

Pleiotropic analysis of lung cancer and blood triglycerides identifies a shared genetic locus

Verena Zuber^{1,2,3,4}; Crystal N. Marconett⁵; Jianxin Shi⁶; Xing Hua⁶; William Wheeler⁷; Chenchen Yang⁵; Lei Song⁶; Anders M. Dale^{8,9,10,11}; Marina Laplana¹²; Angela Risch^{12,13,14}; Aree Witoelar^{1,2}; Wesley K. Thompson¹⁵; Andrew J. Schork^{8,9,16}; Francesco Bettella^{1,2}; Yunpeng Wang^{1,2}; Srdjan Djurovic^{17,18}; Beiyun Zhou¹⁹; Zea Borok¹⁹; Henricus F.M. van der Heijden²⁰; Jacqueline de Graaf²⁰; Dorine Swinkels²¹; Katja K. Aben²²; James McKay²³; Rayjean J. Hung²⁴; Heike Bikeböllner²⁵; Victoria L. Stevens²⁶; Demetrius Albanes⁶; Neil E. Caporaso⁶; Younghun Han²⁷; Yongyue Wei²⁸; Maria Angeles Panadero²⁸; Jose I Mayordomo²⁹; David C. Christiani^{28,31}; Lambertus Kiemeny²⁰; Ole A. Andreassen^{1,2}; Richard Houlston³²; Christopher I. Amos²⁷; Nilanjan Chatterjee⁶; Ite A. Laird-Offringa⁵; Ian G. Mills^{3,33,34†}; Maria Teresa Landi^{6†}

† These authors contributed equally

Correspondent authors: Maria Teresa Landi (landim@mail.nih.gov) and Ian Mills (I.Mills@qub.ac.uk; ian.mills@ncmm.uio.no)

¹NORMENT, KG Jebsen Centre for Psychosis Research, Institute of Clinical Medicine, University of Oslo, Oslo, Norway

²Division of Mental Health and Addiction, Oslo University Hospital, Oslo, Norway

³Prostate Cancer Research Group, Centre for Molecular Medicine Norway, Nordic EMBL Partnership, University of Oslo and Oslo University Hospital, Oslo, Norway

⁴European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, UK.

⁵Departments of Surgery and of Biochemistry and Molecular Biology, Norris Comprehensive Cancer Center, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA

⁶Division of Cancer Epidemiology and Genetics, National Cancer Institute, National Institutes of Health, U.S. Public Health Service, Bethesda, MD 20892, USA

⁷Information Management Services, Inc.; Rockville, MD, 20852; USA

⁸Multimodal Imaging Laboratory, University of California at San Diego, La Jolla, CA, USA

⁹Center for Human Development, University of California at San Diego, La Jolla, CA, USA

¹⁰Department of Radiology, University of California, San Diego, La Jolla, CA, USA

¹¹Department of Neurosciences, University of California, San Diego, La Jolla, CA, USA

¹²Division of Epigenomics and Cancer Risk Factors, German Cancer Research Center (DKFZ), Heidelberg, Germany.

¹³Department of Molecular Biology, University of Salzburg, Salzburg, Austria.

¹⁴Translational Lung Research Center Heidelberg TLRC-H, Member of the German Center for Lung Research DZL, Heidelberg, Germany.

¹⁵Department of Psychiatry, University of California, San Diego, La Jolla, CA, USA

¹⁶Cognitive Sciences Graduate Program, University of California, San Diego, La Jolla, CA, USA

¹⁷Department of Medical Genetics, Oslo University Hospital, Oslo, Norway

¹⁸NORMENT, KG Jebsen Centre for Psychosis Research, Department of Clinical Science, University of Bergen, Bergen, Norway

¹⁹Will Rogers Institute Pulmonary Research Center and Division of Pulmonary, Critical Care and Sleep Medicine, Department of Medicine, and Norris Comprehensive Cancer Center, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA

²⁰Radboud University Medical Center, Radboud Institute for Health Sciences, Nijmegen, The Netherlands

²¹Radboud University Medical Center, Radboud Institute for Molecular Life Sciences, Nijmegen, The Netherlands

- ²²Netherlands Comprehensive Cancer Organization, Utrecht, The Netherlands
- ²³International Agency for Research on Cancer (IARC, World Health Organization (WHO)), Lyon, France.
- ²⁴Lunenfeld-Tanenbaum Research Institute of Mount Sinai Hospital, Toronto, Ontario, Canada.
- ²⁵Zentrum Informatik, Statistik und Epidemiologie, Universitätsmedizin Göttingen
- ²⁶Epidemiology Research Program, American Cancer Society, Atlanta, Georgia, USA.
- ²⁷Center for Genomic Medicine, Department of Community and Family Medicine, Geisel School of Medicine, Dartmouth College, Lebanon, New Hampshire, USA.
- ²⁸Department of Epidemiology and Environmental Health, Harvard School of Public Health
- ²⁹Division of Medical Oncology, Ciudad de Coria Hospital, Coria, Spain
- ³⁰Division of Medical Oncology, University Hospital, Zaragoza, Spain
- ³¹Division of Pulmonary/Critical Care, Department of Medicine Massachusetts General Hospital/Harvard Medical School
- ³²Division of Genetics and Epidemiology, Institute of Cancer Research, Sutton, Surrey, UK.
- ³³Department of Molecular Oncology, Institute of Cancer Research and Department of Urology, Oslo University Hospital, Oslo, Norway
- ³⁴Prostate Cancer UK/Movember Centre of Excellence for Prostate Cancer Research, Centre for Cancer Research and Cell Biology, Queen's University, Belfast, UK

Epidemiologically-related traits may share genetic risk factors and pleiotropic analysis could identify individual loci associated with these traits. Because of their shared epidemiological associations, we conducted pleiotropic analysis of genome-wide association studies of lung cancer (12,160 lung cancer cases and 16,838 controls) and cardiovascular disease risk factors (blood lipids from 188,577 subjects, type 2 diabetes from 148,821 subjects, body mass index from 123,865 subjects, and smoking phenotypes from 74,053 subjects). We found that 6p22.1 (rs6904596, *ZNF184*) was associated with both lung cancer ($P=5.5 \times 10^{-6}$) and blood triglycerides ($P=1.39 \times 10^{-5}$). We replicated the association in 6,097 lung cancer cases and 204,657 controls ($P=2.4 \times 10^{-4}$) and in 71,113 subjects with triglycerides data ($P=0.011$). rs6904596 reached genome-wide significance in lung cancer meta-analysis (odds ratio=1.15, $P_{\text{combined}}=5.2 \times 10^{-9}$). The large sample size provided by the lipid GWAS data and the shared genetic risk factors between the two traits contributed to the uncovering of a hitherto unidentified genetic locus for lung cancer.

Genetic heritability of lung cancer is estimated to be 14% [1], but only a few genetic risk loci have been identified to date in genome-wide association studies (GWAS) of lung cancer in Europeans [2]. Epidemiological studies have shown associations between lung cancer and cardiovascular disease (CVD) risk factors related to the metabolic syndrome [3,4]. There is also substantial evidence that lipid metabolism and innate immunity evolved from common pathways and consequently genes that influence lipid traits may also influence inflammation and subsequent cancer development [5-7]. Lung cancer is also well-known to be strongly associated with

tobacco smoking. Predicated on the hypothesis that investigating shared genetic risk factors across these traits could enhance the possibility to identify new genetic loci for lung cancer, we used quantile-quantile (Q-Q) plots [8] (Online Methods) to assess potential polygenic enrichment of SNPs associated with lung cancer given association with each CVD risk factor or smoking phenotypes (**Figure 1 and Supplementary Figure 1**).

The analysis was based on the TRICL consortium meta-analysis of lung cancer GWAS, including 12,160 lung cancer cases and 16,838 controls [2] (**Supplementary Table 1**); the meta-analysis data of blood lipids from the Global Lipids Genetics Consortium (GLGC, including genetic association with triglycerides (TG), and high and low density lipoproteins-cholesterol (HDL-C and LDL-C)) from 188,577 subjects [9], of Type 2 diabetes (T2D) from 148,821 subjects [10] and of body mass index (BMI) from 123,865 subjects [11] ; and the meta-analysis of cigarettes per day (CPD) and never vs. ever smoking data from the Tobacco, Alcohol and Genetics (TAG) consortium, including 74,053 subjects (**Supplementary Table 2**). The **Supplementary Materials** (available online) contain additional details on the contributing studies, statistical analyses and functional tests.

The Q-Q plots show enrichment between lung cancer and LDL-C and between lung cancer and TG blood lipid traits across multiple p -value thresholds up to 10^{-5} (**Figure 1A-B**) verified by an adaptive permutation procedure (**Supplementary Table 3**). In contrast, we observed no significant enrichment ($P < 0.001$) between lung cancer and HDL, BMI, T2D or smoking phenotypes (the analysis of smoking excluded the SNP markers mapping to chr15:78,686,690-79,231,478, which are known to be associated with lung cancer and smoking [12-13] (**Supplementary Table 3**,

Supplementary Figure 1). Thus we excluded these traits from further analysis.

Cross-phenotype associated loci between lung cancer and TG and between lung cancer and LDL-C were assessed by conjunction false discovery rate (FDR) [8] (**Supplementary Materials**). Because controlling FDR is heavily affected by the number of identified SNPs, we pruned SNPs in linkage disequilibrium (LD) ($r^2 > 0.8$) and excluded the major histocompatibility complex (MHC) (genomic position (hg 19): chr6:29,528,318-33,373,649 [14]), which harbors established lung cancer susceptibility SNPs and is known for long range LD. By controlling conjunction FDR, we identified one genetic locus at 6p22.1, rs6904596, A>G, Minor Allele Frequency in Caucasians=0.094, associated with both lung cancer and blood triglycerides (conjunction FDR=0.0124; $P=5.5 \times 10^{-6}$ for lung cancer; $P=1.39 \times 10^{-5}$ for TG (This locus and additional genetic loci shared between lung cancer and lipid traits are shown in **Supplementary Table 4** and **Supplementary Figures 2-6**). This locus remained significant also using different thresholds for pruning SNPs in LD (**Supplementary Table 5**).

