

1 Genome-wide association study of lung adenocarcinoma in East Asia and comparison with a European
2 population
3

4 Jianxin Shi †,1, Kouya Shiraishi †,2, Jiyeon Choi †,1, Keitaro Matsuo †,3, Tzu-Yu Chen †,4, Juncheng Dai
5 †,5,6, Rayjean J Hung †,7, Kexin Chen †,8, Xiao-Ou Shu †,9, Young Tae Kim¹⁰, Maria Teresa Landi¹,
6 Dongxin Lin¹¹, Wei Zheng⁹, Zhihua Yin¹², Baosen Zhou¹³, Bao Song¹⁴, Jiucun Wang^{15,16}, Wei Jie
7 Seow^{1,17,18}, Lei Song¹, I-Shou Chang¹⁹, Wei Hu¹, Li-Hsin Chien⁴, Qiuyin Cai⁹, Yun-Chul Hong²⁰, Hee
8 Nam Kim²¹, Yi-Long Wu²², Maria Pik Wong²³, Brian Douglas Richardson^{1,24}, Karen M Funderburk¹,
9 Shilan Li^{1,25}, Tongwu Zhang¹, Charles Breeze¹, Zhaoming Wang²⁶, Batel Blechter¹, Bryan A Bassig¹,
10 Jin Hee Kim²⁷, Demetrius Albanes¹, Jason YY Wong¹, Min-Ho Shin²¹, Lap Ping Chung²³, Yang
11 Yang²⁸, She-Juan An²², Hong Zheng⁸, Yasushi Yatabe²⁹, Xu-Chao Zhang²², Young-Chul Kim^{30,31}, Neil
12 E Caporaso¹, Jiang Chang³², James Chung Man Ho³³, Michiaki Kubo³⁴, Yataro Daigo^{35,36}, Minsun
13 Song³⁷, Yukihide Momozawa³⁴, Yoichiro Kamatani³⁸, Masashi Kobayashi³⁹, Kenichi Okubo³⁹,
14 Takayuki Honda⁴⁰, H Dean Hosgood⁴¹, Hideo Kunitoh⁴², Harsh Patel¹, Shun-ichi Watanabe⁴³, Yohei
15 Miyagi⁴⁴, Haruhiko Nakayama⁴⁵, Shingo Matsumoto⁴⁶, Hidehito Horinouchi⁴³, Masahiro Tsuboi⁴⁷,
16 Ryuji Hamamoto⁴⁸, Koichi Goto⁴⁶, Yuichiro Ohe⁴³, Atsushi Takahashi³⁸, Akiteru Goto⁴⁹, Yoshihiro
17 Minamiya⁵⁰, Megumi Hara⁵¹, Yuichiro Nishida⁵¹, Kenji Takeuchi⁵², Kenji Wakai⁵², Koichi Matsuda⁵³,
18 Yoshinori Murakami⁵⁴, Kimihiro Shimizu⁵⁵, Hiroyuki Suzuki⁵⁶, Motonobu Saito⁵⁷, Yoichi Ohtaki⁵⁸,
19 Kazumi Tanaka⁵⁸, Tangchun Wu⁵⁹, Fusheng Wei⁶⁰, Hongji Dai⁸, Mitchell J Machiela¹, Jian Su²², Yeul
20 Hong Kim⁶¹, In-Jae Oh^{30,31}, Victor Ho Fun Lee⁶², Gee-Chen Chang^{63,64,65,66}, Ying-Huang Tsai^{67,68},
21 Kuan-Yu Chen⁶⁹, Ming-Shyan Huang⁷⁰, Wu-Chou Su⁷¹, Yuh-Min Chen⁷², Adeline Seow¹⁷, Jae Yong
22 Park⁷³, Sun-Seog Kweon^{21,74}, Kun-Chieh Chen⁶⁴, Yu-Tang Gao⁷⁵, Biyun Qian⁸, Chen Wu¹¹, Daru
23 Lu^{15,16}, Jianjun Liu^{76,77}, Ann G Schwartz⁷⁸, Richard Houlston⁷⁹, Margaret R Spitz⁸⁰, Ivan P Gorlov⁸⁰,
24 Xifeng Wu⁸¹, Ping Yang⁸², Stephen Lam⁸³, Adonina Tardon⁸⁴, Chu Chen⁸⁵, Stig E Bojesen^{86,87},
25 Mattias Johansson⁸⁸, Angela Risch^{89,90,91}, Heike Bickeböller⁹², Bu-Tian Ji¹, H-Erich Wichmann^{93,94,95},
26 David C Christiani⁹⁶, Gadi Rennert⁹⁷, Susanne Arnold⁹⁸, Paul Brennan⁸⁸, James McKay⁸⁸, John K
27 Field⁹⁹, Sanjay S Shete¹⁰⁰, Loic Le Marchand¹⁰¹, Geoffrey Liu¹⁰², Angeline Andrew¹⁰³, Lambertus A
28 Kiemeny¹⁰⁴, Shan Zienolddiny-Narui¹⁰⁵, Kjell Grankvist¹⁰⁶, Mikael Johansson¹⁰⁷, Angela Cox¹⁰⁸,
29 Fiona Taylor¹⁰⁸, Jian-Min Yuan¹⁰⁹, Philip Lazarus¹¹⁰, Matthew B Schabath¹¹¹, Melinda C Aldrich¹¹²,
30 Hyo-Sung Jeon¹¹³, Shih Sheng Jiang¹⁹, Jae Sook Sung⁶¹, Chung-Hsing Chen¹⁹, Chin-Fu Hsiao⁴, Yoo
31 Jin Jung¹¹⁴, Huan Guo¹¹⁵, Zhibin Hu⁵, Laurie Burdett^{1,116}, Meredith Yeager^{1,116}, Amy Hutchinson^{1,116},
32 Belynda Hicks^{1,116}, Jia Liu^{1,116}, Bin Zhu^{1,116}, Sonja I Berndt¹, Wei Wu¹², Junwen Wang^{117,118}, Yuqing
33 Li¹¹⁹, Jin Eun Choi¹¹³, Kyong Hwa Park⁶¹, Sook Whan Sung¹²⁰, Li Liu¹²¹, Chang Hyun Kang¹¹⁴, Wen-
34 Chang Wang¹²², Jun Xu¹²³, Peng Guan^{12,124}, Wen Tan¹¹, Chong-Jen Yu¹²⁵, Gong Yang⁹, Alan Dart
35 Loon Sihoe¹²⁶, Ying Chen¹⁷, Yi Young Choi¹¹³, Jun Suk Kim¹²⁷, Ho-Il Yoon¹²⁸, In Kyu Park¹¹⁴, Ping
36 Xu¹²⁹, Qincheng He¹², Chih-Liang Wang¹³⁰, Hsiao-Han Hung¹⁹, Roel C.H. Vermeulen¹³¹, Iona
37 Cheng¹³², Junjie Wu^{15,16}, Wei-Yen Lim¹⁷, Fang-Yu Tsai¹⁹, John K.C. Chan¹³³, Jihua Li¹³⁴, Hongyan
38 Chen^{15,16}, Hsien-Chih Lin⁴, Li Jin^{15,16}, Jie Liu¹⁴, Norie Sawada¹³⁵, Taiki Yamaji¹³⁶, Kathleen
39 Wyatt^{1,116}, Shengchao A. Li^{1,116}, Hongxia Ma^{5,6}, Meng Zhu^{5,6}, Zhehai Wang¹⁴, Sensen Cheng¹⁴,
40 Xuelian Li^{12,124}, Yangwu Ren^{12,124}, Ann Chao¹³⁷, Motoki Iwasaki^{135,136}, Junjie Zhu²⁸, Gening Jiang²⁸,
41 Ke Fei²⁸, Guoping Wu⁶⁰, Chih-Yi Chen^{138,139}, Chien-Jen Chen¹⁴⁰, Pan-Chyr Yang¹⁴¹, Jinming Yu¹⁴,
42 Victoria L. Stevens¹⁴², Joseph F. Fraumeni Jr¹, Nilanjan Chatterjee †,1,143,144, Olga Y Gorlova †,80,145,
43 Chao Agnes Hsiung †,4, Christopher I Amos †,80,145, Hongbing Shen †,5,6, Stephen J Chanock †,1,
44 Nathaniel Rothman †,1, Takashi Kohno †,2, Qing Lan †,1
45

46 ¹Division of Cancer Epidemiology and Genetics, National Cancer Institute, Rockville, MD, USA,
47 ²Division of Genome Biology, National Cancer Research Institute, Tokyo, Japan, ³Division of Cancer
48 Epidemiology and Prevention, Aichi Cancer Center Research Institute, Nagoya, Japan, ⁴Institute of
49 Population Health Sciences, National Health Research Institutes, Zhunan, Taiwan, ⁵Department of
50 Epidemiology, School of Public Health, Nanjing Medical University, Nanjing, China, ⁶Jiangsu Key
51 Lab of Cancer Biomarkers, Prevention and Treatment, Collaborative Innovation Center for Cancer
52 Medicine, Nanjing Medical University, Nanjing, China, ⁷Prosserman Centre for Population Health
53 Research, Lunenfeld-Tanenbaum Research Institute, Sinai Health, Toronto, Canada, ⁸Department of
54 Epidemiology and Biostatistics, National Clinical Research Center for Cancer, Key Laboratory of
55 Molecular Cancer Epidemiology of Tianjin, Tianjin Medical University Cancer Institute and Hospital,
56 Tianjin Medical University, Tianjin, China, ⁹Division of Epidemiology, Department of Medicine,
57 Vanderbilt University Medical Center and Vanderbilt-Ingram Cancer Center, Nashville, TN, USA,
58 ¹⁰Cancer Research Institute, Seoul National University College of Medicine, Seoul, Republic of Korea,
59 ¹¹Department of Etiology & Carcinogenesis and State Key Laboratory of Molecular Oncology, Cancer
60 Institute and Hospital, Chinese Academy of Medical Sciences and Peking Union Medical College,
61 Beijing, China, ¹²Department of Epidemiology, School of Public Health, China Medical University,
62 Shenyang, China, ¹³Department of Clinical Epidemiology and Center of Evidence Based Medicine,
63 The First Hospital of China Medical University, Shenyang, China, ¹⁴Department of Oncology,
64 Shandong Cancer Hospital and Institute, Shandong Academy of Medical Sciences, Jinan, China,
65 ¹⁵Ministry of Education Key Laboratory of Contemporary Anthropology, School of Life Sciences,
66 Fudan University, Shanghai, China, ¹⁶State Key Laboratory of Genetic Engineering, School of Life
67 Sciences, Fudan University, Shanghai, China, ¹⁷Saw Swee Hock School of Public Health, National
68 University of Singapore, Singapore, Singapore, ¹⁸Department of Medicine, Yong Loo Lin School of
69 Medicine, National University of Singapore and National University Health System, Singapore,
70 Singapore, ¹⁹National Institute of Cancer Research, National Health Research Institutes, Zhunan,
71 Taiwan, ²⁰Department of Preventive Medicine, Seoul National University College of Medicine, Seoul,
72 Republic of Korea, ²¹Department of Preventive Medicine, Chonnam National University Medical
73 School, Gwangju, Republic of Korea, ²²Guangdong Lung Cancer Institute, Medical Research Center
74 and Cancer Center of Guangdong Provincial People's Hospital, Guangdong Academy of Medical
75 Sciences, Guangzhou, China, ²³Department of Pathology, Queen Mary Hospital, Hong Kong, Hong
76 Kong, ²⁴Department of Biostatistics, Gillings School of Global Public Health, University of North
77 Carolina, Chapel Hill, NC, USA, ²⁵Department of Biostatistics, Bioinformatics & Biomathematics,
78 Georgetown University Medical Center, Washington, DC, USA, ²⁶Department of Computational
79 Biology, St. Jude Children's Research Hospital, Memphis, TN, USA, ²⁷Department of Environmental
80 Health, Graduate School of Public Health, Seoul National University, Seoul, Republic of Korea,
81 ²⁸Shanghai Pulmonary Hospital, Shanghai, China, ²⁹Department of Pathology and Clinical
82 Laboratories, National Cancer Center Hospital, Tokyo, Japan, ³⁰Lung and Esophageal Cancer Clinic,
83 Chonnam National University Hwasun Hospital, Hwasun, Republic of Korea, ³¹Department of
84 Internal Medicine, Chonnam National University Medical School, Gwangju, Republic of Korea,
85 ³²Department of Etiology & Carcinogenesis, Cancer Institute and Hospital, Chinese Academy of
86 Medical Sciences and Peking Union Medical College, Beijing, China, ³³Department of Medicine, The
87 University of Hong Kong, Queen Mary Hospital, Hong Kong, Hong Kong, ³⁴Laboratory for
88 Genotyping Development, RIKEN Center for Integrative Medical Sciences, Yokohama, Japan,
89 ³⁵Center for Antibody and Vaccine Therapy, Research Hospital, Institute of Medical Science, The
90 University of Tokyo, Tokyo, Japan, ³⁶Department of Medical Oncology and Cancer Center, and Center
91 for Advanced Medicine against Cancer, Shiga University of Medical Science, Shiga, Japan,

