN-rate=100%).
288 Rated as "affected" appears in the test set only in the third branch III, covering the remaining fourth of never smokers

289 (n=284 of the 2nd validation set). This third branch requires genotypes of several *AhR/Wnt-markers* assigned to *AHR*,
290 *Axin2*, *DKK2* and/or *SFRP4*. Herein, one in three (n=97 of the 2nd validation set) is truly “affected” and is given a
291 chance to be correctly identified, which appears in 8 LC-cases (TP-rate=9%, TN-rate=88%). We also noted that the
292 histological subtypes are equally distributed between the branches (see Additional file 1: S-Table 7).

293 Figure 4: Decision tree for never smoker

294 Gene Expression

295 *AHR*, *Axin2*, *DKK3* are ubiquitously expressed, with RNA expression detected in many tissues and evidence for protein
296 expression. *Axin2* and *DKK3* are moderately to highly expressed in normal lung tissues according to the Human Pro-
297 tein Atlas. [37] *AhR* is expressed at low levels in macrophage cells of the lung. No expression is reported for other
298 *Wnt/AhR*-genes. (see Additional file 1: S-Figure 9 and S-Table 9). Significant differential expression is listed in *Lung-*
299 *GENS* for *AhR*, *Axin2*, *DKK2*, *DKK3* and *SFRP4* [39] (see Additional file 1: S-Table 8). Further, *AhR* is reported to be
300 abundantly expressed in solid lung tumours, especially in adenocarcinomas. *AhR* overexpression was associated with
301 upregulation of IL-6 secretion, which is critical for lung cancer initiation. [41] Detailed information on gene expres-
302 sion is given in the Appendix. In addition, the *DKK1* serum level was seen as significantly lower in NSCLC and SCLC
303 patients compared to healthy controls. [42] Significant upregulation of *DKK2* expression was found in *APC* (adeno-
304 matous polyposis coli)-mutated non-SCLC lung cancers. [43]

305 Discussion

306 This investigation was intended to discover association of the *Wnt*-genes *DKK4* (8p11.21), *DKK3* (11p15.3), *DKK2*
307 (4q25), *FRZB* (2q32.1, also known as *sFRP3*), *SFRP4* (7p14.1), *Axin2* (17q24.1) and a potential interaction with *AhR-*
308 *genes*, to LC in a large sample of 26,458 individuals of European descent. No marginal association of *AhR/Wnt-mark-*
309 *ers* with overall LC was observed. Interestingly, an accumulation of associated markers was observed splitting the
310 sample by smoking status, where respective markers in ever smokers are assigned to *SFRP4*. On the other hand,
311 association analysis in never smokers reflects complex gene-gene interactions, as markers of several *AhR/Wnt-genes*
312 were found to be genome-wide associated with LC. This complexity is also visible through the decision tree analysis.

313 Our results are in line with findings from northern India [25, 26] and from Istanbul, Turkey [27], both of which are
314 based on much smaller samples (approx. 600 and 270 people, respectively). In these investigations, the interaction
315 of DKK2 and DKK3 with SFRP4 and Axin2 polymorphisms turned out to be the best (of all examined) predictors of LC,
316 especially in smokers. Axin2, but also AHR, FRZB and DKK2, were observed to be complex associated in never smok-
317 ers. Our analysis agrees with both previous studies that complex interaction patterns between the examined genes
318 contribute to overall LC susceptibility or within certain subgroups. However, we have not been able to replicate
319 reported single marker associations directly

320 To discover patterns of *Ahr/Wnt-genes* involved in LC genesis we further changed the focus from significance of
321 association to inclusion in prediction models, and followed two approaches: First, we searched for polygenic risk
322 scores (PRS). Doing so, we add up marker main effects to construct multidimensional scores, optimising model fit
323 (instead of marker preselection by p-value below some threshold), in order to discriminate cases from controls in a
324 somehow ideal way. Complex gene x gene (GxG) interactions are not modelled.

325 Nevertheless, the proportion of *Ahr/Wnt-genes* entering *some of the predictive* models was remarkable large, given
326 that these markers are not, all other candidates however genome-wide significantly associated to LC. This was par-
327 ticularly noticeable for SCLC, since *Ahr/Wnt-markers* contribute more than twice as much to the prediction score as
328 *LC-markers*. It is known, that within current smokers, tobacco consumption is strongest associated to SCLC. [44]
329 Moreover, within never smokers, a stringed defined score is made up from only two *Ahr/Wnt-markers*, assigned to
330 *Axin2* and *SFRP4*. However, the discriminative ability of PRSs for LC, contributing markers with significance for main
331 effect at different levels, is in general poor. The AUC of the BIC^{LC} score for overall LC (0.58 in the test set and 0.55 in
332 the 2nd validation set) corresponds to the AUC=0.54 based on four top *LC-genes* in a simulated population, as given
333 by the GWAS-ROCS Database (<https://gwasrocs.ca/>). This may be due to other overpowering risk factors, since mod-
334 els including e.g. age, sex and smoking variables achieve higher AUCs (0.62 to 0.79). [45]

335 Recently two polygenic risk scores (PRSs) for overall-LC had been developed, validated and assessed with respect to
336 improving eligibility to low-dose computed tomography (LDCT) as the only recommended screening test for lung
337 cancer. Jia et al. [46, 47] build a PRS on 19 genome-wide associated SNPs ($p < 0.5 \cdot 10^{-8}$). Hung et al. [48], integrated
338 their PRS on 128 SNPs (35 “known” LC-related loci, 93 suggestive associated loci selected by LASSO-regression