We tested this SNP for replication in 6,097 lung cancer cases and 204,657 controls from deCODE, Harvard, Holland and Spain (**Supplementary Table 6**). This locus was replicated ($P_{\text{replication}}=2.4 \times 10^{-4}$) and attained genome-wide significance for lung cancer risk in the meta-analysis of discovery and replication data (two-sided $P_{\text{combined}}=5.2 \times 10^{-9}$, $P_{\text{heterogeneity}}=0.91$, **Table 1**). This SNP was also replicated in the association with TG in 71,113 independent samples from deCODE and Holland (two-sided $P=0.011$, $P_{\text{combined}}=1.34 \times 10^{-6}$, **Table 1**).

The SNP association with lung cancer was mostly driven by the squamous cell carcinoma subtype ($P=2.8 \times 10^{-5}$) and not adenocarcinoma ($P=0.06$, **Supplementary**

Table 7).

rs6904596 localizes to 6p22.1 (27,491,299 bp; hg19) and lies 50kb 5' of Zinc Finger Protein 184 (ZNF184). It shows expression-QTL in lung tissue [15] with HLA-DRB3 ($\beta=-6.79$, $P=1.10 \times 10^{-11}$). Additionally, rs7749305, located on chr6:[27,446,566](#) ($r^2=1$ with rs6904596 in HapMap 3 of Caucasian populations), shows suggestive regulatory functions. This SNP showed the strongest association with lung cancer, but was not genotyped in the Global Lipid Consortium GWAS. It lies within a DNaseI hypersensitive region in small airway epithelial cells (SAEC) and A549 adenocarcinoma cells (ENCODE) and lies in a region hypomethylated in primary Alveolar Epithelial Cells (AEC) from our laboratory (**Supplementary Figure 7**). rs7749305 alternate allele C appears to create ATF3 and HIF1A binding sites. Similar findings are evident in adipocytes (ENCODE), extending the pleiotropic association between lung cancer and lipid traits to their function in respective tissue types.

Our study emphasizes that pleiotropic analysis of GWAS data of epidemiologically-related traits can uncover hitherto unidentified genetic associations. Moreover, some GWAS of quantitative traits may be much larger than disease specific GWAS (like in the case of CVD risk factors vs. lung cancer), and thus may improve the likelihood to identify new loci for the disease with the smaller sample size.

Funding

This work was supported by the Intramural Research Program of the Division of Cancer Epidemiology and Genetics, National Cancer Institute, National Institutes of Health, DHHS, Bethesda, MD.

Transdisciplinary Research for Cancer of Lung (TRICL): National Institute of Health U19 CA148127-01 (PI: Amos), Canadian Cancer Society Research Institute (no. 020214, PI: Hung). The Environment and Genetics in Lung Cancer Etiology (EAGLE), Prostate, Lung, Colon, Ovary Screening Trial (PLCO), and Alpha-Tocopherol, Beta-Carotene Cancer Prevention (ATBC) studies were supported by the Intramural Research Program of the National Institutes of Health, National Cancer Institute (NCI), Division of Cancer Epidemiology and Genetics. ATBC was also supported by U.S. Public Health Service contracts (N01-CN-45165, N01-RC-45035, and N01-RC-37004) from the NCI.

PLCO was also supported by individual contracts from the NCI to the University of Colorado Denver (NO1-CN-25514), Georgetown University (NO1-CN-25522), the Pacific Health Research Institute (NO1-CN-25515), the Henry Ford Health System (NO1-CN-25512), the University of Minnesota, (NO1-CN-25513), Washington University (NO1-CN-25516), the University of Pittsburgh (NO1-CN-25511), the University of Utah (NO1-CN-25524), the Marshfield Clinic Research Foundation (NO1-CN-25518), the University of Alabama at Birmingham (NO1-CN-75022), Westat, Inc. (NO1-CN-25476), and the University of California, Los Angeles (NO1-CN-25404). The Cancer Prevention Study-II (CPS-II) Nutrition Cohort was supported by the American Cancer Society. Funding for the Lung Cancer and Smoking study was provided by National Institutes of Health (NIH), Genes, Environment and Health Initiative (GEI) Z01 CP 010200, NIH U01 HG004446, and NIH GEI U01 HG 004438. For the lung study, the GENEVA Coordinating Center provided assistance with genotype cleaning and general study coordination, and the Johns Hopkins University Center for Inherited Disease Research conducted genotyping. Harvard Lung Study: The Harvard Lung Cancer Study is supported by the US National Institutes of Health (R01 CA092824, P50 CA090578, and R01 CA074386).

The authors thank deCODE genetics for contributing GWAS data.

Global Lipids Genetics Consortium (GLGC): Data on the lipid traits were provided by the Global Lipids Consortium through their access portal (<http://csg.sph.umich.edu/abecasis/public/lipids2013/>). The full Consortium acknowledgements are available in the supplementary information of Willer *et al.*

Tobacco, Alcohol and Genetics (TAG) consortium: Data on the smoking traits were provided by the Tobacco, Alcohol and Genetics through their access portal (<http://www.broadinstitute.org/mpg/ricopili/>). The full Consortium acknowledgements are available in the supplementary information of *the Tobacco, Alcohol and Genetics*.

Lung eQTL study: The lung eQTL study at Laval University was supported by the Chaire de pneumologie de la Fondation JD Bégin de l'Université Laval, the Fondation de l'Institut universitaire de cardiologie et de pneumologie de Québec, the Respiratory Health Network of the FRQS, the Canadian Institutes of Health Research

(MOP - 123369), and the Cancer Research Society and Read for the Cure. Y. Bossé is the recipient of a Junior 2 Research Scholar award from the Fonds de recherche Québec – Santé (FRQS).

Lung meQTL study: The meQTL study based on the environment and Genetics in Lung Cancer Etiology (EAGLE) study was supported by the Intramural Research Program of the Division of Cancer Epidemiology and Genetics, National Cancer Institute, NIH, DHHS.

Epigenetics: This study was supported by NIH grant (1 R01 HL114094) to IAL-O and ZB, NIH grants (1 P30 H101258) and (R37HL062569-13) to ZB. ZB is supported by the Ralph Edgington Chair in Medicine. CNM was supported in part by the department of Surgery, USC. Generation of epigenetic data was supported in part by the Norris Comprehensive Cancer Center core grant, award number P30CA014089 from the National Cancer Institute.

This work utilized the computational resources of the NIH HPC Biowulf cluster (<http://hpc.nih.gov>).

Notes

The study funders had no role in the design of the study; the collection, analysis, or interpretation of the data; the writing of the manuscript; nor the decision to submit the manuscript for publication. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Cancer Institute or the National Institutes of Health. The authors have no conflicts of interest to disclose.

References

1. Hemminki K, Lonnstedt I, Vaittinen P, *et al.* Estimation of genetic and environmental components in colorectal and lung cancer and melanoma. *Genet Epidemiol* 2001;20(1):107-116.
2. Wang Y, McKay JD, Rafnar T, *et al.* Rare variants of large effect in BRCA2 and CHEK2 affect risk of lung cancer. *Nat Genet* 2014;46(7):736-41.
3. Braun S, Bitton-Worms K, LeRoith D. The link between the metabolic syndrome and cancer. *Int J Biol Sci* 2011;7(7):1003-15.
4. Kucharska-Newton AM, Rosamond WD, Schroeder JC, *et al.* HDL-cholesterol and the incidence of lung cancer in the Atherosclerosis Risk in Communities (ARIC) study. *Lung Cancer* 2008;61(3):292-300.
5. Kominsky DJ, Campbell EL, Colgan SP. Metabolic shifts in immunity and inflammation. *J Immunol* 2010;184(8):4062-8.
6. Shi J, Chatterjee N, Rotunno M, *et al.* Inherited variation at chromosome 12p13.33, including RAD52, influences the risk of squamous cell lung carcinoma. *Cancer Discov* 2012;2(2):131-9.
7. Yu H, Pardoll D, Jove R. STATs in cancer inflammation and immunity: a leading role for STAT3. *Nat Rev Cancer* 2009;9(11):798-809.
8. Andreassen OA, Zuber V, Thompson WK, *et al.* Shared common variants in prostate cancer and blood lipids. *Int J Epidemiol* 2014; 10.1093/ije/dyu090.
9. Global Lipids Genetics C, Willer CJ, Schmidt EM, *et al.* Discovery and refinement of loci associated with lipid levels. *Nat Genet* 2013;45(11):1274-83.

10. Morris AP, Voight BF, Teslovich TM, *et al.* Large-scale association analysis provides insights into the genetic architecture and pathophysiology of type 2 diabetes. *Nat Genet* 2012;44(9):981-90.
11. Speliotes EK, Willer CJ, Berndt SI, *et al.* Association analyses of 249,796 individuals reveal 18 new loci associated with body mass index. *Nat Genet* 2010;42(11):937-48.
12. Amos CI, Wu X, Broderick P, *et al.* Genome-wide association scan of tag SNPs identifies a susceptibility locus for lung cancer at 15q25.1. *Nat Genet* 2008;40(5):616-22.
13. Thorgeirsson TE, Geller F, Sulem P, *et al.* A variant associated with nicotine dependence, lung cancer and peripheral arterial disease. *Nature* 2008;452(7187):638-42.
14. Shiina T, Hosomichi K, Inoko H, *et al.* The HLA genomic loci map: expression, interaction, diversity and disease. *J Hum Genet* 2009;54(1):15-39.
15. Hao K, Bosse Y, Nickle DC, *et al.* Lung eQTLs to help reveal the molecular underpinnings of asthma. *PLoS Genet* 2012;8(11):e1003029.