92 ³⁷Department of Statistics & Research Institute of Natural Sciences, Sookmyung Women's University,
93 Seoul, Republic of Korea, ³⁸Laboratory for Statistical Analysis, RIKEN Center for Integrative Medical
94 Sciences, Yokohama, Japan, ³⁹Department of Thoracic Surgery, Tokyo Medical and Dental University,
95 Tokyo, Japan, ⁴⁰Department of Respiratory Medicine, Tokyo Medical and Dental University, Tokyo,
96 Japan, ⁴¹Department of Epidemiology and Population Health, Albert Einstein College of Medicine,
97 New York, USA, ⁴²Department of Medical Oncology, Japanese Red Cross Medical Center, Tokyo,
98 Japan, ⁴³Department of Thoracic Surgery, National Cancer Center Hospital, Tokyo, Japan, ⁴⁴Molecular
99 Pathology and Genetics Division, Kanagawa Cancer Center Research Institute, Yokohama, Japan,
100 ⁴⁵Department of Thoracic Surgery, Kanagawa Cancer Center, Yokohama, Japan, ⁴⁶Department of
101 Thoracic Oncology, National Cancer Center Hospital East, Kashiwa, Japan, ⁴⁷Department of Thoracic
102 Surgery, National Cancer Center Hospital East, Kashiwa, Japan, ⁴⁸Division of Medical AI Research
103 and Development, National Cancer Center Research Institute, Tokyo, Japan, ⁴⁹Department of Cellular
104 and Organ Pathology, Graduate School of Medicine, Akita University, Akita, Japan, ⁵⁰Department of
105 Thoracic Surgery, Graduate School of Medicine, Akita University, Akita, Japan, ⁵¹Department of
106 Preventive Medicine, Faculty of Medicine, Saga University, Saga, Japan, ⁵²Department of Preventive
107 Medicine, Nagoya University Graduate School of Medicine, Nagoya, Japan, ⁵³Laboratory of Clinical
108 Genome Sequencing, Department of Computational Biology and Medical Science, Graduate School of
109 Frontier Sciences, The University of Tokyo, Tokyo, Japan, ⁵⁴Division of Molecular Pathology,
110 Institute of Medical Science, The University of Tokyo, Tokyo, Japan, ⁵⁵Department of Surgery,
111 Division of General Thoracic Surgery, Shinshu University School of Medicine Asahi, Nagano, Japan,
112 ⁵⁶Department of Chest Surgery, Fukushima Medical University School of Medicine, Fukushima,
113 Japan, ⁵⁷Department of Gastrointestinal Tract Surgery, Fukushima Medical University School of
114 Medicine, Fukushima, Japan, ⁵⁸Department of Integrative center of General Surgery, Gunma
115 University Hospital, Gunma, Japan, ⁵⁹Institute of Occupational Medicine and Ministry of Education
116 Key Lab for Environment and Health, School of Public Health, Huazhong University of Science and
117 Technology, Wuhan, China, ⁶⁰China National Environmental Monitoring Center, Beijing, China,
118 ⁶¹Department of Internal Medicine, Division of Oncology/Hematology, College of Medicine, Korea
119 University Anam Hospital, Seoul, Republic of Korea, ⁶²Department of Clinical Oncology, The
120 University of Hong Kong, Queen Mary Hospital, Hong Kong, Hong Kong, ⁶³School of Medicine and
121 Institute of Medicine, Chung Shan Medical University, Taichung, Taiwan, ⁶⁴Department of Internal
122 Medicine, Division of Pulmonary Medicine, Chung Shan Medical University Hospital, Taichung,
123 Taiwan, ⁶⁵Institute of Biomedical Sciences, National Chung Hsing University, Taichung, Taiwan,
124 ⁶⁶Department of Internal Medicine, Division of Chest Medicine, Taichung Veterans General Hospital,
125 Taichung, Taiwan, ⁶⁷Department of Respiratory Therapy, Chang Gung University, Taoyuan, Taiwan,
126 ⁶⁸Department of Pulmonary and Critical Care, Xiamen Chang Gung Hospital, Xiamen, China,
127 ⁶⁹Department of Internal Medicine, National Taiwan University Hospital and College of Medicine,
128 Taipei, Taiwan, ⁷⁰Department of Internal Medicine, E-Da Cancer Hospital, I-Shou University and
129 Kaohsiung Medical University, Kaohsiung, Taiwan, ⁷¹Department of Oncology, National Cheng Kung
130 University Hospital, College of Medicine, National Cheng Kung University, Tainan, Taiwan,
131 ⁷²Department of Chest Medicine, Taipei Veterans General Hospital, and school of Medicine, National
132 Yang Ming Chiao Tung University, Taipei, Taiwan, ⁷³Lung Cancer Center, Kyungpook National
133 University Medical Center, Daegu, Republic of Korea, ⁷⁴Jeonnam Regional Cancer Center, Chonnam
134 National University, Hwasun, Republic of Korea, ⁷⁵Department of Epidemiology, Shanghai Cancer
135 Institute, Shanghai, China, ⁷⁶Genome Institute of Singapore, Agency of Science, Technology and
136 Research, Singapore, Singapore, ⁷⁷Yong Loo Lin School of Medicine, National University of
137 Singapore, Singapore, Singapore, ⁷⁸Karmanos Cancer Institute, Detroit, MI, USA, ⁷⁹Division of

138 Genetics and Epidemiology, Institute of Cancer Research, London, UK, ⁸⁰Department of Medicine,
139 Section of Epidemiology and Population Science, Institute for Clinical and Translational Research,
140 Houston, TX, USA, ⁸¹School of Medicine, Zhejiang University, Hangzhou, Zhejiang, China,
141 ⁸²Department of Health Sciences Research, Mayo Clinic, Scottsdale, AZ, USA, ⁸³British Columbia
142 Cancer Agency, Vancouver, BC, Canada, ⁸⁴IUOPA, University of Oviedo and CIBERESP, Spain,
143 ⁸⁵Public Health Sciences Division, Fred Hutchinson Cancer Center, Seattle, Washington, USA,
144 ⁸⁶Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark,
145 ⁸⁷Department of Clinical Biochemistry, Herlev and Gentofte Hospital, Copenhagen University
146 Hospital, Copenhagen, Denmark, ⁸⁸International Agency for Research on Cancer (IARC/WHO), Lyon,
147 France, ⁸⁹German Cancer Research Center (DKFZ), Heidelberg, Germany, ⁹⁰Translational Lung
148 Research Center Heidelberg (TLRC-H), Member of the German Center for Lung Research (DZL),
149 Heidelberg, Germany, ⁹¹University of Salzburg and Cancer Cluster Salzburg, Salzburg, Austria,
150 ⁹²University Medical Center Goettingen, Goettingen, Germany, ⁹³Institute of Medical Informatics,
151 Biometry and Epidemiology, Ludwig Maximilians University, Munich, Germany, ⁹⁴Helmholtz Center
152 Munich, Institute of Epidemiology, Munich, Germany, ⁹⁵Institute of Medical Statistics and
153 Epidemiology, Technical University Munich, Munich, Germany, ⁹⁶Harvard TH Chan School of Public
154 Health, Boston, Massachusetts, USA, ⁹⁷Carmel Medical Center, Israel, ⁹⁸Markey Cancer Center,
155 Lexington, KY, USA, ⁹⁹Liverpool University, Liverpool, UK, ¹⁰⁰The University of Texas MD
156 Anderson Cancer Center, Houston, Texas, USA, ¹⁰¹Epidemiology Program, University of Hawaii
157 Cancer Center, Honolulu, HI, USA, ¹⁰²Princess Margaret Cancer Center, Toronto, ON, Canada,
158 ¹⁰³Norris Cotton Cancer Center, Lebanon, NH, USA, ¹⁰⁴Radboud University Medical Center,
159 Nijmegen, Netherlands, ¹⁰⁵National Institute of Occupational Health, Oslo, Norway, ¹⁰⁶Department of
160 Medical Biosciences, Umeå University, Umeå, Sweden, ¹⁰⁷Department of Radiation Sciences, Umeå
161 University, Umeå, Sweden, ¹⁰⁸University of Sheffield, Sheffield, UK, ¹⁰⁹UPMC Hillman Cancer
162 Center and Department of Epidemiology, School of Public Health, University of Pittsburgh,
163 Pittsburgh, PA, USA, ¹¹⁰Washington State University College of Pharmacy, Spokane, WA, USA,
164 ¹¹¹Department of Cancer Epidemiology, H. Lee Moffitt Cancer Center and Research Institute, Tampa,
165 FL, USA, ¹¹²Department of Thoracic Surgery, Division of Epidemiology, Vanderbilt University
166 Medical Center, Nashville, TN, USA, ¹¹³Cancer Research Center, Kyungpook National University
167 Medical Center, Daegu, Republic of Korea, ¹¹⁴Department of Thoracic and Cardiovascular Surgery,
168 Cancer Research Institute, Seoul National University College of Medicine, Seoul, Republic of Korea,
169 ¹¹⁵Department of Occupational and Environmental Health and Ministry of Education Key Lab for
170 Environment and Health, School of Public Health, Tongji Medical College, Huazhong University of
171 Science and Technology, Wuhan, China, ¹¹⁶Cancer Genomics Research Laboratory, Leidos
172 Biomedical Research Inc., Rockville, MD, USA, ¹¹⁷Department of Biochemistry, Li Ka Shing (LKS)
173 Faculty of Medicine, The University of Hong Kong, Hong Kong, China, ¹¹⁸Centre for Genomic
174 Sciences, Li Ka Shing (LKS) Faculty of Medicine, The University of Hong Kong, Hong Kong, China,
175 ¹¹⁹Department of Human Genetics, Genome Institute of Singapore, Singapore, Singapore,
176 ¹²⁰Department of Thoracic and Cardiovascular Surgery, Seoul National University Bundang Hospital,
177 Seongnam, Republic of Korea, ¹²¹Department of Oncology, Cancer Center, Union Hospital, Huazhong
178 University of Science and Technology, Wuhan, China, ¹²²The Ph.D. Program for Translational
179 Medicine, College of Medical Science and Technology, Taipei Medical University, Taipei, Taiwan,
180 ¹²³School of Public Health, Li Ka Shing (LKS) Faculty of Medicine, The University of Hong Kong,
181 Hong Kong, China, ¹²⁴Key Laboratory of Cancer Etiology and Intervention, University of Liaoning
182 Province, Shenyang, China, ¹²⁵Department of Internal Medicine, National Taiwan University Hospital
183 Hsin-Chu Branch, Hsinchu, Taiwan, ¹²⁶Gleneagles Hong Kong Hospital, Hong Kong, China,