339 model) into the PLCO_{all2014} risk model. Both approaches have been validated using data from the UK Biobank. For
340 both scores, the mean PRS differed only slightly between LC cases and cancer-free controls (Jia: effect size ~ 0.19;
341 Hung: effect size ~ 0.22). For both scores, no substantial increase in discriminability of cases from controls is re-
342 ported, when adding the PRS to existing risk models (Jia: family history – AUC=0.589, family history + PRS –
343 AUC=0.615; Hung: PLCO_{all2014}–AUC=0.828, PLCO_{all2014}+ PRS – AUC=0.832). However, both were able to show that the
344 age at which a smoker crosses the recommended screening threshold of 1.5% for the 5-year LC risk depends on the
345 genetic background, which is sufficiently quantified by the PRS examined. Some smokers will be eligible by <50 years
346 of age, others by > 60 years of age. Hence, constructing reliable PRS, even with small discriminability, may help to
347 improve the performance of LDCT.

348 Two- and multiway GxG interaction can also contribute to LC susceptibility, rather than just markers with observed
349 (marginal) main effects. GxG interaction is in general less commonly investigated, not only because this requires
350 much larger samples. However, Li et al. [49] found RGL1:RAD51B in overall LC and non-SCLC, SYNE1:RNF43 in ade-
351 nocarcinoma and FHIT:TSPAN8 in SqCLC to interactively contribute to LC susceptibility. As in the presented data
352 analysis, the impact of these genes would also have been overlooked considering main effects only. Another reason
353 could be that LC itself is just a generic term of several subcategories that differ in terms of LC initiation and require
354 separate PRSs. [45, 50] A third reason of the poor performance may be due to the exclusively concentration on
355 genetic effects, rather than modelling lifelong interaction with the environment as well. E.g. GxE interaction effects
356 for LC have been observed smoking [51], exposure to asbestos fibres [52, 53] and exposure to radon [54, 55].

357 With this in mind, the data analysis presented shows that the complex interaction of Wnt-related genes has the
358 potential to be part of an adequate risk assessment for never-smokers or in relation to certain histological subtypes
359 of LC.

360 As a second approach, we constructed decision trees, which mainly depict GxG interaction patterns. Although, the
361 ability to discriminate cases from controls is again poor, CHRNA5 was in general the most important first node for
362 overall LC and in many subgroups. *Ahr/Wnt-genes* play a complex but important role in at least one quarter of never
363 smokers, as seen before. Remarkably, *TERT*, which links telomerase activity to *Wnt* signalling, was central in that
364 branch and important for the remaining three quarters of never smoker. This corresponds to a concentration of

365 relevant genes for this subgroup in the *CLPTM1L-TERT* region on chromosome 5, as previously reported by Hung et
366 al.. [56] Our observations confirm the suspicion, that LC in never smokers is a different entity, justified beforehand
367 on differences in epidemiological, clinical and molecular characteristics. [50]

368 We would like to emphasize that this study was not intended to provide a definitive and reliable risk assessment,
369 but rather aimed to examine in depth the LC-relevant complex interaction pattern of *AhR/Wnt-genes* hypothesized by
370 Bahl et al.. Indeed, considering prediction instead of association provides weaker evidence for this, but is valid in
371 view of the large amount of external evidence. The importance of the *Wnt*-signalling pathway and its antagonist's
372 *sFRP*, *DKKs* and *Axin2* for cancer is outlined in the introduction. One can also assume a connection with the molecular
373 functionality, since involved genes are expressed ubiquitously or in lung tissues.

374 Although the large-scale, thoroughly quality checked, and representative sample of genetically proven European
375 descent individuals was used for the presented analysis, some limitations must be noted. We used a rather narrow
376 definition of *AhR/Wnt-genes* to limit the number of possible interactions. An extension to e.g. *EGFR*, *APC*, *FRAT2* or
377 the *CYP*-family would also be justified. We further could have chosen the random forest method as a more contem-
378 porary and robust approach than decision trees, but we would not be able to present our results so illustrative.
379 However, the sample size allowed subgroup analyses, whereby the special importance of *AhR/Wnt-genes* for SCLC
380 and never smokers could be shown.

381 Conclusions

382 The role of markers belonging to *Wnt* signalling and the *AhR* pathway in LC susceptibility may be underrated in main-
383 effects association analysis. Complex interaction patterns in individuals of European decent have moderate predic-
384 tive capacity for LC or subsets thereof, especially in never smokers.

385 List of abbreviations

386	AhR	Aryl hydrocarbon receptor
387	GWAS	genome-wide association studies
388	LC	lung cancer

389	NSCLC	non-small cell lung cancer
390	SCLC	small cell lung cancer
391	SqCLC	squamous cell lung cancer
392	adenoLC	adenocarcinoma lung cancer
393	OR	odds ratio
394	CART	classification and regression tree
395	AUCs of ROC	area under the receiver operation characteristic curve
396	ILCCO	International Lung Cancer Consortium
397	INTEGRAL	Integrative analysis of Lung Cancer Etiology and Risk
398	PRS	polygenic risk scores
399	BIC	Bayesian information criterion
400	AIC	Akaike information criterion
401	CV	cross validation
402	MAF	minor allele frequency
403	TP	true positive rate
404	TN	true negative rate
405	LDCT	low-dose computed tomography

406 **Declarations**

407 **Ethics approval and consent to participate**

408 All participants in this study signed an informed consent, approved by the local internal review board or ethics com-
409 mittee and administered by trained personnel. All consortium research received approval from the Dartmouth Com-
410 mittee for Protection of Human Subjects on 7/30/2014 with id STUDY00023602. All experimental protocols and
411 other methods used comply with institutional, national, or international guidelines.