Table 1. Association of rs6904596 at 6p22.1 with both lung cancer risk and blood triglycerides' levels

| Study | Lung cancer Cases/ Controls | OR | 95% CI | P-value | | Study | Triglycerides Individuals | $\beta\beta$ | 95% CI | P-value |
|------------------------------------------------|------------------------------------|-----------|---------------|------------------------------|--|--------------|----------------------------------|--------------------------------|----------------|-----------------------|
| TRICL | 12160/16838 | 1.15 | (1.08,1.21) | 5.50×10^{-6} | | GLGC | 188577 | 0.0244 | (0.013,0.035) | 1.39×10^{-5} |
| Replication | 6097/204657 | 1.16 | (1.07,1.25) | 2.4×10^{-4} | | Replication | 71113 | 0.0290 | (0.006,0.051) | 1.14×10^{-2} |
| deCODE | 3865/196658 | 1.17 | (1.06,1.30) | 3.05×10^{-3} | | deCODE | 66027 | 0.0200 | (-0.004,0.044) | 0.10 |
| Harvard | 984/ 970 | 1.18 | (0.93,1.50) | 0.171 | | Holland | 5086 | 0.0891 | (0.027,0.151) | 5.12×10^{-3} |
| Holland | 687/ 5158 | 1.15 | (0.96,1.37) | 0.119 | | | | | | |
| Spain | 561/ 1871 | 1.10 | (0.88,1.37) | 0.40 | | Combined | 259690 | 0.0253 | (0.015,0.035) | 1.34×10^{-6} |
| Combined P _Q - heterogeneity* | 18257/ 221495 | 1.15 | (1.10,1.21) | 5.2×10^{-9} 0.92 | | | | | | 0.23 |

* Heterogeneity of effect size across studies was evaluated using the Cochran's Q statistic

Figures

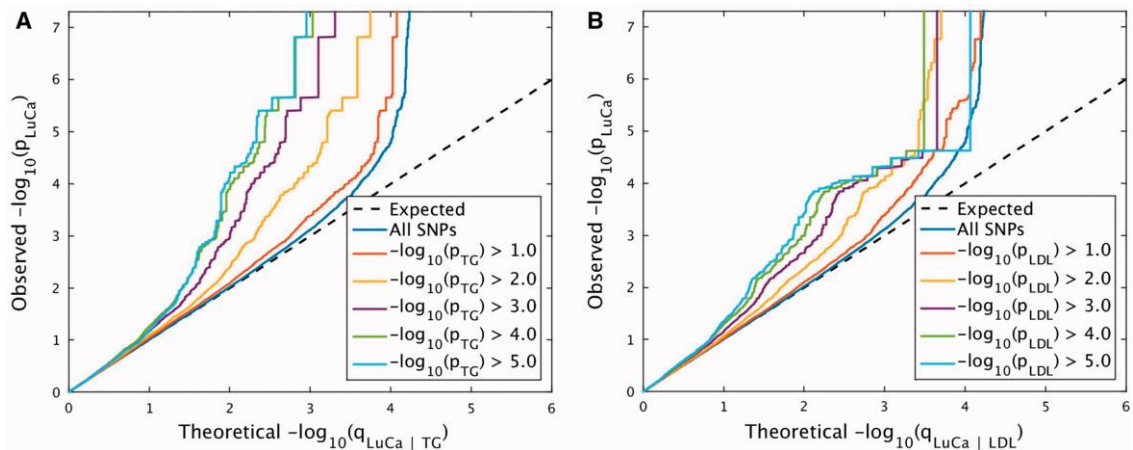


Figure 1. Conditional Q-Q plots: LuCa | CVD factors (TG and LDL-C).

‘Conditional Q-Q plot’ of theoretical vs empirical $-\log_{10} p$ -values (corrected for genomic control $\lambda\lambda$) in lung cancer (LuCa) below the standard GWAS threshold of $-\log_{10} p$ -values equal to 7.3 (equals p -values above 5×10^{-8}) as a function of significance of association with (A) triglycerides (TG) and (B) low-density lipoprotein (LDL-C) at the level of $p < 1$, $p < 0.1$, $p < 0.01$, $p < 0.001$, $p < 0.0001$, $p < 0.00001$ respectively. Dotted lines indicate the theoretical line in case of no association.

Supplementary Material

Supplementary Methods

Data and contributing studies

Input for the genetic epidemiology framework is summary statistics of genome-wide association studies (GWAS). Summary statistics on lung cancer were provided by the TRICL consortium (1) and were generated from a meta-analysis of 12,160 lung cancer cases and 16,838 controls. Further details on the sub-studies contributing to the meta-analysis are given in **Supplementary Table 1**. The histology-specific analyses were based on 3,718 adenocarcinoma (AD) cases and 3,422 squamous cell carcinoma (SQ) cases from the same TRICL consortium.

Data for the metabolic CVD risk factors was generated from meta-analyses of LDL, HDL, and TG (2), BMI (3), and T2D (4). Information on smoking behavior was measured by cigarettes per day (CPD) and never vs. ever smoking (SMOKER) in a large meta-analysis (5). For more details on these studies we refer to **Supplementary Table 2** including references and sample sizes.

We extracted summary statistics (P -values, risk alleles, ORs, and Z -scores) for 2,558,411 common SNPs created as a reference panel from the 1000Genomes data. There was overlap in samples between the lung cancer study and other traits: 2,282 with blood lipids study, 1,959 with the BMI study and 3,179 with the T2D study. We calculated the correlation of two Z -score statistics for each pair of traits due to sample overlap: $r_{LuCa,LDL}=-0.0034$, $r_{LuCa,TG}=-0.0054$, $r_{LuCa,HDL}=0.0170$, $r_{LuCa,BMI}=0.0176$, $r_{LuCa,T2D}=-0.0327$, $r_{LuCa,CPD}=-0.0371$, and $r_{LuCa,Smoker}=-0.0014$, suggesting that correlations due to sample overlap had a negligible impact on statistical inference even without explicit adjustment.

The analysis of lung cancer and cardiovascular risk factors included 483,841 independent individuals. The analysis of nicotine dependence included an additional 74,053

individuals. The replication datasets for lung cancer consisted of 6,097 lung cancer cases and 204,657 controls from deCODE, Harvard, Holland, and Spain; the replication for triglycerides and LDL-C consisted of 71,113 and 45,815 subjects, respectively from deCODE and Holland. Thus, we examined overall 885,576 individuals.

Enrichment analysis

We used conditional quantile-quantile (Q-Q) plots (6-9) after randomly pruning SNPs in linkage disequilibrium (LD) ($r^2 > 0.2$) for visualizing the polygenic enrichment patterns of lung cancer by restricting to the SNP sets that showed the strongest associations in a secondary trait (the CVD risk factors and smoking phenotypes). In particular, a conditional Q-Q plot was generated for SNPs with $P < 1E-1$, $1E-2$, $1E-3$, $1E-4$, and $1E-5$ for the secondary trait. In order to avoid confounding by large LD blocks, we performed a random pruning algorithm to compute the conditional Q-Q plot. Briefly, we defined LD blocks by an r^2 threshold of 0.8 and randomly selected one representative SNP from each LD block to compute a Q-Q plot. The final conditional Q-Q plot was generated by averaging the Q-Q plots from 100 random pruning steps.

We performed permutations to test whether high-ranking SNPs for the tested risk factors were enriched in lung cancer GWAS. Suppose there are N SNPs (denoted as S) after LD pruning. We aimed to test whether a specific set of N_1 SNPs (denoted as S_1) were enriched in lung cancer GWAS. The standard Kolmogorov-Smirnov test or Wilcoxon rank sum test was not sensitive to capture the deviation in the tail of the Q-Q plots. Thus, we designed an adaptive permutation procedure to test the deviation in the tails. We calculated the ranks of all N SNPs according to the association P values. We considered the top α proportion of SNPs in S_1 and calculated the total number of ranks of these SNPs, denoted as $T(\alpha)$. We performed permutations to approximate the statistical significance of $T(\alpha)$ as $p(\alpha)$. Because the test was sensitive to the choice of α , we chose a series of α values ($\alpha_1, \dots, \alpha_m$) and derived the corresponding P-values ($p(\alpha_1), \dots, p(\alpha_m)$). The overall one-sided statistic

for testing the enrichment was defined as $Q = \min(p(\alpha_1), \dots, p(\alpha_m))$ and its statistical significance was evaluated by permutations. In our data analysis we chose $(\alpha_1, \alpha_2, \alpha_3, \alpha_4) = (0.05, 0.1, 0.2, 0.3)$ and ran 10,000 permutations.

Assessing cross-phenotype association by conjunction FDR

In order to identify shared risk loci between two traits we used the conjunction false discovery rate (FDR), a genetic epidemiology framework based on an extension of the FDR into two dimensions. The conditional FDR is an extension of the standard FDR that allows including additional information on the p -value of the same SNP in a secondary trait 2. It is defined as the probability that a given SNP is null given that the p -values for trait 1 and trait 2 are as small or smaller than the observed ones. Low values of conditional FDR can be driven by the first trait only. To detect SNPs associated with both traits at the same time we used the conjunction FDR, which is defined as the probability of being null for either trait, or for both traits simultaneously given that the p -values for the two traits are as small or smaller than the observed ones. Thus, a true discovery is only the case when a SNP is non-null for both traits jointly. For more information on conditional and conjunction FDR and implementation we refer to (6, 8, 9). We aimed to control at a conservative FDR level of 0.05 per pair-wise comparison.