184 ¹²⁷Department of Internal Medicine, Division of Medical Oncology, College of Medicine, Korea
185 University Guro Hospital, Seoul, Republic of Korea, ¹²⁸Department of Internal Medicine, Seoul
186 National University Bundang Hospital, Seongnam, Republic of Korea, ¹²⁹Department of Oncology,
187 Wuhan Iron and Steel (Group) Corporation Staff-Worker Hospital, Wuhan, China, ¹³⁰Department of
188 Pulmonary and Critical Care, Chang Gung Memorial Hospital, Taoyuan, Taiwan, ¹³¹Division of
189 Environmental Epidemiology, Institute for Risk Assessment Sciences (IRAS), Utrecht University,
190 Utrecht, The Netherlands, ¹³²Department of Epidemiology and Biostatistics, University of California,
191 San Francisco, San Francisco, CA, USA, ¹³³Department of Pathology, Queen Elizabeth Hospital, Hong
192 Kong, China, ¹³⁴Qijing Center for Diseases Control and Prevention, Qijing, China, ¹³⁵Division of
193 Cohort Research, National Cancer Center Institute for Cancer Control, National Cancer Center, Tokyo,
194 Japan, ¹³⁶Division of Epidemiology, National Cancer Center Institute for Cancer Control, National
195 Cancer Center, Tokyo, Japan, ¹³⁷Center for Global Health, National Cancer Institute, Bethesda, MD,
196 USA, ¹³⁸Institute of Medicine, Chung Shan Medical University, Taichung, Taiwan, ¹³⁹Division of
197 Thoracic Surgery, Department of Surgery, Chung Shan Medical University Hospital, Taichung,
198 Taiwan, ¹⁴⁰Genomic Research Center, Academia Sinica, Taipei, Taiwan, ¹⁴¹Department of Internal
199 Medicine, National Taiwan University Hospital, Taipei, Taiwan, ¹⁴²Laboratory Services, American
200 Cancer Society, Georgia, USA, ¹⁴³Department of Oncology, School of Medicine, Johns Hopkins
201 University, Baltimore, MD, USA, ¹⁴⁴Department of Biostatistics, Johns Hopkins Bloomberg School of
202 Public Health, Baltimore, MD, USA, ¹⁴⁵Dan L Duncan Comprehensive Cancer Center, Baylor College
203 of Medicine, Houston, TX, USA
204
205

206 † These authors contributed equally to this work

207 ‡ These authors jointly supervised this work

208

209 Correspondence: jianxin.shi@nih.gov or qingl@mail.nih.gov.

210

211 **Abstract**

212

213 Lung adenocarcinoma is the most common type of lung cancer. Known risk variants explain only a
214 small fraction of lung adenocarcinoma heritability. Here, we conducted a two-stage genome-wide
215 association study of lung adenocarcinoma of East Asian ancestry (21,658 cases and 150,676 controls;
216 54.5% never-smokers) and identified 12 novel susceptibility variants, bringing the total number to 28
217 at 25 independent loci. Transcriptome-wide association analyses together with colocalization studies
218 using a Taiwanese lung expression quantitative trait loci dataset (n=115) identified novel candidate
219 genes, including *FADS1* at 11q12 and *ELF5* at 11p13. In a multi-ancestry meta-analysis of East Asian
220 and European studies, four loci were identified at 2p11, 4q32, 16q23, and 18q12. At the same time,
221 most of our findings in East Asian populations showed no evidence of association in European
222 populations. In our studies drawn from East Asian populations, a polygenic risk score based on the 25
223 loci had a stronger association in never-smokers vs. individuals with a history of smoking
224 ($P_{\text{interaction}}=0.0058$). These findings provide new insights into the etiology of lung adenocarcinoma in
225 individuals of East Asian ancestry, which could be important in developing translational applications.

226

227

228

229 **Introduction**

230 Lung adenocarcinoma (LUAD) is the most common histologic subtype of lung cancer and accounts for
231 approximately 40% of lung cancer incidence worldwide^{1, 2, 3}. In studies drawn from East Asian (EA)
232 ancestry, LUAD has been the predominant histologic subtype among females² and has replaced
233 squamous cell carcinoma as the most common subtype in males^{4, 5}. Well established risk factors,
234 namely, tobacco smoking, certain environmental/occupational exposures and lifestyle factors, and
235 family history, contribute to the risk of LUAD^{6, 7, 8}. In addition, multiple genome-wide association
236 studies (GWAS) have identified at least 24 susceptibility loci for LUAD that achieved genome-wide
237 significance, many drawn from studies in EA^{9, 10, 11, 12, 13, 14, 15} and European (EUR)^{16, 17, 18, 19, 20, 21, 22, 23}
238 populations, as well as multi-ancestry meta-analyses^{24, 25}. Of these, 12 loci have been reported at
239 genome-wide significance in GWAS of either never-smokers^{9, 11, 12, 13} or smokers and nonsmokers
240 combined^{10, 14, 15, 24} in EA populations while another two loci were suggested in a multi-ancestry meta-
241 analysis²⁴. We estimated that the known susceptibility variants account for only 13% of the estimated
242 familial risk in EA populations. Accordingly, larger studies are needed to investigate the underlying
243 architecture of susceptibility to LUAD in never-smokers and individuals with a history of smoking and
244 in different ancestral populations. The importance of multi-ancestry analyses is further highlighted by
245 reports of susceptibility loci showing association for LUAD in EA but not in EUR populations¹³.

246 In the current study, we conducted a two-stage GWAS meta-analysis in EA populations using
247 unpublished and previously published data from four studies: the Female Lung Cancer Consortium in
248 Asia (FLCCA), Nanjing Lung Cancer Study (NJLCS)^{10, 24}, National Cancer Center Research Institute
249 (NCC) and Aichi Cancer Center (ACC), with 11,753 cases and 30,562 controls in the discovery set and
250 9,905 cases and 120,114 controls in the replication set. A multi-ancestry meta-analysis of EA and
251 EUR studies^{16, 22} (from the International Lung Cancer Consortium, ILCCO) was performed to identify
252 variants shared by both populations. We also investigated the heterogeneity of effect sizes for
253 susceptibility variants identified in EA and EUR populations^{16, 22} and obtained genome-wide estimates
254 of effect-size correlation. Finally, we evaluated the genetic architecture²⁶ of LUAD, characterized by
255 the number of susceptibility variants and their effect size distribution after normalizing allele
256 frequencies, to investigate the accuracy of genetic risk prediction in the future GWAS in EA
257 populations with increased sample sizes.

258 **Results**

260 **Two-stage GWAS meta-analysis of LUAD in East Asian populations**

261 For the discovery set, we performed a fixed-effect meta-analysis (11,753 cases and 30,562 controls)
262 drawn from EA studies (Table 1, Supplementary Table 1). Details of quality control, imputation and
263 post-imputation filtering are described in Methods. Variants with an imputation quality score ≥ 0.5 and
264 minor allele frequency (MAF) ≥ 0.01 were included for meta-analysis. The estimated genetic
265 correlation between LUAD in never-smokers and individuals with a history of smoking was $r_g = 0.81$
266 (s.e. = 0.16) using linkage disequilibrium (LD) score regression (LDSC)²⁷, which enabled the primary
267 meta-analysis to include the two groups. LDSC analysis suggested little evidence of residual
268 population stratification (LDSC intercept = 1.03). We identified 14 loci achieving genome-wide
269 significance $P < 5 \times 10^{-8}$ (Supplementary Table 2); two were novel at 2p23.3 (rs682888, OR = 0.89, P =

270 4.94×10^{-10}) and at 7q31.33 (rs4268071, OR = 1.39, P = 7.27×10^{-10}). In meta-analysis performed
271 separately for males and females, and for never-smokers and individuals with a history of smoking, no
272 further loci achieved genome-wide significance.

273 In the replication phase, we selected 37 lead variants with $P < 10^{-5}$ in the discovery data that were not
274 previously reported as genome-wide significant in either EA or EUR populations and genotyped them
275 in an independent data set of 9,905 LUAD cases and 120,114 controls from a Japanese population
276 (Table 1, Supplementary Table 1). After combining the discovery and the replication data, we
277 identified a total of 10 novel loci achieving genome-wide significance and a novel variant on the locus
278 at 15q21.2 that was previously reported in EUR populations¹⁶ (Table 2, Manhattan plot in Fig. 1, and
279 regional association plots in Supplementary Fig. 1).

280 Conditional analysis using GWAS summary statistics suggested two additional susceptibility variants
281 rs13167280 (OR = 1.29, P = 4.07×10^{-13}) and rs62332591 (OR = 0.87, P = 3.21×10^{-8}) in the locus at
282 5p15.33 (Table 3, Supplementary Fig. 2); both are in modest LD with previously reported secondary
283 variants in EA populations²⁸ ($R^2=0.27$ between rs13167280 and rs10054203²⁸; $R^2=0.19$ between
284 rs62332591 and rs10054203²⁸). Another variant, rs12664490 (OR = 0.81, P = 1.24×10^{-10}), was
285 conditionally significant in a locus previously reported in EA at 6p21.1 (Table 3, Supplementary Fig.
286 3), adding another novel variant (12 novel variants in total).

287 A previous multi-ancestry meta-analysis conducted by Dai *et al.*²⁴ that included Chinese samples and
288 EUR samples from the ILCCO study identified three SNPs for LUAD, one of which achieved genome-
289 wide significance and the other two were suggestive in their analysis restricted to the Chinese
290 subgroup²⁴ (see Supplementary Table 3). In the meta-analysis of the Chinese samples in Dai *et al.*²⁴

291 with our independent EA samples, all three variants exceeded the threshold of genome-wide
292 significance without issues of heterogeneity (Supplementary Table 3).

293 Overall, our study identified 12 novel susceptibility variants bringing the total to 28 genetic variants at
294 25 loci that have been identified to date in EA populations (Supplementary Table 4, Fig. 1). Assuming
295 a familial risk estimate of 1.84 for first-degree relatives²⁹, the 25 independent susceptibility variants for
296 LUAD (Supplementary Table 4) captured 16.2% of the familial relative risk in EA populations.
297 Moreover, we found no evidence that the SNP associations differed between the samples from the
298 Mainland of China and those from outside of the Mainland of China, or between Han Chinese and
299 Japanese, the two largest ancestry populations in our study (Supplementary Table 5).

300 We further examined whether the novel variants identified in this study were associated with smoking
301 behaviors (i.e., smoking status, cigarettes per day, initiation age and cessation) or chronic obstructive
302 pulmonary disease in the Biobank Japan Project³⁰ (BBJ). We found no evidence that these variants
303 were implicated in these traits in this cohort (Supplementary Table 6). A previous GWAS in EUR
304 populations found variants (e.g., rs55781567) at the 15q25.1 *CHRNA5* locus associated with tobacco
305 smoking and lung cancer risk only in individuals with a history of smoking (OR=1.33, P= 1.83×10^{-78} ,
306 MAF=0.39)^{16, 19, 31, 32}. However, this variant did not achieve genome-wide significance in our EA data
307 (OR=1.37, P=0.001 for individuals with a history of smoking; OR=1.05, P=0.44 for never-smokers),
308 likely because of a low MAF=0.03, and no other variant in LD with this SNP showed a substantial
309 association.