412 **Consent for publication**

413 Not applicable

414 **Availability of data and materials**

415 The data that support the findings of this study are available from ILCCO/INTEGRAL but restrictions apply to the
416 availability of these data, which were used under license for the current study, and so are not publicly available. Data
417 are however available from the authors upon reasonable request and with permission of ILCCO/INTEGRAL.

418 **Competing interests**

419 The authors declare that they have no competing interests

420 **Funding**

421 The National Institutes of Health (7U19CA203654-02/ 397 114564-5111078 Integrative Analysis of Lung Cancer Eti-
422 ology and Risk) supported this work. CARET is funded by the National Cancer Institute, National Institutes of Health
423 through grants U01 CA063673, UM1 CA167462, R01 CA 111703, RO1 CA 151989, U01 CA167462 and funds from the
424 Fred Hutchinson Cancer Research Center. Other individual funding for participating studies and members of INTE-
425 GRAL-ILCCO are listed elsewhere [10, 30]. The funders had no role in study design, data collection and analysis,
426 decision to publish, or preparation of the manuscript.

427 **Authors' contributions**

428 A. R. designed the investigation, carried out parts of the formal analysis and wrote the main manuscript text. N.M.
429 carried out parts of the formal analysis, prepared figures and critical reviewed and revised the manuscript. B.W.
430 carried out parts of the formal analysis. R.J.H. and C.I.A. coordinate the research activity of the consortium, including

431 data curation and funding acquisition. H.B: supervised the investigation, including funding acquisition. R.J.H., H.B.,
 432 L.L.M., C.I.A. and L.A.K. critical reviewed and revised the manuscript. A.C., A.H., A.R., A.S.A., A.T., C.C., D.A., D.C.C.,
 433 E.J.D., F.T., G.E.G., G.F.-T., G.L., G.R., H.B., J.A.D., J.K.F., L.A.K., L.L.M., M.B.S., M.C.A., M.J., M.P.A.D., M.T.L., N.E.C.,
 434 P.B., P.J.W., P.L., R.J.H., S.E.B., S.L., S.M.A., S.S.S., S.Z., T.L. and T.R.M. collected and provided study materials and
 435 data.

436 Acknowledgements

437 Not applicable

438 Supplementary Information

439 “Additional file 1.pdf” contains additional information on a) which marker were extracted from ILCCO OncoArray
 440 repository, b) single marker association, c) Polygenic Risk Scores (PRS) and Decision Trees and d) gene expression in
 441 normal tissue.

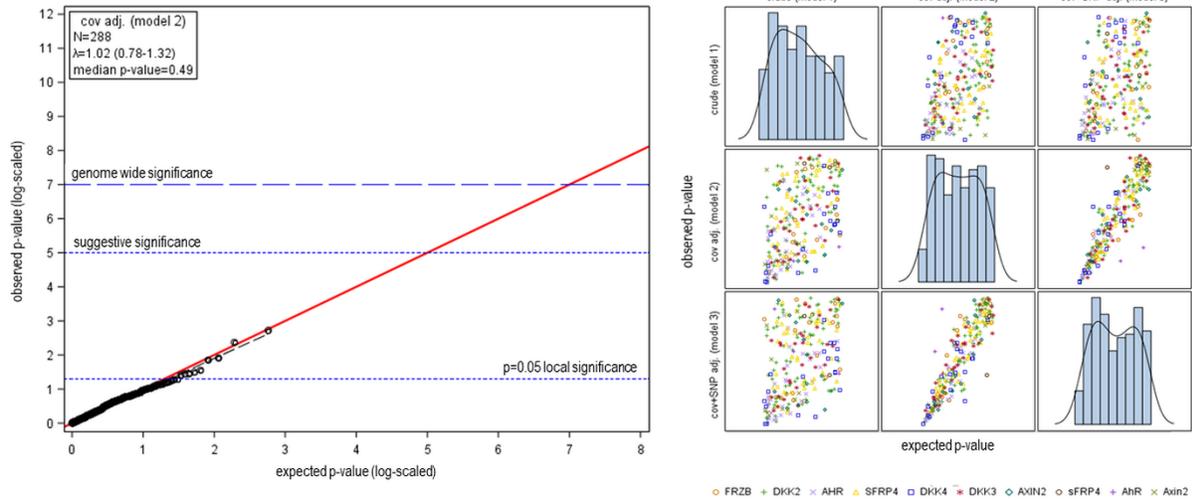
442 Tables and figures

443 **Table 1 Smoking by LC status and subgroups**

		Total N	Never smoker		former		Ever smoker		ever [§]	
			never n	%	n	%	current n	%	n	%
control	age ≤ 55 yrs.	2,762	951	34%	698	25%	896	32%	217	8%
	age >55 yrs.	9,628	2,960	31%	3,572	37%	2,568	27%	528	5%
	all	12,390	3,911	32%	4,270	34%	3,464	28%	745	6%
case	SqCLC	3,692	138	4%	1,257	34%	2,158	58%	139	4%
	SCLC	1,450	48	3%	383	26%	965	67%	54	4%
	other LC	3,629	405	11%	1,200	33%	1,820	50%	204	6%
	AdenoLC	5,297	740	14%	1,989	38%	2,401	45%	167	3%
	age ≤ 55 yrs.	2,765	281	10%	452	16%	1,945	70%	87	3%
	age >55 yrs.	11,303	1,050	9%	4,377	39%	5,399	48%	477	4%
	all	14,068	1,331	9%	4,829	34%	7,344	52%	564	4%
	total	26,458	5,242	20%	9,099	34%	10,808	41%	1,309	5%