It has been shown that the estimation of the FDR can be impacted by correlation among the summary statistics (10). In order to address LD among SNPs we performed random pruning. First, we defined loci by an r^2 threshold of 0.8. Within each block we selected one SNP randomly by chance, where every SNP had equal chance of being selected. This random pruning procedure was repeated 100 times and the final computation of the empirical distribution was averaged over the 100 pruning procedures. Thus, the impact of large LD blocks was reduced on the estimation of the FDR.

eQTL in lung tissue

To identify expression quantitative trait loci we first used data from an analysis of genotype

and matched expression data from the Genotype-Tissue Expression Consortium (GTEx) (11) of n=111 healthy lung tissue samples. Additionally, we performed an analysis of genotype and matched expression data on healthy lung tissue in a study of n=1, 111 individuals using methodology previously described by Hao *et al.* (12).

meQTL in lung tissue

To identify methylation quantitative trait loci (meQTL) we applied a methodological approach on genotype and matched methylation data on n=210 non-tumor lung tissue samples as described previously by Shi *et al.*, (13), and in addition applied that analytical approach to a second dataset of non-tumor lung tissue samples obtained from The Cancer Genome Atlas (TCGA) lung cancer initiative (14).

Chromatin-level annotation

Publically available chromatin immunoprecipitation (ChIP)-seq data from the ENCODE and ROADMAP consortiums was used to determine relationship between SNPs and chromatin state in small airway epithelial cells (SAEC), adipocytes, and the A549 lung adenocarcinoma cell line, using the UCSC Genome browser, version hg19. Histone marks relating to active promoters were visualized with BigWig density tracks to determine if the identified loci were functionally active in the given cell types. Publically available ChIP-seq data on transcription factor binding was displayed using the UCSC genome browser.

To fully examine DNA methylation in the region surrounding the SNPs, we examined whole genome bisulfite sequencing (WGBS) data that we had obtained from purified primary alveolar epithelial cells (AEC) obtained as previously described (15, 16). Libraries were plated using the Illumina cBot and run on the Hi-Seq 2000 according to manufacturer's instructions using HSCS v 1.5.15.1. Rep 1 underwent Paired End 100 cycling; rep 2 underwent Paired End 75 cycling. Image analysis and base calling were carried out using RTA 1.13.48.0, deconvolution and fastq file generation was carried out using CASAVA_v1.7.1a5. Alignment to the genome was carried out using bsmap V 2.5. Aligned

.bam files were visualized using IGVviewer V2.3.40 (Broad Institute, Cambridge MA) with alignments colored by Bisulfite mode “CG”. AEC WGBS has been made publically available through GEO record GSE65319 along with complete sample preparation description.

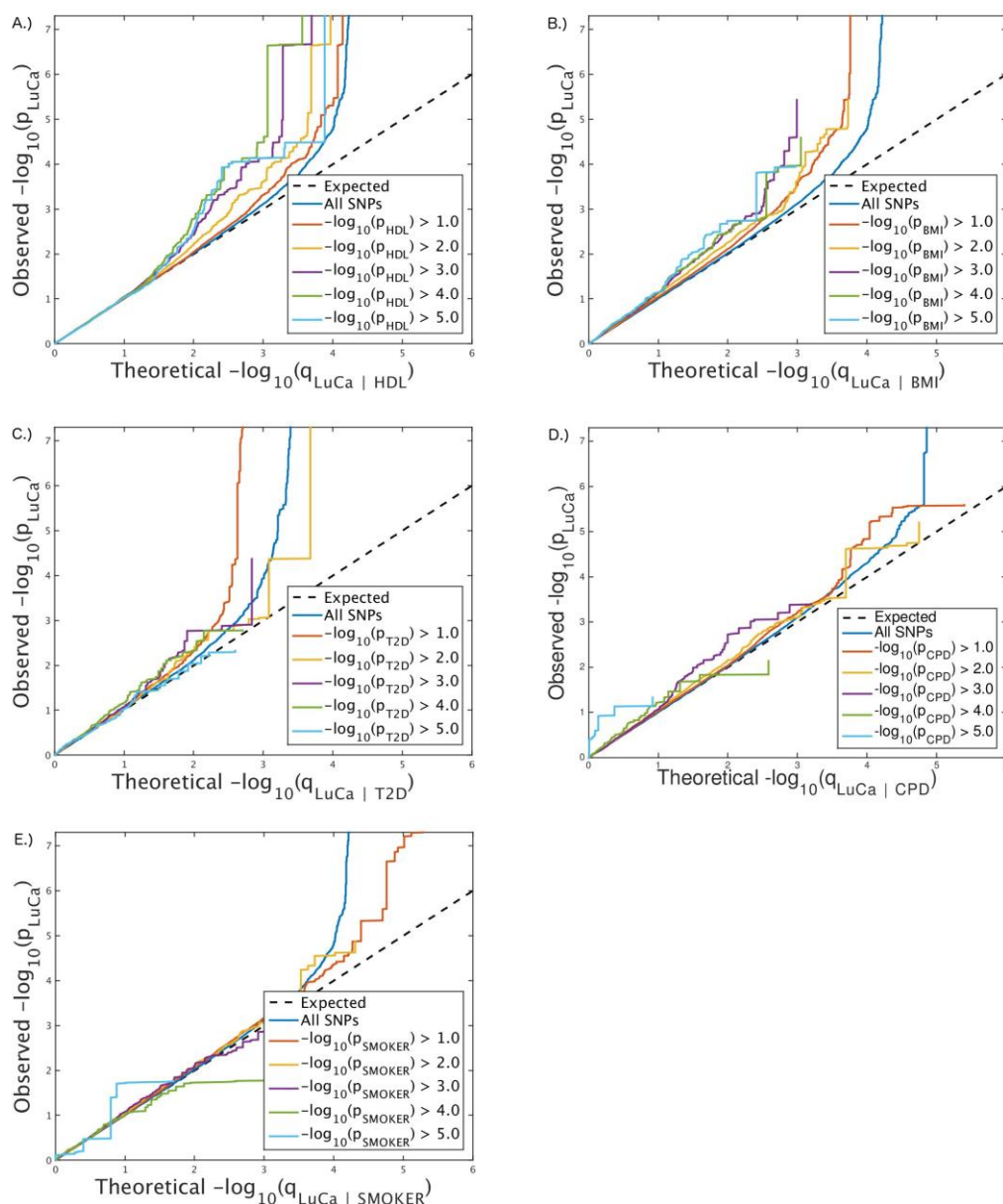
References

1. Wang Y, McKay JD, Rafnar T, *et al.* Rare variants of large effect in BRCA2 and CHEK2 affect risk of lung cancer. *Nat Genet* 2014;46(7):736-41.
2. Global Lipids Genetics C, Willer CJ, Schmidt EM, *et al.* Discovery and refinement of loci associated with lipid levels. *Nat Genet* 2013;45(11):1274-83.
3. Speliotes EK, Willer CJ, Berndt SI, *et al.* Association analyses of 249,796 individuals reveal 18 new loci associated with body mass index. *Nat Genet* 2010;42(11):937-48.
4. Morris AP, Voight BF, Teslovich TM, *et al.* Large-scale association analysis provides insights into the genetic architecture and pathophysiology of type 2 diabetes. *Nat Genet* 2012;44(9):981-90.
5. Tobacco, Genetics C. Genome-wide meta-analyses identify multiple loci associated with smoking behavior. *Nat Genet* 2010;42(5):441-7.
6. Schork AJ, Thompson WK, Pham P, *et al.* All SNPs Are Not Created Equal: Genome-Wide Association Studies Reveal a Consistent Pattern of Enrichment among Functionally Annotated SNPs. *PLoS Genet* 2013;9(4):e1003449.
7. Andreassen OA, McEvoy LK, Thompson WK, *et al.* Identifying Common Genetic Variants in Blood Pressure Due to Polygenic Pleiotropy With Associated Phenotypes. *Hypertension* 2014; 10.1161/HYPERTENSIONAHA.113.02077.
8. Andreassen OA, Thompson WK, Schork AJ, *et al.* Improved detection of common variants associated with schizophrenia and bipolar disorder using pleiotropy-informed conditional False Discovery Rate. *PLoS Genet* 2013;9(4):e1003455.
9. Andreassen OA, Zuber V, Thompson WK, *et al.* Shared common variants in prostate cancer and blood lipids. *Int J Epidemiol* 2014; 10.1093/ije/dyu090.
10. Schwartzman A, Lin X. The effect of correlation in false discovery rate estimation. *Biometrika* 2011;98(1):199-214.
11. Consortium GT. Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* 2015;348(6235):648-60.
12. Hao K, Bosse Y, Nickle DC, *et al.* Lung eQTLs to help reveal the molecular underpinnings of asthma. *PLoS Genet* 2012;8(11):e1003029.
13. Shi J, Marconett CN, Duan J, *et al.* Characterizing the genetic basis of methylome diversity in histologically normal human lung tissue. *Nat Commun* 2014;5:3365.
14. Cancer Genome Atlas Research N. Comprehensive molecular profiling of lung adenocarcinoma. *Nature* 2014;511(7511):543-50.
15. Ballard PL, Lee JW, Fang X, *et al.* Regulated gene expression in cultured type II cells of adult human lung. *American journal of physiology. Lung cellular and molecular physiology* 2010;299(1):L36-50.
16. Marconett CN, Zhou B, Rieger ME, *et al.* Integrated transcriptomic and epigenomic analysis of primary human lung epithelial cell differentiation. *PLoS genetics* 2013;9(6):e1003513.