310 **Fine mapping and functional analyses of GWAS loci**

311 To prioritize candidate variants for functional follow-up from each of the LUAD GWAS loci, we
312 performed Bayesian fine mapping using FINEMAP³³ (Methods). Fine mapping of the genome-wide

313 significant loci from the discovery set nominated 95% credible set variants for 9 loci with a median of
314 63 variants per locus (Supplementary Data 1). For the 12 novel variants identified from the combined
315 discovery and replication datasets as well as conditional analysis, we then performed variant
316 annotation analysis. High-LD variants for these signals ($R^2 \geq 0.8$ with the lead SNP in the 1000
317 Genomes, phase 3, EA) included those located in predicted promoters or enhancers in lung tissues/cells
318 (RegulomeDB³⁴, Haploreg³⁵ v4.1, and FORGE2³⁶; Supplementary Data 2), which can be tested in
319 future experimental studies.

320 To further characterize the functionality of the prioritized susceptibility genes that could explain the
321 new GWAS loci, eQTL colocalization and transcriptome-wide association study (TWAS) analyses
322 were conducted. Initial stratified LD score regression³⁷ using GTEx data (Supplementary Fig. 4;
323 Supplementary Data 3) indicated that LUAD heritability drawn from EA populations are enriched in
324 lung tissue-specific genes and chromatin features compared to other tissues (aggregated rank test $P =$
325 1.36×10^{-2} and 7.7×10^{-3} , respectively; Supplementary Data 3). Accordingly, we performed eQTL
326 analyses using the Taiwanese dataset of adjacent normal lung tissues from 115 never-smoking lung
327 cancer patients (LCTCNS) (Methods; Supplementary data 4). We performed colocalization analyses of
328 eQTL genes using eCAVIAR³⁸ and HyPrColoc³⁹. A notable finding was the colocalization of *FADS1*
329 at 11q12.2 (rs174559, posterior probability = 0.91) (Fig. 2; Supplementary Data 5), particularly since
330 rs174559 was in LD with a recently identified functional variant (rs174557) regulating allelic *FADS1*
331 expression in liver cells⁴⁰. *FADS1* encodes fatty acid desaturase 1, which is a key enzyme in the
332 metabolism of polyunsaturated fatty acids and plays a key role in inflammatory diseases⁴¹. Higher
333 *FADS1* levels in the lung tissues were associated with LUAD risk, which is consistent with its role in
334 increasing the proliferation and migration of laryngeal squamous cell carcinoma through activation of

335 the Akt/mTOR pathway⁴². Among the known loci, colocalization identified *TP63* at 3q28 and
336 *ACVR1B* at 12q13.13 (Supplementary Data 5).

337 We then performed a TWAS using LCTCNS eQTL dataset. TWAS identified *FADS1* as a
338 susceptibility gene from the 11q12.2 locus (TWAS $P=3.01\times 10^{-6}$) validating the finding from the
339 colocalization analysis. We further identified *ELF5* (TWAS $P=1.89\times 10^{-8}$) as a novel gene from a locus
340 (at 11p13) not originally passing the genome-wide significance threshold based on a single variant test
341 in our EA discovery GWAS (Supplementary data 6, Methods). For these two loci, we also performed
342 TWAS conditional analysis to assess whether genetically predicted expression of these genes explain
343 most of the GWAS signal. When GWAS signal was conditioned on predicted expression of *ELF5*,
344 most of the signal disappeared, adding support for *ELF5* as the main susceptibility gene in this locus
345 (Supplementary Fig. 5A). *ELF5* encodes E74-like factor 5, a key transcription factor of alveologenesis
346 of mammary glands⁴³. Lower levels of *ELF5* were associated with LUAD risk in the TWAS. Similarly,
347 when GWAS signal was conditioned on predicted expression of *FADS1*, the strongest part of the signal
348 disappeared (Supplementary Fig. 5B). We further performed TWAS analysis using GTEx lung eQTL
349 dataset (v8, n = 515, ~85% Europeans) and identified five genes from four loci (Supplementary Data
350 6). While identification of *ELF5* was common between two datasets, GTEx identified four unique
351 genes from three known loci (*DCBLD1*, *MPZL3*, *JAML*, and *LINC00674*). Notably, *FADS1* was
352 identified only by ancestry-matched LCTCNS eQTL dataset even with a ~4 times smaller sample size.

353 An investigation of the local environment of susceptibility loci revealed further plausible candidate
354 genes that could be pursued in laboratory follow-up. For instance, rs137884934 on 3q22.3 maps to
355 *PIK3CB* encoding an isoform of p110 catalytic subunit of Class IA PI3K⁴⁴. Previous studies have
356 shown that PI3K/Akt/mTOR signaling pathway plays an important role in the development and
357 progression of non-small cell lung cancer⁴⁵. Moreover, rs764014 on 15q21.3 is located adjacent to

358 *NEDD4*, which is a negative regulator of tumor suppressor PTEN⁴⁶, which encodes a lipid phosphatase
359 which counteracts the growth promoting effect of PI3K pathway⁴⁷.

360 **Multi-ancestry meta-analysis in East Asian and European populations**

361 To identify variants shared by EA and EUR populations, we performed a fixed effect, multi-ancestry
362 GWAS meta-analysis including data from samples in EA (11,753 cases and 30,562 controls) and samples
363 from EUR populations (11,273 cases and 55,483 controls). We identified four additional loci
364 (Supplementary Table 7) with similar effect sizes in the two populations: rs1130866 (2p11.2, OR = 1.08,
365 $P = 1.56 \times 10^{-8}$), rs2320614 (4q32.2, OR = 1.08, $P = 6.51 \times 10^{-9}$), rs34638657 (16q23.3, OR = 1.09, $P =$
366 2.19×10^{-9}) and rs638868 (18q12.1, OR=1.08, $P=3.6 \times 10^{-8}$). Regional association plots are shown in
367 Supplementary Fig. 6. A multi-ancestry meta-analysis stratified by smoking status did not reveal loci
368 specific to never-smokers or individuals with a history of smoking (sample size information in
369 Supplementary Table 8).

370 Among the four loci, rs1130866 at 2p11.2 is a missense variant (Ile131Thr) of *SFTPB*, encoding
371 surfactant protein B. Pulmonary surfactant lines the alveoli of lung to reduce the surface tension and is
372 essential for lung function, and increasing circulating level of pro-SFTPB suggested increased lung
373 cancer risk based on prediagnostic samples⁴⁸. Notably, two other novel variants, rs34638657 at 16q23.3
374 (*MPHOSPH6*)^{49, 50} and rs2320614 at 4q32.2 (*NAFI*)⁵¹, are on or near genes implicated in telomere
375 biology. Together with other known or new loci (rs2736100 *TERT*, rs4268071 *POT1*, rs75031349
376 *RTEL1*^{52, 53}, rs7902587 *OBFC1*⁵⁴, rs35446936 *TERC*) (Supplementary data 7), our findings further
377 support the role of telomere biology in LUAD.

378

379 **Mendelian randomization analysis of telomere length**

380 We performed a Mendelian randomization (MR) analysis to investigate a potential causal relationship
381 between telomere length and the risk of LUAD. The MR analysis was based on 46 independent
382 variants identified in a recent multi-ancestry GWAS of telomere length in the TOPMed study⁵⁵,
383 cumulatively accounting for 3.74% of telomere length variance (Methods). Since genetic effects on
384 telomere length showed no evidence of heterogeneity across populations in the TOPMed study, we
385 used the genetic effects estimated based on all populations in the TOPMed study. Our MR analysis
386 was based on MR-PRESSO⁵⁶, a robust approach that estimates causal effects after removing variants
387 detected with evidence of pleiotropic effects. Genetically predicted longer telomere length was
388 significantly associated with increased risk of LUAD with similar ORs (per one standard deviation
389 change in genetically increased telomere length) between the two populations: OR = 2.61 (95% CI =
390 2.08, 3.28, P = 8.14×10^{-10}) in EA populations, OR = 2.67 (95% CI = 2.07, 3.43, P = 7.14×10^{-9}) in
391 EUR populations, consistent with previous MR reports^{57, 58, 59} as well as a study of white blood cell
392 DNA telomere length and lung cancer risk in multiple prospective cohorts⁶⁰. MR analyses stratified by
393 smoking status showed similar results between never-smokers and individuals with a history of
394 smoking (Supplementary Table 9). We performed sensitivity analyses using genetic effects estimated
395 based on Asian and European populations in the TOPMed study separately and found similar results
396 (Supplementary Table 9).

397 **Comparing the genetics of LUAD in EA and EUR populations**

398 We systematically compared the effect size in EA vs. EUR populations of 38 susceptibility variants for
399 LUAD. These included 12 variants identified in the current study, 26 variants previously reported in
400 EA^{10, 11, 13, 14, 15, 61} and/or EUR^{16, 19, 20} populations, and results of multi-ancestry meta-analyses
401 combining data from EA and EUR²⁴ populations (Supplementary Data 8). As expected, 11 SNP

402 associations that were independently identified in both populations and through multi-ancestry analysis
403 were very similar (Figs 3A, B, C). In contrast, out of the 19 SNP associations initially identified in EA
404 populations, two had $MAF < 0.01$, 11 showed no evidence of association within EUR populations at
405 $P < 0.05$ (Fig. 3D and Fig. 3E, Supplementary Data 8), and 11 associations were significantly different
406 between the two populations with $FDR < 0.05$. Similar population differences were observed among
407 never-smokers and individuals with a history of smoking (Supplementary Fig. 7). For variants with
408 $MAF > 0.01$ in both populations, the lack of association in EUR populations did not seem to be driven
409 by low MAF or lower statistical power, as MAFs in both populations for most variants were similar
410 and GWAS in both populations had adequate power to detect at least some evidence of association
411 (Supplementary Data 9). Further, evaluation of gene region plots that spanned 500 kb for these loci
412 within EUR populations showed no or very weak evidence of association for other variants in the
413 region as well as the lead variants from the EA populations (Supplementary Figs 8A-J), with one
414 exception (Supplementary Fig. 8K). For SNPs initially identified in EUR populations, there was
415 evidence of association for 5 variants in EA populations (Fig. 3F, Supplementary Fig. 9) although all
416 variants were attenuated in the EA compared to the EUR population and one variant had MAF less
417 than 1% in EA; moreover, two variants were significantly weaker (Supplementary Data 8,
418 Supplementary Fig. 9). Similar patterns were observed among never-smokers and individuals with
419 smoking history (Supplementary Fig. 7).

420 We used LDSC²⁷ to evaluate the heritability and genetic correlation between individuals with a history
421 of smoking and never-smokers within each population and POPCORN⁶² across populations. The
422 genetic correlation was weaker between never-smokers in EA and EUR populations compared to
423 individuals with a history of smoking (Supplementary Fig. 10) although power was limited given the

424 relatively small sample sizes within each group (Supplementary Table 8). Larger sample sizes are
425 needed to estimate these characteristics more precisely.

426 **Polygenic risk score and gene-smoking interaction analysis**

427 We investigated whether the polygenic risk score (PRS), which was based on the cumulative effect of
428 25 independent susceptibility loci for LUAD in EA (Supplementary Table 4), interacted with smoking
429 status to influence the risk of LUAD, given previous evidence of gene-environment interaction^{63, 64}.
430 Since only summary statistics were available for some datasets (instead of individual genotype data),
431 we developed a statistical method for testing the multiplicative smoking-PRS interaction using the
432 summary statistics for the susceptibility variants (Methods). Compared to the middle quintile that
433 represents the average risk in the general population, the top quintile had OR of 2.07 (95% CI = 1.99,
434 2.15) for never-smokers and 1.80 (95% CI = 1.70, 1.89) for individuals with a history of smoking
435 ($P_{\text{interaction}} = 0.0058$, Fig. 4, Supplementary Fig. 11), providing statistical evidence that the association
436 between PRS and LUAD risk was higher for never-smokers. Moreover, we tested for the presence of
437 multiplicative interactions between smoking status and each individual susceptibility variant in the
438 PRS and found five variants with stronger associations in never-smokers than in individuals with a
439 history of smoking ($P < 0.05$) (Supplementary Table 2).