444 [§]..as recorded; SCLC: small cell lung cancer, SqCLC : squamous cell lung cancer, AdenoLC: adenocarcinoma of the lung, other LC: other histo-
 445 logical subtypes;

446 **Figure 1: Association of *AhR/Wnt*-marker**



447

448 Left panel: QQ-Plot for model 2 (adjusted for sex, age and smoking status and the first three principal components); right panel: matrix of p-
 449 values generated by model 1 (crude), model 2 (adjusted for sex, age and smoking status and the first three principal components) and model 3
 450 (additionally adjusted for 6 selected *LC*-markers), genome-wide significance: p-value ≤ 10⁻⁷, suggestive significance: 10⁻⁷ < p-value ≤ 10⁻⁵, nominal
 451 significance: 10⁻⁵ < p-value ≤ 0.05.

452 **Table 2 Significantly associated *AhR/Wnt*-markers within never and ever smokers**

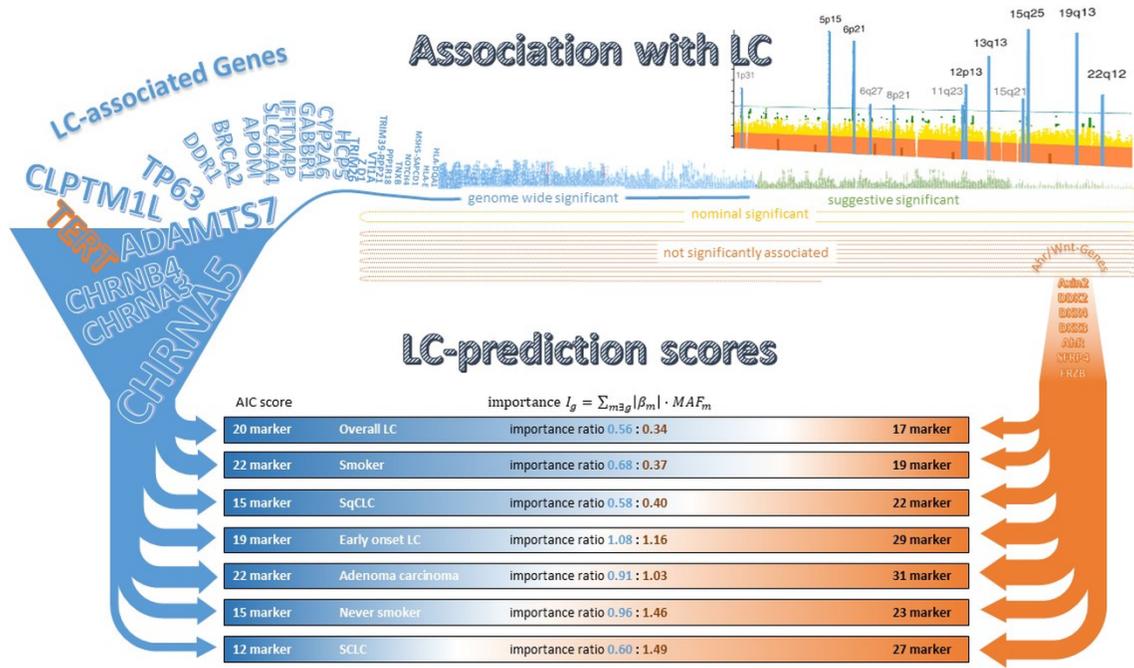
SNP	Cyto band	MAF	gene	model 2			model 1	model 3	
				p-value	OR	95%-CI	OR	OR	
never smoker									
imputed	rs202198518 [§]	7p21.1	14%	<i>AHR</i>	3.4 10 ⁻¹³	0.72	0.66-0.79	0.71	0.90 n.s.
imputed	rs2237297 [§]		14%		9.9 10 ⁻¹⁴	0.71	0.65-0.78	0.71	0.90 n.s.
imputed	rs1133683	17q24.1	42%	<i>Axin2</i>	1.0 10 ⁻¹²	1.27	1.19-1.35	1.27	0.95 n.s.
imputed	rs2240307		5%		7.7 10 ⁻²⁴	0.41	0.34-0.49	0.40	0.62 n.s.
imputed	rs35285779 [§]		9%		3.2 10 ⁻²²	0.58	0.52-0.65	0.58	1.10 n.s.
imputed	rs35415678 [§]		9%		3.7 10 ⁻¹⁹	0.62	0.56-0.69	0.62	1.10 n.s.
imputed	rs288326	2q32.1	10%	<i>FRZB</i>	2.5 10 ⁻⁸	1.42	1.25-1.60	1.41	0.98 n.s.
imputed	rs17037102	4q25	15%	<i>DKK2</i>	7.4 10 ⁻¹⁵	0.69	0.63-0.76	0.69	1.09 n.s.
ever smoker									
genotyped	rs12532321	7p14.1	45%	<i>SFRP4</i>	1.3 10 ⁻⁹	1.14	1.09-1.19	1.15	1.13 s.s.
genotyped	rs7811872		36%		1.3 10 ⁻⁸	0.88	0.84-0.92	0.88	0.88 gw.s.
genotyped	rs10226308		42%		1.8 10 ⁻⁸	0.88	0.85-0.92	0.89	0.89 gw.s.
genotyped	rs10488617		42%		1.6 10 ⁻⁸	0.88	0.85-0.92	0.89	0.89 gw.s.
genotyped	rs2722278		16%		5.6 10 ⁻¹⁰	1.20	1.13-1.27	1.16	1.20 gw.s.
genotyped	rs2722279		11%		9.0 10 ⁻⁹	1.22	1.14-1.31	1.17	1.23 gw.s.
genotyped	rs7811420		43%		7.9 10 ⁻⁸	0.89	0.85-0.93	0.89	0.89 gw.s.
imputed	rs2073664	8p11.21	9%	<i>DKK4</i>	9.4 10 ⁻¹¹	1.20	1.14-1.27	1.15	1.08 s.s.