Supplementary Figures

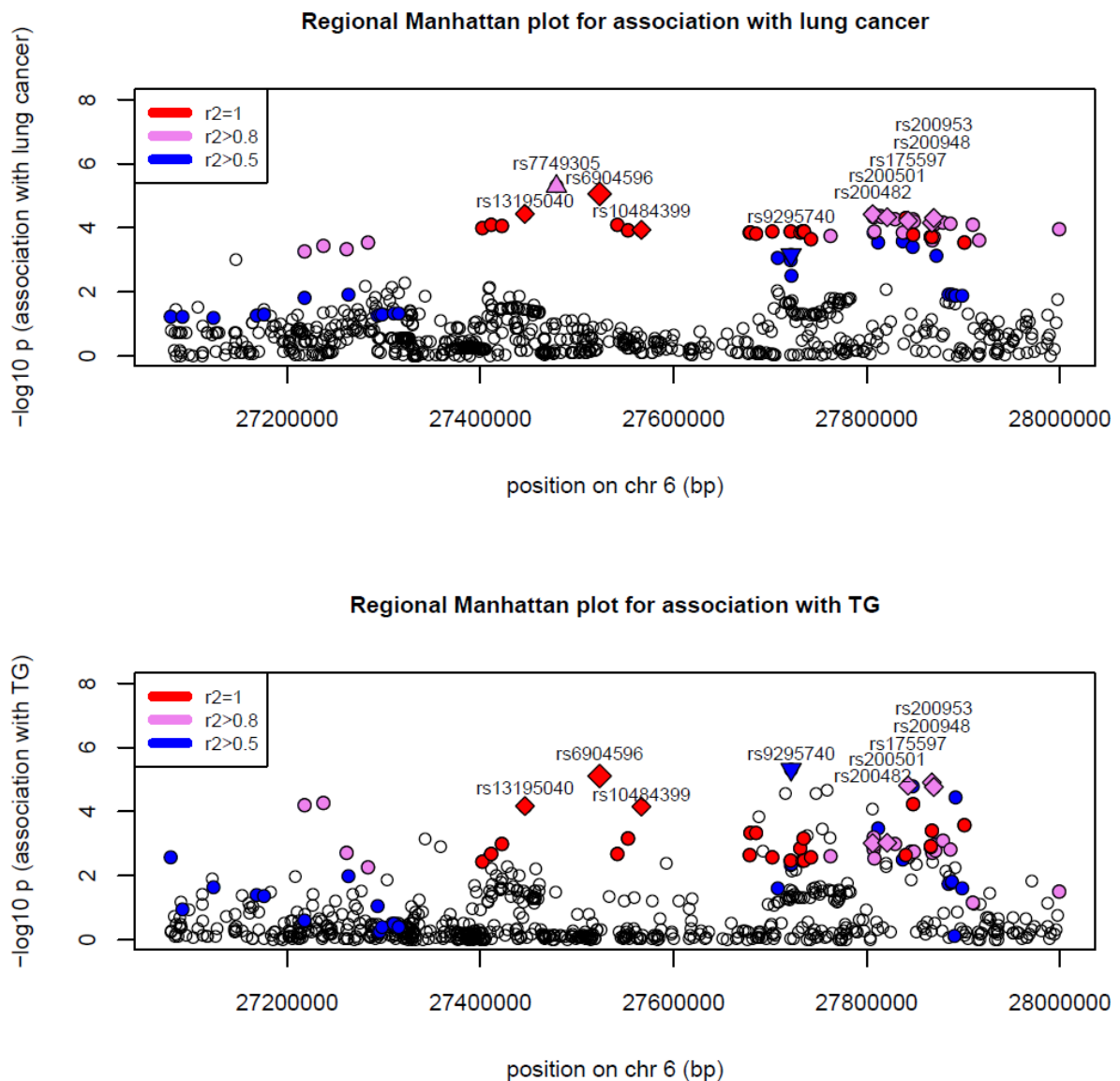
Supplementary Figure 1. Conditional Q-Q plots: LuCa | further CVD factors and smoking traits

‘Conditional Q-Q plot’ of theoretical vs empirical $-\log_{10} p$ -values (corrected for genomic control λ) in lung cancer (LuCa) below the standard GWAS threshold of $-\log_{10} p$ -values equal to 7.3 (equals p -values above 5×10^{-8}) as a function of statistical significance of association with (A) HDL cholesterol (HDL-C), (B) body mass index (BMI), (C) type 2 diabetes (T2D), (D) cigarettes per day (CPD), and (E) ever versus never smoking (SMOKER), at the level of $p < 1$, $p < 0.1$, $p < 0.01$, $p < 0.001$, $p < 1.00 \times 10^{-4}$, $p < 1.00 \times 10^{-5}$ respectively. Dotted lines indicate the theoretical line in case of no association. For the analysis of CPD we removed SNPs mapping to the nicotinic acetylcholine receptors (chr15:78,686,690-79,231,478) to exclude this well-established pleiotropic locus between LuCa and CPD. All statistical tests were two-sided.



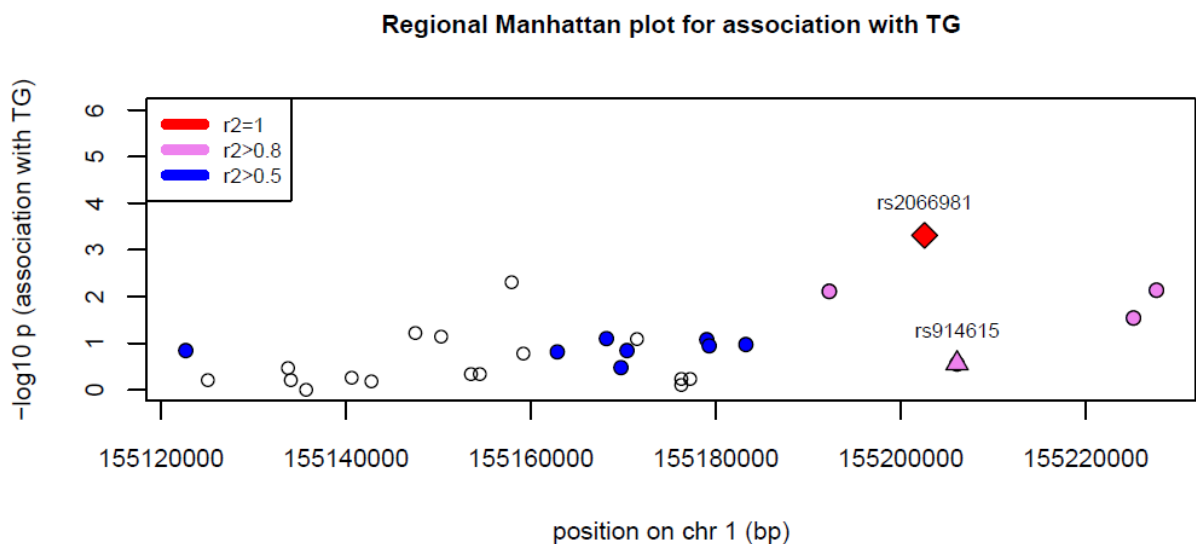
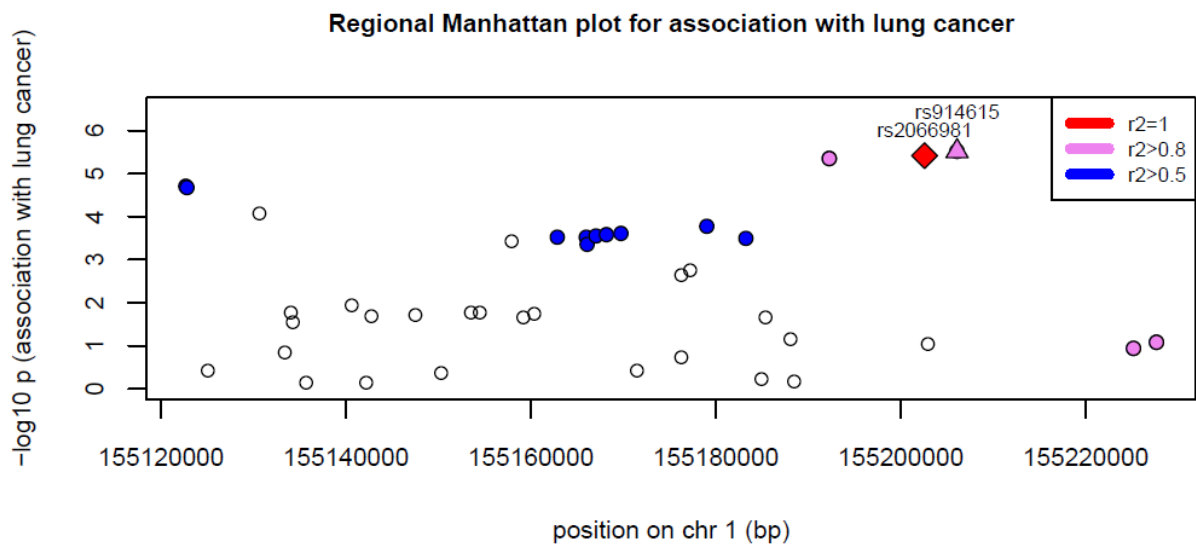
Supplementary Figure 2. Local Manhattan plot for cross-phenotype association of locus 6p22.1 between lung cancer and triglycerides (TG); top shared variant rs6904596

On display are the $\log_{10} p$ -values for the association of SNPs with lung cancer (upper panel) and TG (lower panel) based on the SNPs genomic position. All variants with a shared association between lung cancer and TG (conjunction FDR < 0.05) are symbolized with diamonds. The leading variant for lung cancer (rs7749305) is symbolized with a triangle (top up) and the leading variant (rs9295740) for TG with a triangle top down. rs7749305 was not included in the Global Lipid Consortium GWAS. Linkage disequilibrium (LD) based on HapMap v.3 in Caucasians is color-coded: variants with $r^2=1$ with the lead variant rs6904596 are colored in red, $0.8 \leq r^2 < 1$ are colored in violet, and $0.5 \leq r^2 < 0.8$ are colored in blue.