440 **Genetic architecture, performance of PRS and sample size requirements in EA populations**

441 To further investigate the underlying genetic architecture of susceptibility (Methods) to LUAD⁶⁵ in
442 EA populations, we performed a GENESIS²⁶ analysis based on the GWAS summary statistics for our
443 larger never-smoking dataset. We estimated that approximately 2,275 (s.e.=1,167) susceptibility
444 variants are independently associated with LUAD, suggesting that LUAD is a highly polygenic disease
445 and most of the susceptibility variants have very small effect sizes. Based on the estimated parameters,

446 we investigated how the performance of a PRS, measured as the area under the receiver operating
447 characteristic curve (AUC), depended on the sample size of the training GWAS (Fig. 5). The AUC is
448 predicted to be 60.7% (95% CI = 56.6%, 64.8%) at the current sample size and will increase to 66.9%
449 (95% CI = 62.5%, 71.3%) when the sample size increases to 70,000 cases with one control per case
450 and 68.4% (95% CI = 64.0%, 72.8%) with 1,000,000 controls. Of note, even a small increase of AUC
451 value for a PRS can help identify many more subjects at risk⁶⁶.

452 **Discussion**

453 We conducted the largest GWAS of LUAD in an EA population to date and identified 12 novel
454 susceptibility variants achieving genome-wide significance. In addition, two variants identified from a
455 previous multi-ancestry meta-analysis achieved genome-wide significance as well in EA alone after we
456 combined the reported summary data with our independent data. In total, including the previously
457 described genetic variants, 28 variants at 25 loci have reached genome-wide significance for LUAD in
458 EA populations, representing major progress in elucidating the genetic basis of LUAD. Finally, a
459 multi-ancestry meta-analysis identified four additional loci in the combined EA and EUR populations,
460 with consistent effects in both.

461 Our eQTL colocalization and TWAS analyses using an ancestry-matched lung eQTL dataset (EA
462 population) identified novel LUAD susceptibility genes including *FADS1* and *ELF5*.

463 Importantly, *FADS1* is regulated by sterol-response element-binding proteins (SREBPs)⁶⁷, which
464 govern lipid metabolism in alveolar type II (ATII) cells⁶⁸. *ELF5* is also expressed in tissues with
465 glandular/secretory epithelial cells including salivary gland and lung^{69, 70} and 3.2% of lung alveolar
466 type II cells express *ELF5* in GTEx single-cell expression data. Identification of *FADS1* and *ELF5* in
467 our study suggests a role for alveolar lineage-specific genes and pathways in LUAD susceptibility.

468 Notably, the missense variant (Ile131Thr), rs1130866, in *SFTPB* identified through the multi-ancestry
469 analysis was a protein quantitative trait locus (pQTL) for SFTPB in blood⁷¹, where the LUAD risk-
470 associated A allele (Ile131) is correlated with increased SFTPB levels. Importantly, the genomic
471 region encompassing rs1130866 presents weak LD and high SNP density, consistent with the presence
472 of a recombination hot spot⁷², and therefore fine-mapping inspecting low-frequency variants in the
473 region is warranted. Our TWAS analyses using both ancestry-matched and ancestry-discordant lung
474 eQTL datasets identified both common and unique genes from each dataset, highlighting potential
475 benefits of an eQTL dataset of larger sample size and the importance of an ancestry-matched eQTL
476 dataset, even at a smaller sample size, in detecting susceptibility genes.

477 We evaluated the presence of a gene-environment interaction with tobacco smoking in our EA data.
478 We found that the association between a PRS (constructed by the lead variants at the 25 loci with
479 genome-wide significance in EA) and LUAD in never-smokers was statistically significantly stronger
480 than in individuals with a history of smoking (Fig. 4). This finding, together with our recent paper
481 showing a stronger association of PRS for LUAD risk in non-coal users than in coal users⁷³, provides
482 evidence that genetic susceptibility may vary by exposure patterns in EA populations.

483 We systematically compared top GWAS findings that had been initially reported in one or the other or
484 both populations. After accounting for differences in MAFs and statistical power as well as the local
485 LD pattern of each locus (500 kb each side of the lead variant), we found that a substantial number of
486 the associations initially reported in EA populations showed no signal in EUR populations. It might
487 reflect causal variants for these loci not being tagged well in the EUR populations. This might also
488 suggest important differences between EA and EUR in the genetic architecture of LUAD samples,
489 which could be caused by differential environmental exposures. Finally, this observation is also
490 consistent with distinct tumor molecular characteristics (e.g., *EGFR* mutation prevalence was higher in

491 Asians than EUR populations) observed in LUAD suggesting different etiologies influenced by genetic
492 and/or environmental factors^{13, 74, 75}.

493 Our genetic architecture analysis suggested that LUAD is a highly polygenic disease. Expanding
494 GWAS of LUAD will continue to identify many risk variants albeit with smaller effect sizes.
495 Moreover, our analysis predicts that the AUC of PRS for EA never-smokers could be improved to
496 66.9% for a GWAS training dataset with 70,000 cases and 70,000 controls that could be further
497 increased with a greater number of controls. Thus, an expanded GWAS in the future can lead to the
498 substantial improvement in knowledge about the underlying genetic architecture of LUAD; increased
499 understanding of how known or suspected lung cancer environmental risk factors interact with genetic
500 susceptibility; and assessment of the potential clinical utility of risk models integrating both genetic
501 and non-genetic risk factors^{76, 77}.

502 There are several limitations in the current study. First, the discovery phase included subjects of
503 diverse EA populations (Mainland China 38.2%, Japan 45.9%) and the replication phase only included
504 subjects from Japan. However, our data did not show evidence of heterogeneity in effect sizes for
505 susceptibility variants between Han Chinese and Japanese populations or across geographic locations
506 (Supplementary Table 5), suggesting a minimal impact for using a single EA population for
507 replication. Second, we were underpowered to conduct formal heritability correlation analyses to
508 compare the genetic architecture in EA and EUR populations stratified by smoking status; larger
509 studies will be needed to conclusively characterize differences. Furthermore, completely elucidating
510 the genetic basis of ancestry differences requires detailed information about age of onset, family
511 history and exposures. Finally, rs4268071 (Table 2) achieved genome-wide significance in the
512 discovery data but replication data were not available. While the significance was primarily driven by
513 Japanese samples (MAF=0.04 in Japanese and <1% in other populations), there was no evidence of

514 heterogeneity in effect estimates across EA populations. Replication is warranted to further establish
515 its etiological role.

516 In conclusion, we identified 12 novel variants in a GWAS of LUAD in EA populations as well as 4
517 novel variants in a multi-ancestry meta-analysis of EA and EUR populations. Colocalization and
518 TWAS analyses using an ancestry-matched lung tissue eQTL dataset identified candidate susceptibility
519 genes with suggested roles in alveolar lineage. At the same time, a large majority of variants identified
520 in the EA GWAS showed no evidence of association in EUR populations. Larger samples sizes with
521 data on environmental risk factors will be needed to further characterize the etiologic differences
522 between these populations. Finally, our genetic architecture analysis suggests that the performance and
523 the clinical utility of the PRS will be substantially improved by larger GWAS in the future.

524

525

526 **Methods**

527 **Ethics statement**

528 All participants provided informed consent according to protocols that were evaluated and approved by
529 the internal review boards of the contributing centers. Protocols used to generate new, unpublished
530 data presented in this paper were approved by the National Cancer Center Institutional Review Board,
531 Japan and the Aichi Cancer Center Ethics Committee, Japan.

532
533 **Overview of study**

534
535 We conducted a two-phase GWAS meta-analysis of LUAD in EA populations, including Female Lung
536 Cancer Consortium in Asia (FLCCA), Nanjing Lung Cancer Study (NJLCS)^{10, 24}, National Cancer
537 Center of Japan (NCC) Research Institute and Aichi Cancer Center (ACC). For the FLCCA study,
538 details of the study design, participating studies, case ascertainment, genotyping, and quality controls
539 have been described in detail⁹. Briefly, this international consortium is composed of Asian women who
540 never smoked and resided in Mainland China, Hong Kong, Singapore, Taiwan, South Korea and Japan
541 at the time of recruitment. All were genotyped using the Illumina 660W, 370K and 610Q microarrays.
542 The NCC study included lung cancer patients from NCC and BioBank Japan (BBJ) and non-cancer
543 controls from the Japan Public Health Center-based Prospective Study and the Japan Multi-
544 Institutional Collaborative Cohort Study, genotyped by Illumina HumanOmniExpress and
545 HumanOmni1-Quad genotyping platforms. The ACC study included lung cancer patients from the
546 Aichi Cancer Center, Kyoto University, Okayama University and Hyogo College of Medicine and
547 non-cancer controls from the Nagahama Study and the Aichi Cancer center. Samples were genotyped
548 by Illumina 610k and Illumina660k platforms^{15, 78}. The NJLCS study at the Nanjing Medical
549 University was based on meta-analysis of three studies: the Nanjing GWAS with subjects from

550 Nanjing and Shanghai, the Beijing study with subjects from Beijing and Wuhan (genotyped by
551 Affymetrix Genome-Wide Human SNP Array 6.0) and the Oncoarray GWAS^{10, 79, 80}.
552 The replication study included cases from multiple sources (BBJ, NCC, Kanagawa Cancer Center,
553 Akita University Hospital, Tokyo Medical and Dental University, Hospital and Gunma University
554 Hospital, and Fukushima Medical University School of Medicine) and non-cancer controls from
555 BioBank Japan. Cases were genotyped using the Invader assay and the control samples in BioBank
556 Japan were genotyped using the Illumina HumanOmniExpress genotyping platform.
557 For the multi-ancestry meta-analyses of LUAD and cross-population comparison of top GWAS
558 findings with both never-smokers and individuals with a history of smoking, we used 11,273 cases and
559 55,483 controls of European ancestry in the Integrative Analysis of Lung Cancer Etiology and Risk
560 team of the International Lung Cancer Consortium (INTEGRAL-ILCCO)¹⁶ (Supplementary Table 8).
561 For the multi-ancestry analysis and cross-population comparisons of smokers, we used European
562 samples genotyped with the OncoArray platform in the ILCCO study (Supplementary Table 8). For the
563 multi-ancestry and cross-population comparisons analysis of never-smokers, we used the GWAS of
564 European never-smoking subjects from Hung et al. (2019)²¹.