453 MAF: minor allele frequency; model 1: crude odds ratio (OR); model 2: adjusted for sex, age and smoking status and the first three principal
 454 components; model 3: OR additional adjusted for 6 selected *LC*-markers. ^{gw.s.} genome-wide significant (p-value ≤ 10⁻⁷); ^{s.s.} suggestive significant
 455 (10⁻⁷ < p-value ≤ 10⁻⁵); ^{n.s.} not significant (p>0.05). Only markers are listed for which genome-wide significance (p-value ≤ 10⁻⁷) was achieved.

456 ^{§, §} pair of markers in LD (R²>0.8 in Populations of European decent)

457 **Figure 2: Comparison of score composition**

458

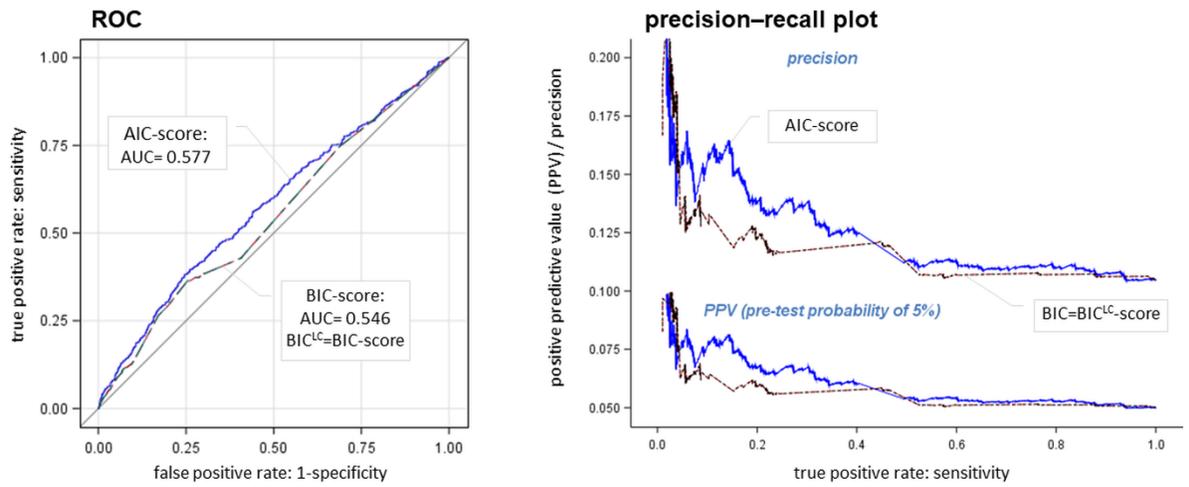


459

460 LC: lung cancer; AIC score: score of a logistic regression model with variant selection according to the Akaike information criterion (AIC); MAF_m:
 461 minor allele frequency of variant (marker) m; β_m regression parameter of variant m; LC-associated genes: previously reported as associated to
 462 LC or one of its histological subtypes; Ahr/Wnt-genes: selected genes assigned to Wnt-signalling, including AhR; Smoker: ever, former and cur-
 463 rent smoker; SCLC: small cell lung cancer, SqCLC : squamous cell lung cancer, Early onset LC: aged ≤55 years; TERT is framed in orange because
 464 telomerase activity is related to Wnt signalling.