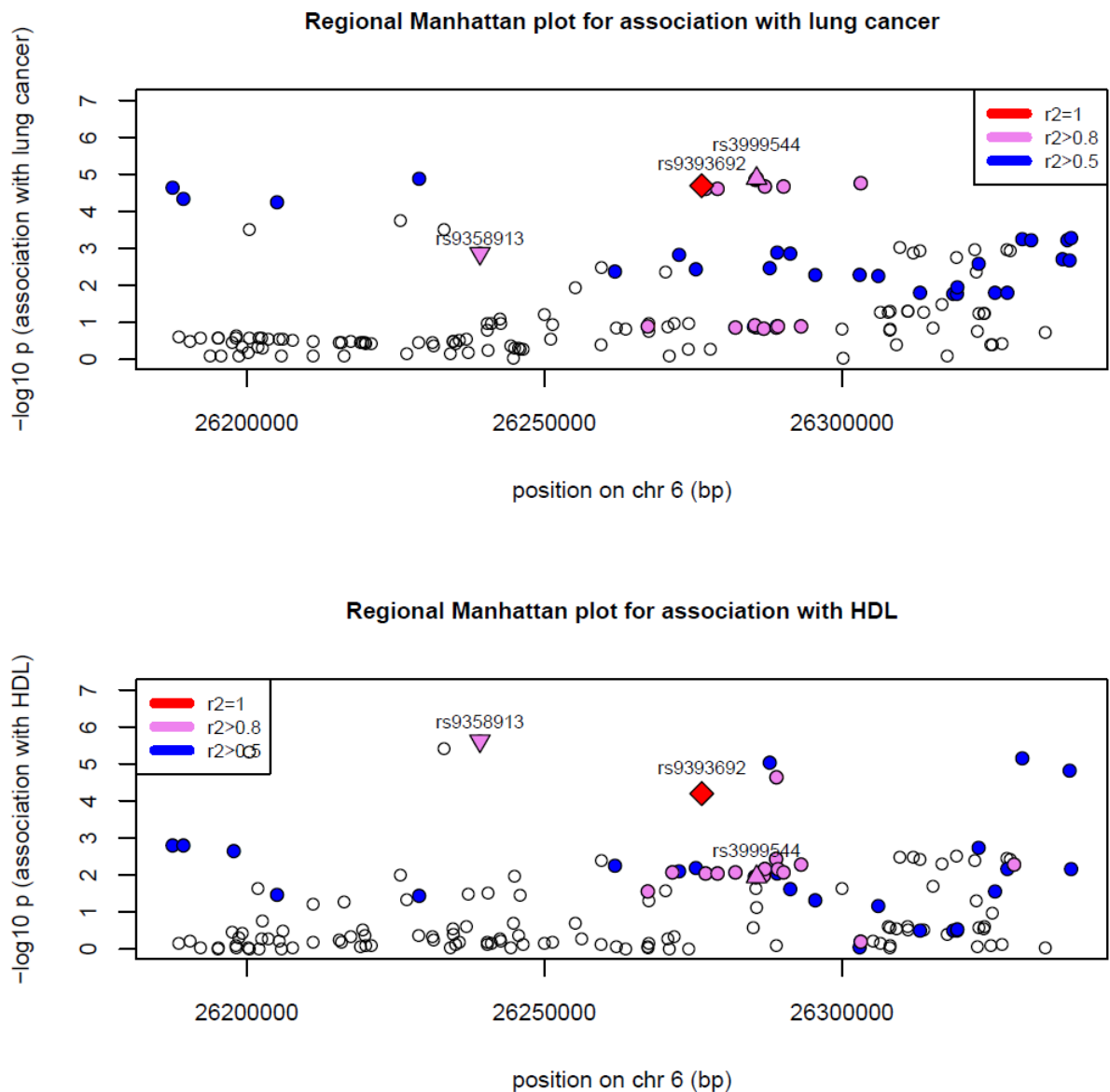
Supplementary Figure 3. Local Manhattan plot for cross-phenotype association of locus 1q22 between lung cancer and triglycerides (TG); top shared variant rs2066981

On display are the $\log_{10} p$ -values for the association of SNPs with lung cancer (upper panel) and TG (lower panel) based on the SNPs genomic position. All variants with a shared association between lung cancer and TG (conjunction FDR < 0.05) are symbolized with diamonds. The leading variant for lung cancer (rs914615) is symbolized with a triangle (top up) and the leading variant for TG was the conjunction top variant rs2066981. Linkage disequilibrium (LD) based on HapMap v.3 in Caucasians is color-coded: variants with $r^2=1$ with the lead variant rs2066981 are colored in red, $0.8 \leq r^2 < 1$ are colored in violet, and $0.5 \leq r^2 < 0.8$ are colored in blue.



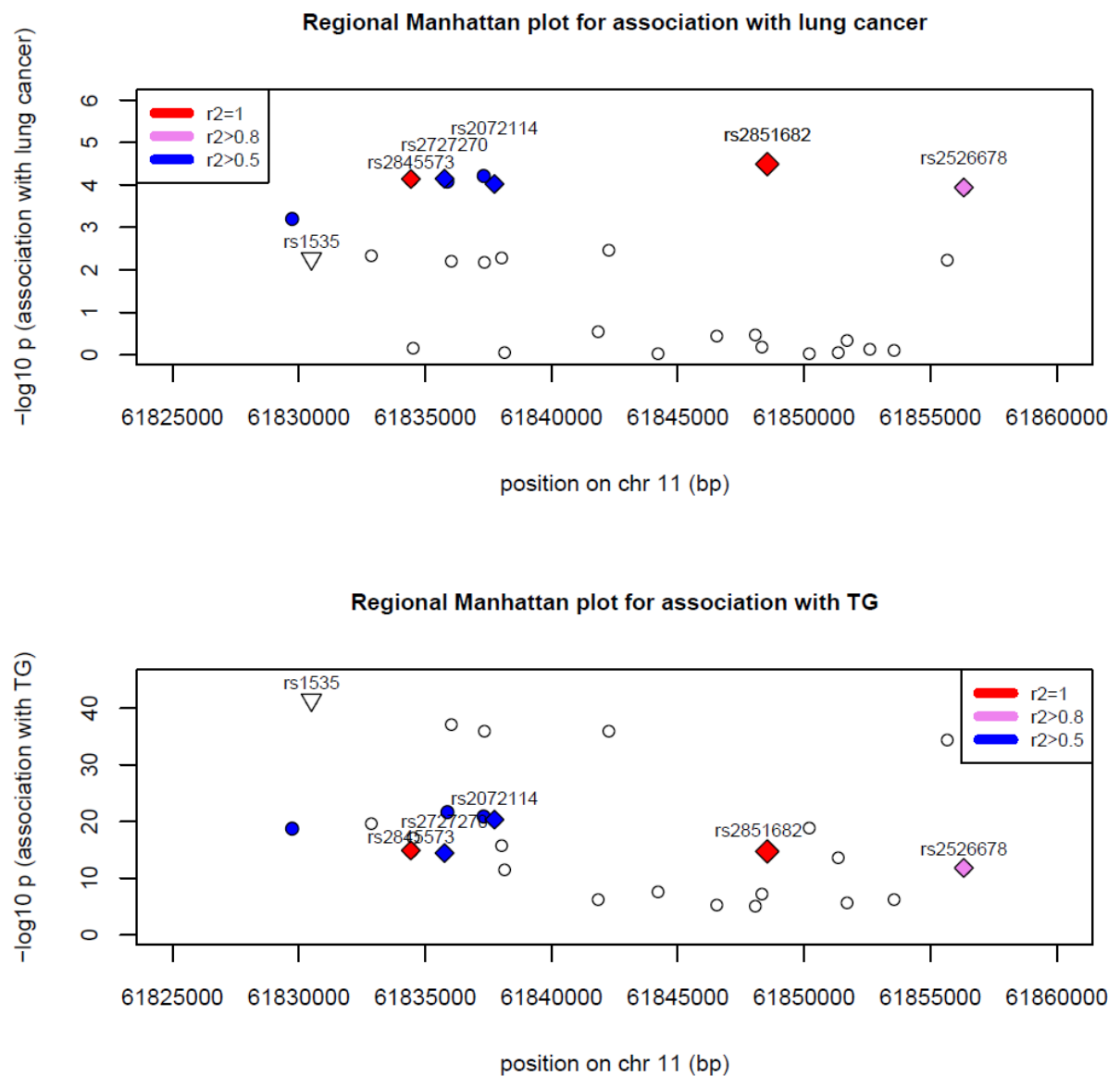
Supplementary Figure 4. Local Manhattan plot for cross-phenotype association of locus 6p22.2 between lung cancer and high-density lipoprotein (HDL); top shared variant rs9393692

On display are the $\log_{10} p$ -values for the association of SNPs with lung cancer (upper panel) and HDL (lower panel) based on the SNPs genomic position. All variants with a shared association between lung cancer and HDL (conjunction FDR < 0.05) are symbolized with diamonds. The leading variant for lung cancer (rs3999544) is symbolized with a triangle (top up) and the leading variant (rs9358913) for HDL with a triangle top down. Linkage disequilibrium (LD) based on HapMap v.3 in Caucasians is color-coded: variants with $r^2=1$ with the lead variant rs9393692 are colored in red, $0.8 \leq r^2 < 1$ are colored in violet, and $0.5 \leq r^2 < 0.8$ are colored in blue.