565 **Quality control, imputation and association analysis in EA populations**

566 For each study, SNPs with minor allele frequency (MAF) < 0.01, Hardy-Weinberg Equilibrium (HWE)
567 p-value < 10⁻⁶ in controls were removed; subjects with missing rate > 3%, sex discrepancy, or
568 displaying non-East Asian ancestry based on principal component analysis scores were removed.
569 Moreover, for any pairs of subjects estimated to be related with identity by descent $\text{pihat} > 0.10$ using
570 PLINK (V2.0), we removed one subject. Imputation was performed using IMPUTE2 and the 1000
571 Genomes Project East Asian samples (Phase 3) as reference. After imputation, SNPs with imputation
572 quality score ≥ 0.5 were used for association analysis in each study. Logistic regression under an

573 additive model was performed using SNPTest (V2) or PLINK2 based on imputed genotypic dosage
574 data adjusting for smoking (if both smokers and never smokers were present) and PCA scores to
575 control for population stratification. Meta-analysis was performed using inverse-variance weighted
576 fixed effects methods. All p-values were two-sided. We consider the following variants as novel for
577 the GWAS in EA: (1) the lead variant with $p < 5 \times 10^{-8}$ in a locus that has not been previously reported
578 in either EA or EUR populations, or (2) a secondary variant with $p < 5 \times 10^{-8}$ conditioning on the lead
579 variant in a previously reported locus in either EA or EUR populations with the requirement that the
580 LD $R^2 \leq 0.2$ between the secondary and the lead variants in both populations.
581 LDSC²⁷ was used to estimate the heritability attributed to genome-wide common variants and to assess
582 the potential inflation due to insufficient correction of population stratification. LDSC was also used to
583 estimate the genetic correlation of LUAD between never-smokers and individuals with a history of
584 smoking in each population. We used POPCORN⁶² to estimate the genetic correlation between EA and
585 EUR populations because LD patterns are expected to be different. To account for the difference of
586 allele frequencies in the two populations, we also used POPCORN to estimate the cross-population
587 genetic-impact correlation that was defined as the correlation of population specific phenotypic
588 variance explained by each SNP.

589 **Conditional analysis and fine mapping**

590 To identify independently associated SNPs at an established susceptibility locus, we performed
591 conditional analysis using software Genome-wide Complex Trait Analysis (GCTA)⁸¹ based on the
592 GWAS meta-analysis summary results of EA populations. LD for the conditional analysis was
593 calculated using a reference population of 4,544 controls from the FLCCA study to achieve a desirable
594 accuracy. Here, genotypes for FLCCA were imputed using IMPUTE2 and the 1000 Genomes Project
595 (Phase 3) reference samples with EA ancestry. SNPs with imputation quality < 0.5 were excluded from

596 the reference set for conditional analysis. Conditional analysis was restricted to 14 loci with lead SNPs
597 achieving genome-wide significance in the discovery-phase meta-analysis. We did not perform
598 conditional analyses for other new SNPs that did not achieve genome-wide significance in the
599 discovery-phase meta-analysis because secondary SNPs would not survive multiple testing correction.
600 Conditional analysis was restricted to SNPs less than 500kb from the lead SNP of each locus. To
601 identify multiple potentially independent SNPs in one locus, we performed stepwise conditional
602 analysis using GCTA. All SNPs identified with $P < 5 \times 10^{-8}$ and the lead SNP of the locus were put into
603 one model to derive the joint estimate of ORs, appropriately adjusting for LD among all SNPs. Only
604 SNPs with p-value $< 5 \times 10^{-8}$ in both conditional and joint analyses were considered to be independently
605 associated SNPs.

606 For 11 out of the 14 loci with genome-wide significance in the discovery phase, we performed a
607 Bayesian fine-mapping analysis using FINEMAP³³ to nominate 95% credible set variants using the
608 same set of imputed genotypes of 4,544 FLCCA subjects as an LD reference. We did not perform fine-
609 mapping analysis for two loci in MHC regions, because of the complex and extensive LD patterns in
610 this region. We also excluded the locus at 7q31 because the lead SNP, rs4268071, had MAF $<1\%$ in our
611 LD reference population. MAF of this variant is 4% in the Japanese populations (45.8% of cases and
612 74.5% of controls in the discovery set) but $<1\%$ in other EA populations included in our study. For
613 FINEMAP analysis, we tested the variants within ± 500 kb of the lead SNP and set the number of
614 maximum causal variants as the number of independent signals ($P \leq 10^{-5}$) observed in the conditional
615 analysis for each locus.

616 **Proportion of familial risk explained**

617 We considered a set of identified variants for LUAD. For SNP t , we defined p_t as the frequency of
618 the risk allele and OR_t as the estimated per-allele odds ratio. Under a multiplicative model, the

619 fraction of the familial risk explained by the set of SNPs was calculated as $\sum_t \log(\lambda_t) / \log(\lambda_0)$,
620 where λ_0 is the observed familial risk to the first degree of LUAD cases and λ_t is the familial risk due
621 to the t^{th} SNP:

$$622 \quad \lambda_t = \frac{p_t OR_t^2 + (1-p_t)}{(p_t OR_t + 1 - p_t)^2} \quad (1)$$

623 **Heritability partitioning in functional classes and tissue-specific analyses**

624 Stratified LD score regression (sLDSC)⁸² was conducted to identify functional annotations enriched for
625 LUAD heritability using summary statistics from the discovery phase of meta-analysis in EA
626 populations. In addition to the functional annotations provided by the sLDSC package, we also
627 analyzed the gene sets defined by smoking studies: differentially expressed genes in peripheral blood
628 mononuclear cells upon nicotine treatment (“PBMC nicotine” gene set) from Moyerbrailean et al.⁸³,
629 those in non-tumorous lungs between current- and never-smokers (“Lung smoking” gene set) from
630 Bosse et al.⁸⁴, and those in normal bronchial airway epithelial cells between current- and never-
631 smokers (“Airway smoking” gene set) from Beane et al.⁸⁵. An annotation was considered to be
632 significantly enriched for LUAD heritability if $FDR < 0.05$.

633 We then performed sLDSC to prioritize relevant tissue types (lung, blood/immune, and brain/CNS)
634 using tissue-specific expressed genes from GTEx v6p (53 tissue types) and other public expression
635 datasets (152 tissue types), as well as tissue-specific chromatin annotations from EnTEX (111
636 annotations in 26 tissue types) and Roadmap dataset (378 annotations in 85 tissue types) as described
637 by Finucane and colleagues³⁷. We used GTEx v6p expression data based on a comparison with v8
638 data, where a median of 83% of tissue-specific differentially expressed genes were shared between two
639 versions. In general, we did not find significant enrichment for individual annotations after adjusting
640 for the multiple testing. To increase the power of prioritizing relevant tissues (lung, blood/immune, and
641 brain/CNS), we performed an aggregated analysis to test if p-values from one tissue (e.g., lung) tended

642 to be smaller than those from the other two tissue groups (blood/immune, and brain/CNS) using the
643 Wilkinson rank test.

644 **eQTL colocalization analysis and TWAS**

645 EA lung eQTL dataset is based on a cohort of 115 never-smoking LUAD patients from Taiwan,
646 referred to as LCTCNS (Lung cancer tissue cohort of never-smokers). Expression array data was
647 obtained for non-tumor lung tissues of these patients using the Illumina WG-DASL HumanRef-8 v3 or
648 HumanHT-12 v4 BeadChip (Illumina Inc.) (Gene Expression Omnibus accession number
649 GSE46539)⁸⁶. Genotype data from buffy coat DNA was obtained using the Illumina Human 660W
650 Quad BeadChip. A systematic quality control for the genotype data was performed as previously
651 described¹² (SNPs were excluded if call rate < 90%, MAF < 5%, or $P < 0.0001$ based on the Hardy-
652 Weinberg equilibrium test. Samples were excluded if call rate < 90%, sex discrepancies based on the X
653 chromosome heterozygosity, contaminated samples with high heterozygosity scores, or first or second-
654 degree relatives), and imputation was carried out using Minimac4 (V4.0.3) with the 1000 Genomes
655 reference set (all populations). For eQTL analysis, expression data was processed for background
656 correction as previously described⁸⁶. Briefly, we kept the probes that are present in both the BeadChip
657 platforms and further removed those with low expression levels (detection $p > 0.05$). Based on the data
658 at the remaining 24,216 probes, we applied model-based background correction. Log₂-transformed
659 expression levels of 24,216 probes were then used to obtain 20 latent factors based on probabilistic
660 estimation of expression residuals (PEER) while specifying batch, sex, age, medical operation status,
661 RNA integrity number, and RNA input quantity as known confounders. The expression residuals from
662 PEER were then inverse rank transformed to the standard normal distribution (the inverse rank
663 transformed residuals) and were used as the dependent variable in the expression levels for eQTL
664 analysis. eQTL analysis was conducted for 29 GWAS lead SNPs (all EA loci including discovery,

665 replication, and conditional signals plus new loci from the multi-ancestry GWAS). In LCTCNS, all
666 these SNPs have a MAF of > 0.01 . For each GWAS lead SNP, its association with each probe located
667 within ± 500 kb of the SNP was tested using an additive linear model where the dependent variable
668 was the expression level as described above and the independent variable was the effect allele count.
669 Based on the resulting p-values of these eQTL analyses for all 29 SNPs, the corresponding Benjamini–
670 Hochberg FDR was calculated. Colocalization analysis was performed using eCAVIAR³⁸ and
671 HyPrColoc³⁹ via ezQTL platform for eight GWAS lead SNP-eQTL gene pairs displaying FDR < 0.05
672 in LCTCNS (Supplementary Data 5). For each of these eight SNP-probe pairs, we further examined
673 the association between the probe and SNPs within ± 100 kb of the lead SNP using Matrix eQTL to
674 obtain the summary statistics as an input to ezQTL for colocalization analysis using HyPrColoc and
675 eCAVIAR. For loci on MHC regions, ± 10 kb window was used for computational efficiency of
676 colocalization analyses. LD matrix was obtained from 1000 Genomes EA populations. For HyPrColoc,
677 posterior probability of > 0.7 was used as a cutoff for colocalization. For eCAVIAR analysis,
678 colocalization posterior probability (CLPP) score > 0.01 was used as a cutoff for colocalization.
679 For TWAS, we adopted FUSION⁸⁷ using LCTCNS or GTEx v8 lung eQTL data and summary
680 statistics of EA discovery GWAS meta-analysis. We computed weights using the elastic-net regression
681 (enet) model for 24,216 expression probes (LCTCNS) or 24,687 genes (GTEx v8 lung) and *cis*-SNPs
682 within 500 kb of the gene for each probe. LD matrix was obtained from 1000 Genomes EA
683 populations. We performed association analysis for 1,875 expression probes (LCTCNS) or 5,534 genes
684 (GTEx v8 lung) with cross-validation cutoff of $R^2 > 0.05$ based on the elastic net model. We defined a
685 significant transcriptome-wide association as TWAS $P < 2.6 \times 10^{-5}$ ($0.05/1,875$; LCTCNS) or $P < 9 \times$
686 10^{-6} ($0.05/5,534$; GTEx v8 lung) based on Bonferroni correction. For two loci passing this cutoff from
687 LCTCNS analysis (*ELF5* and *FADSI*), we further performed conditional analysis as implemented in

688 FUSION by conditioning the GWAS signal on the predicted expression of the probe with the best
689 TWAS P-value.

690
691 **Mendelian randomization**

692 We performed MR analysis to investigate the potential causal relationship between telomere length
693 and the risk of LUAD. MR analysis was based on 46 common SNPs identified in a recent multi-
694 ancestry meta-analysis of telomere length in the TOPMed⁵⁵ study. The original paper identified 48
695 variants associated with telomere length that collectively explained 4.35% of telomere length variance;
696 two of them at the *TERT* locus were excluded using the LD filter $R^2 < 0.05$ that together explained
697 0.61% of the telomere length variance; the remaining 46 variants included in our MR analysis
698 explained 3.74% of telomere length variance. Because there was no significant heterogeneity of effect
699 sizes on telomere length across populations (Table S4 in Taub et al.⁵⁵), the primary MR analyses were
700 based on the estimated effect sizes combining all samples in the TOPMed study in a joint model for
701 telomere length. Analyses were based on MR PRESSO⁸⁸, a powerful and robust approach designed to
702 deal with widespread horizontal pleiotropy. This approach uses a formal test framework to (1) detect
703 the presence of horizontal pleiotropy, (2) detect variant outliers, (3) evaluate distortion, and (4) re-
704 estimate causal effect sizes after removing potentially problematic variants. According to simulations,
705 this approach is best suited when horizontal pleiotropy occurs in $< 50\%$ of instruments. This approach
706 identified 5-7 outlier variants in our data. The estimated β from MR analysis was converted as OR,
707 interpreted as risk increase per standard deviation (640 base pairs⁸⁹) increase of the genetic predicted
708 telomere length.