465 **Figure 3: ROC and precision-recall-plot: SCLC**

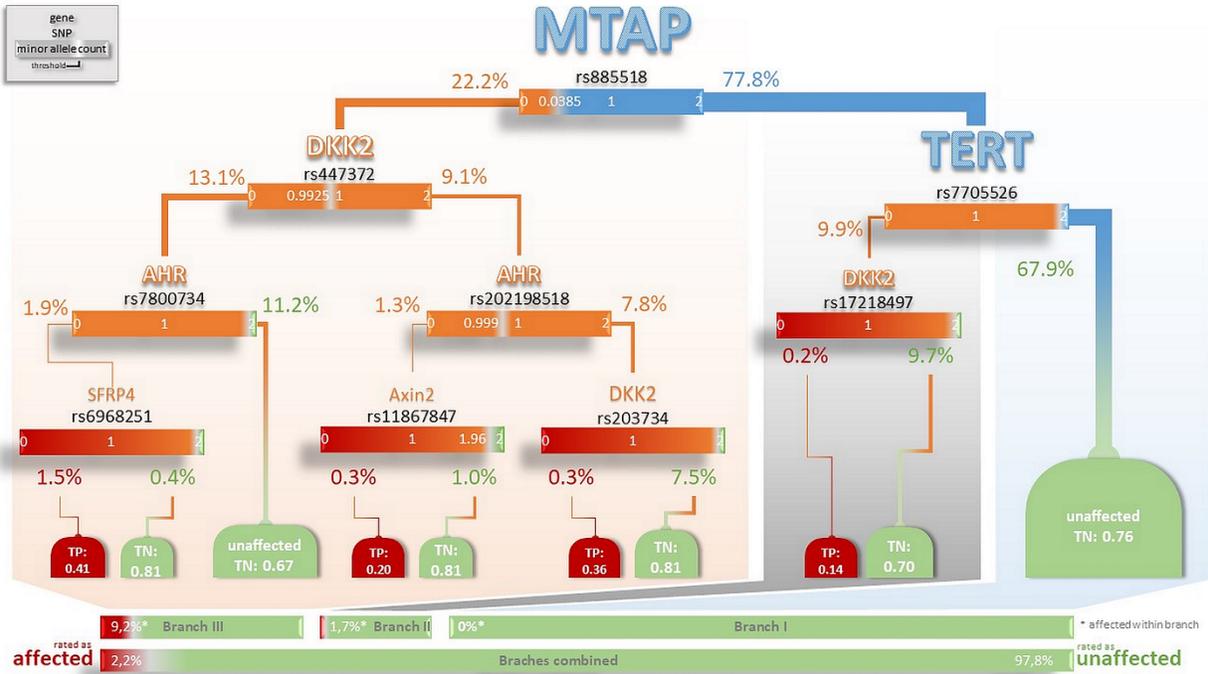
Small cell lung cancer (SCLC)



466

467 The diagnostic performance of the AIC-score compared to the BIC/BIC^{LC}-score in the test-set is presented. Left panel: ROC (receiver operation
468 characteristics); right panel: corresponding precision-recall plot; precision = (true positive cases) / (true positive cases + false positive controls),
469 positive predictive value (PPV) = (sensitivity x pre-test-probability) / [(sensitivity x pre-test-probability) + (1-specificity x 1-pre-test-probability)]
470 for a pre-test-probability of 5%.

471 **Figure 4: Decision tree for never smoker**



472

473 Node information: gene name, marker; split information below the node: threshold for minor allele count; blue split nodes: LC-genes, orange
 474 split nodes: Ahr/Wnt-genes, ; TERT is framed in orange because telomerase activity is related to Wnt signalling; decision nodes and bars: green
 475 for unaffected; red for affected, TN true negative rate, TP true positive red; the size of gene names, lines and decision notes is proportional to
 476 the size of the respective (sub)sample

477 References

- 478 1. Siegel RL, Miller KD, Jemal A (2016) Cancer statistics, 2016. *CA Cancer J Clin* 66:7–30
- 479 2. Ferlay J, Ervik M, Lam F, Colombet M, Mery L, Piñeros M, Znaor A, Soerjomataram I, Bray F Global Cancer
480 Observatory: Cancer Today. <https://gco.iarc.fr/today/data/factsheets/cancers/15-Lung-fact-sheet.pdf>. Ac-
481 cessed 1 Jan 2020
- 482 3. Zhu W, Wang H, Zhu D (2021) Wnt/ β -catenin Signaling Pathway in Lung Cancer. *Medicine in Drug Discovery*
483 100113
- 484 4. McKay JD, Hung RJ, Han Y, et al (2017) Large-scale association analysis identifies new lung cancer susceptibility
485 loci and heterogeneity in genetic susceptibility across histological subtypes. *Nat Genet* 49:1126–1132
- 486 5. Timofeeva MN, Hung RJ, Rafnar T, et al (2012) Influence of common genetic variation on lung cancer risk:
487 meta-analysis of 14 900 cases and 29 485 controls. *Human molecular genetics* 21:4980–95
- 488 6. Rosenberger A, Sohns M, Friedrichs S, et al (2017) Gene-set meta-analysis of lung cancer identifies pathway
489 related to systemic lupus erythematosus. *PLOS ONE* 12:e0173339
- 490 7. Brenner DR, Brennan P, Boffetta P, et al (2013) Hierarchical modeling identifies novel lung cancer susceptibility
491 variants in inflammation pathways among 10,140 cases and 11,012 controls. *Human genetics* 132:579–89
- 492 8. Ji X, Bossé Y, Landi MT, et al (2018) Identification of susceptibility pathways for the role of chromosome
493 15q25.1 in modifying lung cancer risk. *Nature Communications* 9:3221
- 494 9. Truong T, Sauter W, McKay JD, et al (2010) International Lung Cancer Consortium: coordinated association
495 study of 10 potential lung cancer susceptibility variants. *Carcinogenesis* 31:625–633
- 496 10. Wang Y, Wei Y, Gaborieau V, et al (2015) Deciphering associations for lung cancer risk through imputation and
497 analysis of 12 316 cases and 16 831 controls. *European Journal of Human Genetics* 23:1723–1728
- 498 11. Feng Y, Wang Y, Liu H, et al (2018) Novel genetic variants in the P38MAPK pathway gene ZAK and susceptibility
499 to lung cancer. *Mol Carcinog* 57:216–224
- 500 12. Landi MT, Chatterjee N, Yu K, et al (2009) A genome-wide association study of lung cancer identifies a region
501 of chromosome 5p15 associated with risk for adenocarcinoma. *American journal of human genetics* 85:679–
502 91
- 503 13. Zhan T, Rindtorff N, Boutros M (2016) Wnt signaling in cancer. *Oncogene* 36:1461–1473.
504 <https://doi.org/10.1038/onc.2016.304>
- 505 14. Buniello A, MacArthur JAL, Cerezo M, et al (2019) The NHGRI-EBI GWAS Catalog of published genome-wide
506 association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res* 47:D1005–D1012
- 507 15. (2020) GWAS catalog. In: The NHGRI-EBI Catalog of human genome-wide association studies.
508 <https://www.ebi.ac.uk/gwas/>. Accessed 22 Dec 2020
- 509 16. Kerdidani D, Chouvardas P, Arjo AR, et al (2019) Wnt1 silences chemokine genes in dendritic cells and induces
510 adaptive immune resistance in lung adenocarcinoma. *Nat Commun*. [https://doi.org/10.1038/s41467-019-
511 09370-z](https://doi.org/10.1038/s41467-019-09370-z)