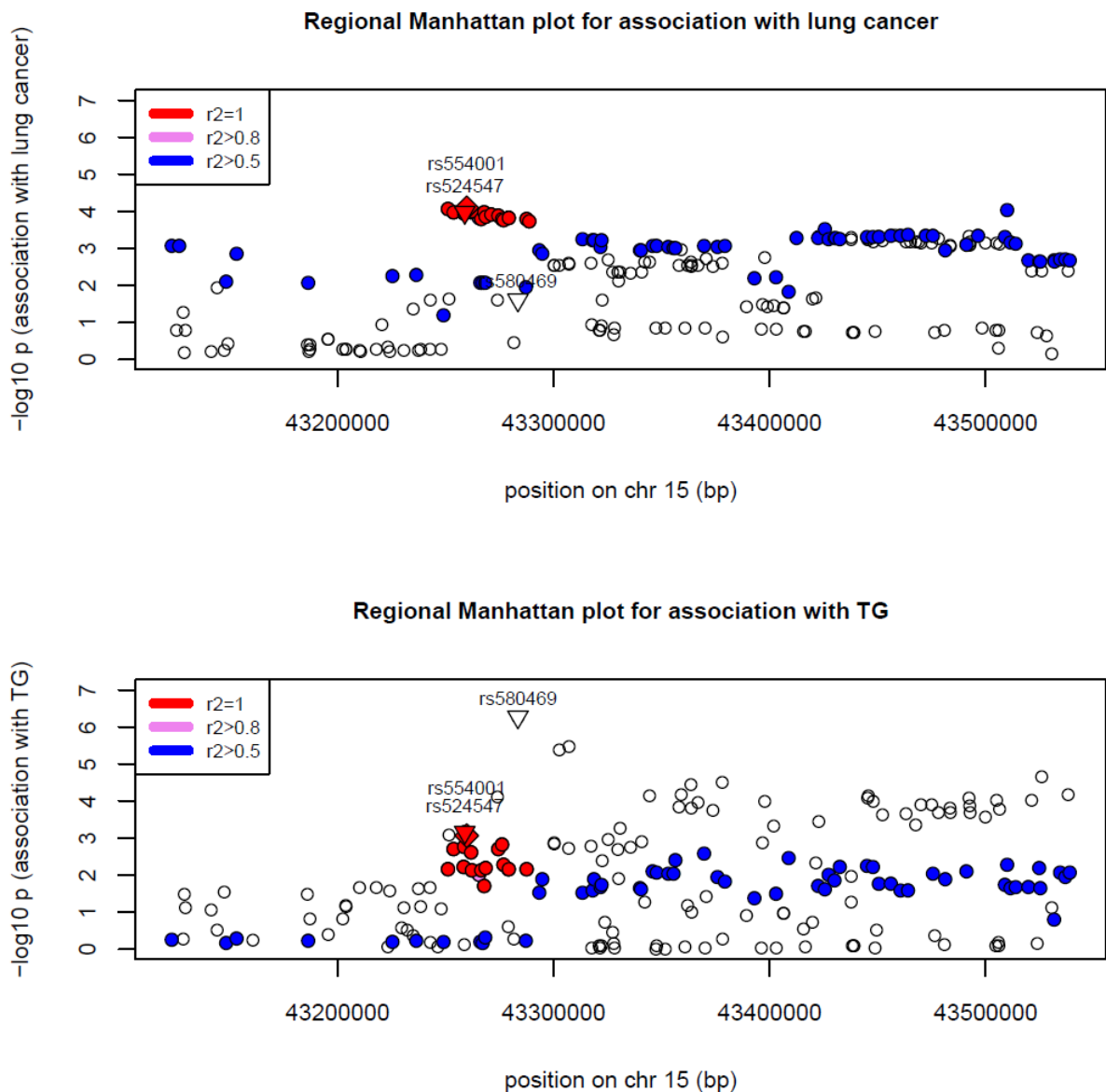
Supplementary Figure 5. Local Manhattan plot for cross-phenotype association of locus 11q12.2 between lung cancer and triglycerides (TG); top shared variant rs2851682

On display are the $\log_{10} p$ -values for the association of SNPs with lung cancer (upper panel) and TG (lower panel) based on the SNPs genomic position. All variants with a shared association between lung cancer and TG (conjunction FDR < 0.05) are symbolized with diamonds. The leading variant for lung cancer is the conjunction variant rs2851682; the leading variant for TG (rs1535) is symbolized with a triangle top down. Linkage disequilibrium (LD) based on HapMap v.3 in Caucasians is color-coded: variants with $r^2=1$ with the lead variant rs2851682 are colored in red, $0.8 \leq r^2 < 1$ are colored in violet, and $0.5 \leq r^2 < 0.8$ are colored in blue.



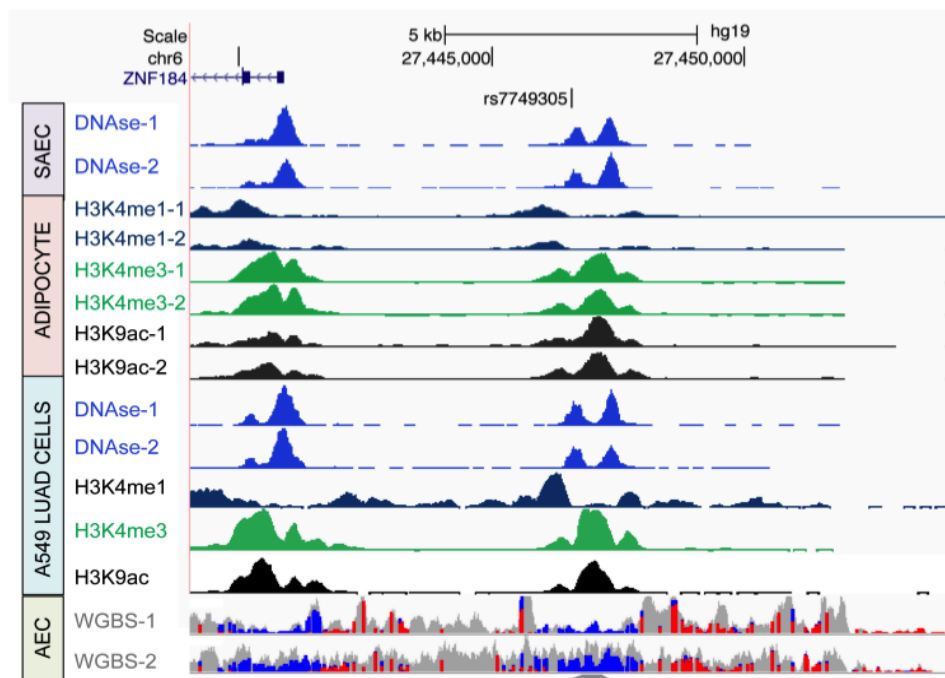
Supplementary Figure 6. Local Manhattan plot for cross-phenotype association of locus 15q15.2 between lung cancer and triglycerides (TG); top shared variant rs554001

On display are the $\log_{10} p$ -values for the association of SNPs with lung cancer (upper panel) and TG (lower panel) based on the SNPs genomic position. All variants with a shared association between lung cancer and TG (conjunction FDR < 0.05) are symbolized with diamonds. The leading variant for lung cancer is the conjunction variant rs554001; the leading variant for TG (rs580469) is symbolized with a triangle top down. Linkage disequilibrium (LD) based on HapMap v.3 in Caucasians is color-coded: variants with $r^2=1$ with the lead variant rs554001 are colored in red, $0.8 < r^2 < 1$ are colored in violet, and $0.5 < r^2 < 0.8$ are colored in blue.



Supplementary Figure 7. rs7749305 lies in a regulatory element in lung and fat tissues/cell lines.

Epigenetic annotation of rs7749305 on chromosome 6 in a frame of 10KB up and downstream using peaks from tissues relevant for lung cancer and blood lipid traits. At the top, the UCSC Genome browser image shows the *ZNF184* locus and to the right below that the position of rs7749305. Below that we show: in small airway epithelial cells (SAEC), DNase hypersensitive sites; in adipocytes, duplicate lanes of ChIP-seq marks for H3K4me1, H3K4me3 and H3K9ac; in A549 cells, duplicate DNase hypersensitive sites, and ChIP-seq marks for H3K4me1, H3K4me3 and H3K9ac; and in our own human alveolar epithelial cells (AEC), duplicate whole genome bisulfite sequencing (WGBS; red=methylated CpGs, blue=unmethylated CpGs). At the bottom, the position of the SNP and its effect on ATF3 and HIF1A binding sites is indicated.



Ref: GTAAATTCTGTTTAAAGATGTGAAGGGAGCCTT

Alt: GTAAATTCTGTTTAAAGACGTGAAGGGAGCCTT

ATF3:

HIF1A:

Supplementary Tables

Supplementary Table 1. TRICL lung cancer studies in the discovery dataset*

| Study | All | | Squamous cell carcinoma | | Adenocarcinoma | |
|--------|-------|----------|-------------------------|----------|----------------|----------|
| | cases | controls | cases | controls | cases | controls |
| UK | 1952 | 5200 | 611 | 5200 | 465 | 5200 |
| MDACC | 1150 | 1134 | 306 | 1134 | 619 | 1134 |
| IARC | 2533 | 3791 | 911 | 2968 | 517 | 2824 |
| NCI | 5713 | 5736 | 1447 | 5736 | 1841 | 5736 |
| SLRI | 331 | 499 | 50 | 499 | 90 | 499 |
| GERMAN | 481 | 478 | 97 | 478 | 186 | 478 |
| Total | 12160 | 16838 | 3422 | 16015 | 3718 | 15871 |

* UK = Institute of Cancer Research, London; MADACC = University of Texas MD

Anderson Cancer Center; IARC = International Agency for Research on Cancer; NCI =

National Cancer Institute; SLRI = Samuel Lunenfeld Research Institute, Toronto; GERMAN

= German Cancer Research Center (DKFZ).