709 **Testing the interaction between polygenic risk score and smoking status**

710 We investigated whether the PRS, which was calculated based on 25 independent SNPs associated
711 with LUAD in EA populations (Supplementary Table 4, excluding three variants identified by

712 conditional analysis), interacted with smoking status for LUAD risk. Because we have only GWAS
 713 summary statistics instead of individual-level data for smokers and never-smokers, we developed a
 714 statistical method for testing the interaction using summary statistics separately from smokers and
 715 never-smokers. Suppose that we have n^{1+} smoking cases, n^{0+} never-smoking cases, n^{1-} smoking
 716 controls and n^{0-} never-smoking controls. Let x_{it}^{s+} and x_{jt}^{s-} be the genotype of SNP t for the i^{th} case
 717 and the j^{th} control, where $s = 1$ indicates smokers and 0 indicates never-smokers. Given smoking
 718 status s , we define $PRS_i^{s+} = \sum_{t=1}^T \beta_t x_{it}^{s+}$ and $PRS_j^{s-} = \sum_{t=1}^T \beta_t x_{jt}^{s-}$ as the PRS for cases and controls,
 719 respectively. For smokers ($s = 1$), the association between PRS and disease risk can be quantified as:

$$726 \quad \Delta_1 = \frac{1}{n^{1+}} \sum_{i=1}^{n^{1+}} PRS_i^{1+} - \frac{1}{n^{1-}} \sum_{j=1}^{n^{1-}} PRS_j^{1-}, \quad (2)$$

720 the difference of average PRS between cases and controls. Similarly, we define Δ_0 to be the difference
 721 of average PRS between cases and controls for never-smokers. Testing the PRS*smoking interaction
 722 can be done using $Z = \frac{\Delta_1 - \Delta_0}{\sqrt{\text{var}(\Delta_1^2) + \text{var}(\Delta_0^2)}}$. Under the null hypothesis of no interaction for all variants,

723 $Z \sim N(0,1)$ asymptotically. Assuming SNPs are independent, we derive $Z = \sum_{t=1}^T (w_t^1 z_t^1 - w_t^0 z_t^0)$,
 724 where z_t^s is the z-score for testing association for SNP t in subjects with smoking status s . The weight
 725 is given as

$$727 \quad w_t^s = \frac{\beta_t \sqrt{\frac{(\sigma_t^{s+})^2}{n_+^s} + \frac{(\sigma_t^{s-})^2}{n_-^s}}}{\sqrt{\sum_{t=1}^T \beta_t^2 \left(\frac{(\sigma_t^{1+})^2}{n_+^1} + \frac{(\sigma_t^{1-})^2}{n_-^1} + \frac{(\sigma_t^{0+})^2}{n_+^0} + \frac{(\sigma_t^{0-})^2}{n_-^0} \right)}}. \quad (3)$$

728 Here, $(\sigma_t^{s+})^2$ and $(\sigma_t^{s-})^2$ are the genotypic variances for SNP t in cases and controls, respectively.
 729 We note that both discovery and replication data are included for testing PRS smoking interaction

730 novel variants included in our PRS to maximize the power of statistical testing. In particular, only the
731 discovery data were available and included for previously identified variants; both discovery and
732 replication data were included for new variants to increase the statistical power. To do this, w_t^s was
733 modified to have SNP-specific sample sizes. All analyses were done using R (x64 4.1.0).

734 **GENESIS analysis for projecting yield of future expanded studies**

735 The genetic architecture of a disease is defined as the number of susceptibility SNPs and the
736 distribution of their effect sizes²⁶. When these parameters are estimated, one can estimate the number
737 of variants achieving genome-wide significance and the accuracy of a polygenic risk model trained
738 using a GWAS with a given sample size. In the current study, we estimated the genetic architecture
739 using GENESIS (GENetic ESTimation and Inference in Structured samples)²⁶ based on the GWAS
740 summary statistics with LD scores calculated based on the genotypes of the subjects of EA ancestry in
741 the 1000 Genomes Project. Since GENESIS requires a large sample size to derive reliable estimates,
742 we performed analysis only for never-smokers in EA. The three-component model

743 $\beta_m \sim \pi p_1 N(0, \sigma_1^2) + \pi p_2 N(0, \sigma_2^2) + (1 - \pi) \delta_0$ best fit the never-smoker data in EA, where β_m
744 represents effects sizes, π denotes the fraction of truly associated variants in the genome, δ_0 denotes
745 the point mass at zero, σ_i^2 denotes the variance of effect sizes for the i^{th} component, πp_i ($i = 1, 2$)
746 represents the fraction of variants with effect size following $N(0, \sigma_i^2)$. Based on this estimated genetic
747 architecture, we calculated the expected number of variants reaching genome-wide significance for a
748 given GWAS and calculated the expected area under the receiver operating characteristic curve (AUC)
749 for an additive polygenic risk prediction model built based on a discovery GWAS for a given sample
750 size. The uncertainty of the AUC was induced by the uncertainty in the estimated parameters in
751 GENESIS ($\Gamma = (\pi, p_1, p_2, \sigma_1^2, \sigma_2^2)$) because of the limited sample size in our summary data. We used a
752 resampling approach to estimate the standard error of AUC. Briefly, we randomly simulated 1000 sets

753 of parameters Γ^k given the estimated $\hat{\Gamma}$ and the estimated covariance matrix, and calculated AUC_k for
754 each simulated parameter Γ^k for a given sample size. The standard error was calculated based on the
755 1000 sets of AUC values.

756

757 **Data Availability**

758 The GWAS data for the FLCCA study is available at dbGap under accession phs000716.v1.p1
759 (https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000716.v1.p1). The
760 GWAS data for the Japanese populations are available at the Integrative Disease Omics Database
761 (<https://integbio.jp/dbcatalog/en/record/nbdc00071>) under accession code GWAS031 and BioBank
762 Japan (<https://biobankjp.org/en/>). The GWAS data for the European populations contributing to this
763 study are available at dbGap under accession phs000877.v1.p1 (Transdisciplinary Research Into
764 Cancer of the Lung (TRICL), [https://www.ncbi.nlm.nih.gov/projects/gap/cgi-](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000876.v2.p1)
765 [bin/study.cgi?study_id=phs000876.v2.p1](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000876.v2.p1)), phs001273.v3.p2 (Oncoarray Consortium,
766 https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001273.v3.p2). The
767 expression data of LCTCNS (Lung cancer tissue cohort of never-smokers) were publicly available at
768 Gene Expression Omnibus under accession number GSE46539. The expression and eQTL data from
769 GTEx (v6 and v8) were obtained from <https://gtexportal.org/home/datasets>. Full TWAS results are
770 included in Supplementary Data 6. The summary statistics for the meta-analysis in East Asian
771 populations with $p \leq 10^{-4}$ are included in Supplementary Data 10.

772

773

774

775

776

777

778

779

780 **References.**

- 781
- 782 1. Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A. Global cancer statistics
783 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185
784 countries. *CA Cancer J Clin* **68**, 394-424 (2018).
785
- 786 2. Cheng TY, Cramb SM, Baade PD, Youlden DR, Nwogu C, Reid ME. The International
787 Epidemiology of Lung Cancer: Latest Trends, Disparities, and Tumor Characteristics. *J Thorac*
788 *Oncol* **11**, 1653-1671 (2016).
789
- 790 3. Barta JA, Powell CA, Wisnivesky JP. Global Epidemiology of Lung Cancer. *Ann Glob Health*
791 **85**, (2019).
792
- 793 4. Cao M, Chen W. Epidemiology of lung cancer in China. *Thorac Cancer* **10**, 3-7 (2019).
794
- 795 5. Kinoshita FL, Ito Y, Nakayama T. Trends in Lung Cancer Incidence Rates by Histological
796 Type in 1975-2008: A Population-Based Study in Osaka, Japan. *J Epidemiol* **26**, 579-586
797 (2016).
798
- 799 6. Landi MT, *et al.* Tracing Lung Cancer Risk Factors Through Mutational Signatures in Never-
800 Smokers. *Am J Epidemiol* **190**, 962-976 (2021).
801
- 802 7. Sisti J, Boffetta P. What proportion of lung cancer in never-smokers can be attributed to known
803 risk factors? *Int J Cancer* **131**, 265-275 (2012).
804
- 805 8. Schottenfeld D, Fraumeni JF. *Cancer epidemiology and prevention*, 3rd edn. Oxford University
806 Press (2006).
807
- 808 9. Lan Q, *et al.* Genome-wide association analysis identifies new lung cancer susceptibility loci in
809 never-smoking women in Asia. *Nat Genet* **44**, 1330-1335 (2012).
810
- 811 10. Hu Z, *et al.* A genome-wide association study identifies two new lung cancer susceptibility loci
812 at 13q12.12 and 22q12.2 in Han Chinese. *Nat Genet* **43**, 792-796 (2011).
813
- 814 11. Wang Z, *et al.* Meta-analysis of genome-wide association studies identifies multiple lung
815 cancer susceptibility loci in never-smoking Asian women. *Hum Mol Genet* **25**, 620-629 (2016).
816
- 817 12. Hsiung CA, *et al.* The 5p15.33 locus is associated with risk of lung adenocarcinoma in never-
818 smoking females in Asia. *PLoS Genet* **6**, (2010).
819
- 820 13. Seow WJ, *et al.* Association between GWAS-identified lung adenocarcinoma susceptibility loci
821 and EGFR mutations in never-smoking Asian women, and comparison with findings from
822 Western populations. *Hum Mol Genet* **26**, 454-465 (2017).
823
- 824 14. Miki D, *et al.* Variation in TP63 is associated with lung adenocarcinoma susceptibility in
825 Japanese and Korean populations. *Nat Genet* **42**, 893-896 (2010).
826

- 827 15. Shiraishi K, *et al.* A genome-wide association study identifies two new susceptibility loci for
828 lung adenocarcinoma in the Japanese population. *Nat Genet* **44**, 900-903 (2012).
829
- 830 16. McKay JD, *et al.* Large-scale association analysis identifies new lung cancer susceptibility loci
831 and heterogeneity in genetic susceptibility across histological subtypes. *Nat Genet* **49**, 1126-
832 1132 (2017).
833
- 834 17. Wang Y, *et al.* Rare variants of large effect in BRCA2 and CHEK2 affect risk of lung cancer.
835 *Nat Genet* **46**, 736-741 (2014).
836
- 837 18. Landi MT, *et al.* A genome-wide association study of lung cancer identifies a region of
838 chromosome 5p15 associated with risk for adenocarcinoma. *Am J Hum Genet* **85**, 679-691
839 (2009).
840
- 841 19. Amos CI, *et al.* Genome-wide association scan of tag SNPs identifies a susceptibility locus for
842 lung cancer at 15q25.1. *Nat Genet* **40**, 616-622 (2008).
843
- 844 20. Wang Y, *et al.* Common 5p15.33 and 6p21.33 variants influence lung cancer risk. *Nat Genet*
845 **40**, 1407-1409 (2008).
846
- 847 21. Hung RJ, *et al.* Lung Cancer Risk in Never-Smokers of European Descent is Associated With
848 Genetic Variation in the 5p15.33 TERT-CLPTM1L1 Region. *J Thorac Oncol* **14**, 1360-1369
849 (2019).
850
- 851 22. Hung RJ, *et al.* A susceptibility locus for lung cancer maps to nicotinic acetylcholine receptor
852 subunit genes on 15q25. *Nature* **452**, 633-637 (2008).
853
- 854 23. Wu C, *et al.* Genetic variants on chromosome 15q25 associated with lung cancer risk in
855 Chinese populations. *Cancer Res* **69**, 5065-5072 (2009).
856
- 857 24. Dai J, *et al.* Identification of risk loci and a polygenic risk score for lung cancer: a large-scale
858 prospective cohort study in Chinese populations. *Lancet Respir Med* **7**, 881-891 (2019).
859
- 860 25. Byun J. HY, Li Y., Xia J., Long E., Choi J. Trans-ethnic genome-wide meta-analysis of 35,732
861 cases and 34,424 controls identifies novel genomic cross-ancestry loci contributing to lung
862 cancer susceptibility. (2021).
863
- 864 26. Zhang Y, Qi G, Park JH, Chatterjee N. Estimation of complex effect-size distributions using
865 summary-level statistics from genome-wide association studies across 32 complex traits. *Nat*
866 *Genet* **50**, 1318-1326 (2018).
867
- 868 27. Bulik-Sullivan B, *et al.* An atlas of genetic correlations across human diseases and traits.
869 *Nature Genetics* **47**, 1236-1241 (2015).
870
- 871 28. Dong J, *et al.* Fine mapping of chromosome 5p15.33 identifies novel lung cancer susceptibility
872 loci in Han Chinese. *Int J Cancer* **141**, 447-456 (2017).