- 512 17. Fang L, Cai J, Chen B, et al (2015) Aberrantly expressed miR-582-3p maintains lung cancer stem cell-like traits
513 by activating Wnt/beta-catenin signalling. *Nat Commun* 6:8640
- 514 18. Amos CI, Wu X, Broderick P, et al (2008) Genome-wide association scan of tag SNPs identifies a susceptibility
515 locus for lung cancer at 15q25.1. *Nat Genet* 40:616–622
- 516 19. Yuan Y, Lu C, Xue L, Ge D (2014) Association between TERT rs2736100 polymorphism and lung cancer suscep-
517 tibility: evidence from 22 case–control studies. *Tumor Biol* 35:4435–4442
- 518 20. Akiyama T (2000) Wnt/beta-catenin signaling. *Cytokine Growth Factor Rev* 11:273–82
- 519 21. Schneider AJ, Branam AM, Peterson RE (2014) Intersection of AHR and Wnt Signaling in Development, Health,
520 and Disease. *International Journal of Molecular Sciences* 15:17852–17885
- 521 22. Chang JT, Chang H, Chen PH, Lin SL, Lin P (2007) Requirement of aryl hydrocarbon receptor overexpression for
522 CYP1B1 up-regulation and cell growth in human lung adenocarcinomas. *Clin Cancer Res* 13:38–45
- 523 23. Lin P, Chang H, Tsai WT, Wu MH, Liao YS, Chen JT, Su JM (2003) Overexpression of aryl hydrocarbon receptor
524 in human lung carcinomas. *Toxicol Pathol* 31:22–30
- 525 24. Wang CK, Chang H, Chen PH, Chang JT, Kuo YC, Ko JL, Lin P (2009) Aryl hydrocarbon receptor activation and
526 overexpression upregulated fibroblast growth factor-9 in human lung adenocarcinomas. *Int J Cancer* 125:807–
527 15
- 528 25. Bahl C, Singh N, Behera D, Sharma S (2017) High-order gene interactions between the genetic polymorphisms
529 in Wnt and AhR pathway in modulating lung cancer susceptibility. *Personalized Medicine*.
530 <https://doi.org/10.2217/pme-2017-0018>
- 531 26. Bahl C, Singh N, Behera D, Sharma S (2017) Association of polymorphisms in Dickkopf (DKK) gene towards
532 modulating risk for lung cancer in north Indians. *Future Oncol* 13:213–232
- 533 27. Yilmaz M, Donmez G, Kacan T, Sari I, Akgül Babacan N, Sari M, Kilickap S (2015) Significant Association Between
534 Polymorphisms of Wnt Antagonist Genes and Lung Cancer: *Journal of Investigative Medicine* 1
- 535 28. Amos CI, Dennis J, Wang Z, et al (2017) The OncoArray Consortium: A Network for Understanding the Genetic
536 Architecture of Common Cancers. *Cancer Epidemiol Biomarkers Prev* 26:126–135
- 537 29. Machiela MJ, Chanock SJ (2015) LDlink: a web-based application for exploring population-specific haplotype
538 structure and linking correlated alleles of possible functional variants. *Bioinformatics* 31:3555–3557
- 539 30. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ (2015) Second-generation PLINK: rising to the
540 challenge of larger and richer datasets. *GigaScience* 4:1–16
- 541 31. Purcell S, Neale B, Todd-Brown K, et al (2007) PLINK: a tool set for whole-genome association and population-
542 based linkage analyses. *American journal of human genetics* 81:559–75
- 543 32. Stone M (1977) An Asymptotic Equivalence of Choice of Model by Cross-Validation and Akaike’s Criterion.
544 *Journal of the Royal Statistical Society Series B (Methodological)* 39:44–47
- 545 33. Fang Y (2009) Asymptotic Equivalence between Cross-Validations and Akaike Information Criteria in Mixed-
546 Effects Models. *JDS* 15–21