Supplementary Table 2. Discovery dataset

| Disease/Trait | N | # SNPs | Reference |
|-------------------------------------------------------------------------------------|----------------|-----------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Lung Cancer (LgCa) | 28,998 | 2,433,836 | Wang Y, McKay JD, Rafnar T, et al. Rare variants of large effect in BRCA2 and CHEK2 affect risk of lung cancer; <i>Nat Genet</i> 2014; 46, 736–741_DETAILS IN SUPPL.Table 2 |
| Body Mass Index (BMI) | 123,865 | 2,400,377 | Speliotes EK, Willer CJ, Berndt SI, et al. Association analyses of 249,796 individuals reveal 18 new loci associated with body mass index. <i>Nat Genet</i> 2010;42:937-48. |
| Type 2 Diabetes (TD2) | 149,821 | 94,012 | Morris AP, Voight BF, Teslovich TM, et al. Large-scale association analysis provides insights into the genetic architecture and pathophysiology of type 2 diabetes. <i>Nature genetics</i> 2012; 44: 981-90; http://diagram-consortium.org |
| Low density lipoprotein (LDL) | 188,577 | 2,491,989 | Willer CJ, Schmidt EM, et al. Discovery and refinement of loci associated with lipid levels. <i>Nat Genet</i> 2013;45, 1274–1283; http://csg.sph.umich.edu//abecasis/public/lipids2013/ |
| Triglycerides (TG) | 188,577 | 2,487,152 | |
| High density lipoprotein (HDL) | 188,577 | 2,492,237 | |
| Total lipid traits contributing to the meta-analysis | 491,261 | | |
| Total lung cancer and CVD traits' subjects removing the overlapping subjects | 483,841 | | |
| Cigarettes per day (CPD) | 74,053 | 2,397,337 | The Tobacco and Genetics Consortium Genome-wide meta-analyses identify multiple loci associated with smoking behavior. <i>Nat Genet</i> 2010; 42, 441–447 |
| Total subjects including those in CPD analysis | 557,894 | | |
| eQTL lung | 1,111 | | Hao, K, Bosse, Y, Nickle, DC, et al. Lung eQTLs to help reveal the molecular underpinnings of asthma. (2012) <i>PLoS Genet</i> , 8(11), e1003029. doi: 10.1371/journal.pgen.1003029 |
| meQTL lung | 244 | | Shi, J, Marconett, CN, Duan, J, et al. Characterizing the genetic basis of methylome diversity in histologically normal human lung tissue. <i>Nat Commun</i> (2014) 5, 3365 |

Supplementary Table 3. Testing for enrichment by permutation*

| Enrichment <i>p</i>-value† | HDL | LDL | TG | BMI | T2D | CPD | SMOKER |
|---------------------------------------|------------|------------|-----------|------------|------------|------------|---------------|
| <i>p</i> <0.001 | 0.15 | 0.001 | <0.0001 | 0.36 | 0.28 | 0.46 | 0.10 |
| <i>p</i> <0.0001 | 0.06 | <0.0001 | <0.0001 | 0.37 | 0.35 | 0.55 | 0.55 |

* Smoking phenotypes analysis excluded the chrom. 15q25 locus (chr15:78686690-79231478). LDL = Low density lipoprotein cholesterol; HDL = High density lipoprotein cholesterol; TG = Triglycerides; BMI = Body mass index; T2D = Type 2 Diabetes; CPD = Cigarettes per day; SMOKER = ever vs. never smoker.

† Permutation test (one-sided) for enrichment, for further details see supplementary material.

Supplementary Table 4. Summary data for cross-phenotype associated loci for lung cancer and lipid traits*

| SNP | Gene | Band | Ref. A. | Alt. A. | MAF | LuCa&LDL | LuCa&TG | LuCa&HDL | p_LuCa [†] | p_LDL [†] | p_TG [†] | p_HDL [†] | z_LuCa [‡] | z_LDL [‡] | z_TG [‡] | z_HDL [‡] |
|-------------------------|-------------------|---------|---------|---------|------|----------|----------|----------|---------------------|--------------------|-------------------|--------------------|---------------------|--------------------|-------------------|--------------------|
| rs2066981 | THBS3 HIST1H2B | 1q22 | G | A | 0.57 | 2.70E-01 | 3.17E-02 | 6.51E-01 | 3.86E-06 | 7.40E-03 | 4.85E-04 | 0.04 | 4.62E+00 | -2.68E+00 | 3.49E+00 | 2.01E+00 |
| rs9393692 | I | 6p22.2 | A | G | 0.57 | 6.07E-01 | 1.00E+00 | 3.64E-02 | 2.03E-05 | 0.05 | 0.81 | 6.31E-05 | -4.26E+00 | -1.95E+00 | 2.46E-01 | -4.00E+00 |
| rs6904596 [§] | ZNF184 | 6p22.1 | A | G | 0.08 | 9.33E-01 | 1.24E-02 | 2.59E-02 | 8.70E-06 | 0.6.3 | 7.79E-06 | 3.19E-04 | 4.45E+00 | -4.74E-01 | 4.47E+00 | 3.60E+00 |
| rs2851682 | MYRF | 11q12.2 | A | G | 0.11 | 2.91E-02 | 2.49E-02 | 4.13E-02 | 3.23E-05 | 5.99E-18 | 1.89E-15 | 1.13E-09 | -4.16E+00 | 8.63E+00 | 7.95E+00 | 6.09E+00 |
| rs554001 [¶] | TGM5 | 15q15.2 | G | A | 0.33 | 1.00E+00 | 4.49E-02 | 3.65E-01 | 7.98E-05 | 0.8.6 | 8.63E-04 | 0.01 | -3.95E+00 | 1.74E-01 | 3.33E+00 | 2.49E+00 |

*Independent genetic loci ($r^2 < 0.2$) with a conjunction FDR < 0.05 in both lung cancer (LuCa) and in low density lipoprotein (LDL); in both lung cancer and in triglycerides (TG); or in both lung cancer and in high density lipoprotein (HDL) after removing the major histocompatibility complex (MHC) on chromosome 6 (genomic position (hg 19): chr6:29528318-33373649). The table shows the SNPs with the strongest association in each LD block based on the minimum conjunction FDR (LuCa&LDL, LuCa&TG, LuCa&HDL). LuCa = Lung cancer; LDL = Low density lipoprotein cholesterol; HDL = High density lipoprotein cholesterol; TG = Triglycerides.

[†]p = SNP association test P-value (two-sided) from a two-sided Wald test

[‡]z = SNP effect size (z-score) with respect to reference allele

[§]Only this locus has been replicated in both lung cancer and lipid traits (triglycerides)

^{||}This locus was previously identified in association with lipid traits (Global Lipids Genetics Consortium et al. Nat Genet 45, 1274-83, 2013)

[¶]This locus was previously identified in association with lung cancer risk (Rafnar, T. et al. Cancer Res 71, 1356-61, 2011)

Supplementary Table 5. Sensitivity analysis: p -values for conjunction FDR for lung cancer and triglyceride levels of the 6p22.1 locus using different r^2 thresholds for pruning SNPs in LD

| r^2 | 0.8 | 0.7 | 0.6 | 0.5 | 0.4 | 0.3 | 0.2 |
|-----------------|----------|----------|----------|----------|----------|----------|----------|
| conjunction FDR | 1.24E-02 | 1.01E-02 | 8.86E-03 | 6.81E-03 | 5.77E-03 | 4.16E-03 | 2.78E-03 |

Supplementary Table 6. Replication dataset*

| Study | Phenotype | Cases | Controls |
|------------------------|---------------------------------|-----------------|-----------------|
| Lung Cancer DeCODE | Lung cancer | 3865 | 196658 |
| | Lung adenocarcinoma | 1434 | 198663 |
| | Lung squamous cell carcinoma | 784 | 171059 |
| Harvard | Lung cancer | 984 | 970 |
| | Lung adenocarcinoma | 597 | 970 |
| | Lung squamous cell carcinoma | 216 | 970 |
| Holland | Lung cancer | 687 | 5158 |
| | Lung adenocarcinoma | 250 | 5158 |
| | Lung squamous cell carcinoma | 251 | 5158 |
| Spain | Lung cancer | 561 | 1871 |
| | Lung adenocarcinoma | 97 | 1871 |
| | Lung squamous cell carcinoma | 167 | 1871 |
| Total | | 6097 | 204657 |
| Study | Phenotype | Subjects | |
| Lipid traits deCODE | LDL* | 40724 | |
| | Triglycerides | 66027 | |
| Holland | LDL* | 5091 | |

*LDL= Low density lipoprotein cholesterol

Supplementary Table 7. Summary data for genetic loci and lead SNP rs6904596 with reference allele A and alternative allele G at 6p22.1 in lung adenocarcinoma and squamous cell carcinoma.

| Study | OR | 95% CI | p-value* |
|-----------------------------|-----------|---------------|-----------------|
| Adenocarcinoma | | | |
| TRICL | 1.05 | (0.95,1.15) | 0.354 |
| DECODE | 1.09 | (0.93,1.29) | 0.293 |
| Harvard | 1.14 | (0.87,1.50) | 0.343 |
| Dutch | 1.22 | (0.93,1.61) | 0.15 |
| Spain | 0.96 | (0.58,1.60) | 0.89 |
| Replication | 1.12 | (0.99,1.26) | 0.07 |
| Combined | 1.08 | (1.00,1.17) | 0.061 |
| <i>p</i> for heterogeneity† | | | 0.85 |
| Squamous cell carcinoma | | | |
| TRICL | 1.21 | (1.11,1.33) | 4.50E-05 |
| DECODE | 1.11 | (0.88,1.38) | 0.38 |
| Harvard | 1.15 | (0.78,1.71) | 0.48 |
| Holland | 1.05 | (0.79,1.40) | 0.73 |
| Spain | 1.18 | (0.81,1.71) | 0.4 |
| Replication | 1.11 | (0.95,1.29) | 0.18 |
| Combined | 1.18 | (1.09,1.28) | 2.80E-05 |
| <i>p</i> for heterogeneity† | | | 0.90 |

* P-values were derived from a two-sided Wald test. The reference group for the odds

ratio (OR) in the lung cancer study were healthy controls without lung cancer.

CI=confidence interval.

† Heterogeneity of effect size across studies was evaluated using the Cochran's Q statistic. The test is defined as one-sided.