- 547 34. DeLong ER, DeLong DM, Clarke-Pearson DL (1988) Comparing the areas under two or more correlated receiver
548 operating characteristic curves: a nonparametric approach. *Biometrics* 44:837–845
- 549 35. Andri S (2021) DescTools: Tools for Descriptive Statistics.
- 550 36. Therneau T, Atkinson B, Ripley B (2019) rpart: Recursive partitioning for classification, regression and survival
551 trees.
- 552 37. Uhlén M, Fagerberg L, Hallström BM, et al (2015) Tissue-based map of the human proteome. *Science*.
553 <https://doi.org/10.1126/science.1260419>
- 554 38. The Human Protein Atlas. <https://www.proteinatlas.org/>. Accessed 5 Feb 2021
- 555 39. LungGENS. <https://research.cchmc.org/pbge/lunggens/>. Accessed 5 Feb 2021
- 556 40. Du Y, Guo M, Whitsett JA, Xu Y (2015) “LungGENS”: a web-based tool for mapping single-cell gene expression
557 in the developing lung. *Thorax* 70:1092–1094
- 558 41. Xue P, Fu J, Zhou Y (2018) The Aryl Hydrocarbon Receptor and Tumor Immunity. *Front Immunol*.
559 <https://doi.org/10.3389/fimmu.2018.00286>
- 560 42. Xu H, Wu J, Chen B, et al (2014) Serum Dickkopf-1 (DKK1) is significantly lower in patients with lung cancer but
561 is rapidly normalized after treatment. *Am J Transl Res* 6:850–856
- 562 43. Shen T, Chen Z, Qiao J, Sun X, Xiao Q (2019) Neutralizing monoclonal antibody against Dickkopf2 impairs lung
563 cancer progression via activating NK cells. *Cell Death Discovery* 5:1–12
- 564 44. Lee PN, Forey BA, Coombs KJ (2012) Systematic review with meta-analysis of the epidemiological evidence in
565 the 1900s relating smoking to lung cancer. *BMC cancer* 12:385
- 566 45. Katki HA, Kovalchik SA, Petito LC, Cheung LC, Jacobs J, Jemal A, Berg CD, Chaturvedi AK (2018) Implications of
567 Nine Risk Prediction Models for Selecting Ever-Smokers for Computed Tomography Lung Cancer Screening.
568 *Annals of internal medicine*. <https://doi.org/10.7326/M17-2701>
- 569 46. Jia G, Lu Y, Wen W, Long J, Liu Y, Tao R, Li B, Denny JC, Shu X-O, Zheng W (2020) Evaluating the Utility of
570 Polygenic Risk Scores in Identifying High-Risk Individuals for Eight Common Cancers. *JNCI Cancer Spectr*.
571 <https://doi.org/10.1093/jncics/pkaa021>
- 572 47. Jia G, Wen W, Massion PP, Shu X-O, Zheng W (2021) Incorporating Both Genetic and Tobacco Smoking Data
573 to Identify High-Risk Smokers for Lung Cancer Screening. *Carcinogenesis*. <https://doi.org/10.1093/carcin/bgab018>
- 575 48. Hung RJ, Warkentin MT, Brhane Y, et al (2021) Assessing Lung Cancer Absolute Risk Trajectory Based on a
576 Polygenic Risk Model. *Cancer Res* 81:1607–1615
- 577 49. Li Y, Xiao X, Bossé Y, et al (2019) Genetic interaction analysis among oncogenesis-related genes revealed novel
578 genes and networks in lung cancer development. *Oncotarget* 10:1760–1774
- 579 50. Sun S, Schiller JH, Gazdar AF (2007) Lung cancer in never smokers--a different disease. *Nature reviews Cancer*
580 7:778–90

- 581 51. Saccone NL, Culverhouse RC, Schwantes-An T-H, et al (2010) Multiple Independent Loci at Chromosome
582 15q25.1 Affect Smoking Quantity: a Meta-Analysis and Comparison with Lung Cancer and COPD. *PLOS Genet-*
583 *ics* 6:e1001053
- 584 52. Liu CY, Stucker I, Chen C, Goodman G, McHugh MK, D'Amelio AM Jr, Etzel CJ, Li S, Lin X, Christiani DC (2015)
585 Genome-wide Gene-Asbestos Exposure Interaction Association Study Identifies a Common Susceptibility Var-
586 iant on 22q13.31 Associated with Lung Cancer Risk. *Cancer Epidemiol Biomarkers Prev* 24:1564–73
- 587 53. Wei S, Wang L-E, McHugh MK, Han Y, Xiong M, Amos CI, Spitz MR, Wei QW (2012) Genome-wide gene-envi-
588 ronment interaction analysis for asbestos exposure in lung cancer susceptibility. *Carcinogenesis* 33:1531–1537
- 589 54. Rosenberger A, Hung RJ, Christiani DC, et al (2018) Genetic modifiers of radon-induced lung cancer risk: a
590 genome-wide interaction study in former uranium miners. *Int Arch Occup Environ Health*.
591 <https://doi.org/10.1007/s00420-018-1334-3>
- 592 55. Lorenzo-González M, Ruano-Ravina A, Torres-Durán M, et al (2019) Residential radon, genetic polymorphisms
593 in DNA damage and repair-related. *Lung Cancer* 135:10–15
- 594 56. Hung RJ, Spitz MR, Houlston RS, et al (2019) Lung Cancer Risk in Never-Smokers of European Descent is Asso-
595 ciated With Genetic Variation in the 5p15.33 TERT-CLPTM1L1 Region. *J Thorac Oncol* 14:1360–1369
- 596 57. Du Y, Kitzmiller JA, Sridharan A, et al (2017) Lung Gene Expression Analysis (LGEA): an integrative web portal
597 for comprehensive gene expression data analysis in lung development. *Thorax* 72:481–484