- 873
874 29. Matakidou A, Eisen T, Houlston RS. Systematic review of the relationship between family
875 history and lung cancer risk. *Br J Cancer* **93**, 825-833 (2005).
876
- 877 30. Nagai A, *et al.* Overview of the BioBank Japan Project: Study design and profile. *J Epidemiol*
878 **27**, S2-S8 (2017).
879
- 880 31. Landi MT, *et al.* A genome-wide association study of lung cancer identifies a region of
881 chromosome 5p15 associated with risk for adenocarcinoma. *Am J Hum Genet* **85**, 679-691
882 (2009).
883
- 884 32. Timofeeva MN, *et al.* Influence of common genetic variation on lung cancer risk: meta-
885 analysis of 14 900 cases and 29 485 controls. *Hum Mol Genet* **21**, 4980-4995 (2012).
886
- 887 33. Benner C, Spencer CC, Havulinna AS, Salomaa V, Ripatti S, Pirinen M. FINEMAP: efficient
888 variable selection using summary data from genome-wide association studies. *Bioinformatics*
889 **32**, 1493-1501 (2016).
890
- 891 34. Boyle AP, *et al.* Annotation of functional variation in personal genomes using RegulomeDB.
892 *Genome Res* **22**, 1790-1797 (2012).
893
- 894 35. Ward LD, Kellis M. HaploReg v4: systematic mining of putative causal variants, cell types,
895 regulators and target genes for human complex traits and disease. *Nucleic Acids Res* **44**, D877-
896 881 (2016).
897
- 898 36. Breeze CE, *et al.* Integrative analysis of 3604 GWAS reveals multiple novel cell type-specific
899 regulatory associations. *Genome Biol* **23**, 13 (2022).
900
- 901 37. Finucane HK, *et al.* Heritability enrichment of specifically expressed genes identifies disease-
902 relevant tissues and cell types. *Nat Genet* **50**, 621-629 (2018).
903
- 904 38. Hormozdiari F, *et al.* Colocalization of GWAS and eQTL Signals Detects Target Genes. *Am J*
905 *Hum Genet* **99**, 1245-1260 (2016).
906
- 907 39. Foley CN, *et al.* A fast and efficient colocalization algorithm for identifying shared genetic risk
908 factors across multiple traits. *Nat Commun* **12**, 764 (2021).
909
- 910 40. Pan G, *et al.* PATZ1 down-regulates FADS1 by binding to rs174557 and is opposed by
911 SP1/SREBP1c. *Nucleic Acids Res* **45**, 2408-2422 (2017).
912
- 913 41. Glaser C, Heinrich J, Koletzko B. Role of FADS1 and FADS2 polymorphisms in
914 polyunsaturated fatty acid metabolism. *Metabolism* **59**, 993-999 (2010).
915
- 916 42. Zhao R, *et al.* FADS1 promotes the progression of laryngeal squamous cell carcinoma through
917 activating AKT/mTOR signaling. *Cell Death Dis* **11**, 272 (2020).
918

- 919 43. Oakes SR, *et al.* The Ets transcription factor Elf5 specifies mammary alveolar cell fate. *Genes*
920 *Dev* **22**, 581-586 (2008).
921
- 922 44. Tan AC. Targeting the PI3K/Akt/mTOR pathway in non-small cell lung cancer (NSCLC).
923 *Thorac Cancer* **11**, 511-518 (2020).
924
- 925 45. Cheng H, Shcherba M, Pendurti G, Liang Y, Piperdi B, Perez-Soler R. Targeting the
926 PI3K/AKT/mTOR pathway: potential for lung cancer treatment. *Lung Cancer Manag* **3**, 67-75
927 (2014).
928
- 929 46. Amodio N, *et al.* Oncogenic role of the E3 ubiquitin ligase NEDD4-1, a PTEN negative
930 regulator, in non-small-cell lung carcinomas. *Am J Pathol* **177**, 2622-2634 (2010).
931
- 932 47. Cantley LC, Neel BG. New insights into tumor suppression: PTEN suppresses tumor formation
933 by restraining the phosphoinositide 3-kinase/AKT pathway. *Proc Natl Acad Sci U S A* **96**,
934 4240-4245 (1999).
935
- 936 48. Sin DD, *et al.* Pro-surfactant protein B as a biomarker for lung cancer prediction. *J Clin Oncol*
937 **31**, 4536-4543 (2013).
938
- 939 49. Lehner B, Sanderson CM. A protein interaction framework for human mRNA degradation.
940 *Genome Res* **14**, 1315-1323 (2004).
941
- 942 50. Moon DH, *et al.* Poly(A)-specific ribonuclease (PARN) mediates 3'-end maturation of the
943 telomerase RNA component. *Nat Genet* **47**, 1482-1488 (2015).
944
- 945 51. Stanley SE, *et al.* Loss-of-function mutations in the RNA biogenesis factor NAF1 predispose to
946 pulmonary fibrosis-emphysema. *Sci Transl Med* **8**, 351ra107 (2016).
947
- 948 52. Vannier JB, Pavicic-Kaltenbrunner V, Petalcorin MI, Ding H, Boulton SJ. RTEL1 dismantles T
949 loops and counteracts telomeric G4-DNA to maintain telomere integrity. *Cell* **149**, 795-806
950 (2012).
951
- 952 53. Sarek G, Vannier JB, Panier S, Petrini JHJ, Boulton SJ. TRF2 recruits RTEL1 to telomeres in S
953 phase to promote t-loop unwinding. *Mol Cell* **57**, 622-635 (2015).
954
- 955 54. Miyake Y, *et al.* RPA-like mammalian Ctc1-Stn1-Ten1 complex binds to single-stranded DNA
956 and protects telomeres independently of the Pot1 pathway. *Mol Cell* **36**, 193-206 (2009).
957
- 958 55. Margaret A. Taub MPC, Rebecca Keener, Kruthika R. Iyer, Joshua S. Weinstock, Lisa R.
959 Yanek, John Lane, Tyne W. Miller-Fleming, Jennifer A. Brody, Laura M. Raffield, Caitlin P.
960 McHugh, Deepti Jain, Stephanie M. Gogarten, Cecelia A. Laurie, *et al.* Novel genetic
961 determinants of telomere length from a trans-ethnic analysis of 109,122 whole genome
962 sequences in TOPMed. *Cell Genomics* **2**, 100084 (2022).
963

- 964 56. Verbanck M, Chen CY, Neale B, Do R. Publisher Correction: Detection of widespread
965 horizontal pleiotropy in causal relationships inferred from Mendelian randomization between
966 complex traits and diseases. *Nat Genet* **50**, 1196 (2018).
967
- 968 57. Machiela MJ, *et al.* Genetic variants associated with longer telomere length are associated with
969 increased lung cancer risk among never-smoking women in Asia: a report from the female lung
970 cancer consortium in Asia. *Int J Cancer* **137**, 311-319 (2015).
971
- 972 58. Telomeres Mendelian Randomization C, *et al.* Association Between Telomere Length and Risk
973 of Cancer and Non-Neoplastic Diseases: A Mendelian Randomization Study. *JAMA Oncol* **3**,
974 636-651 (2017).
975
- 976 59. Zhang C, *et al.* Genetic determinants of telomere length and risk of common cancers: a
977 Mendelian randomization study. *Hum Mol Genet* **24**, 5356-5366 (2015).
978
- 979 60. Seow WJ, *et al.* Telomere length in white blood cell DNA and lung cancer: a pooled analysis of
980 three prospective cohorts. *Cancer Res* **74**, 4090-4098 (2014).
981
- 982 61. Timofeeva MN, *et al.* Influence of common genetic variation on lung cancer risk: meta-
983 analysis of 14 900 cases and 29 485 controls. *Hum Mol Genet* **21**, 4980-4995 (2012).
984
- 985 62. Brown BC, Asian Genetic Epidemiology Network Type 2 Diabetes C, Ye CJ, Price AL, Zaitlen
986 N. Transethnic Genetic-Correlation Estimates from Summary Statistics. *Am J Hum Genet* **99**,
987 76-88 (2016).
988
- 989 63. Zhang R, *et al.* A genome-wide gene-environment interaction analysis for tobacco smoke and
990 lung cancer susceptibility. *Carcinogenesis* **35**, 1528-1535 (2014).
991
- 992 64. Li Y, *et al.* Genome-wide interaction study of smoking behavior and non-small cell lung cancer
993 risk in Caucasian population. *Carcinogenesis* **39**, 336-346 (2018).
994
- 995 65. Zhang YD, *et al.* Assessment of polygenic architecture and risk prediction based on common
996 variants across fourteen cancers. *Nat Commun* **11**, 3353 (2020).
997
- 998 66. Maas P, *et al.* Breast Cancer Risk From Modifiable and Nonmodifiable Risk Factors Among
999 White Women in the United States. *JAMA Oncol* **2**, 1295-1302 (2016).
1000
- 1001 67. Angelidis I, *et al.* An atlas of the aging lung mapped by single cell transcriptomics and deep
1002 tissue proteomics. *Nat Commun* **10**, 963 (2019).
1003
- 1004 68. Plantier L, *et al.* Activation of sterol-response element-binding proteins (SREBP) in alveolar
1005 type II cells enhances lipogenesis causing pulmonary lipotoxicity. *J Biol Chem* **287**, 10099-
1006 10114 (2012).
1007

- 1008 69. Choi YS, Chakrabarti R, Escamilla-Hernandez R, Sinha S. Elf5 conditional knockout mice
1009 reveal its role as a master regulator in mammary alveolar development: failure of Stat5
1010 activation and functional differentiation in the absence of Elf5. *Dev Biol* **329**, 227-241 (2009).
1011
- 1012 70. Chakrabarti R, *et al.* Elf5 inhibits the epithelial-mesenchymal transition in mammary gland
1013 development and breast cancer metastasis by transcriptionally repressing Snail2. *Nat Cell Biol*
1014 **14**, 1212-1222 (2012).
1015
- 1016 71. Emilsson V, *et al.* Co-regulatory networks of human serum proteins link genetics to disease.
1017 *Science* **361**, 769-773 (2018).
1018
- 1019 72. Hamvas A, *et al.* Comprehensive genetic variant discovery in the surfactant protein B gene.
1020 *Pediatr Res* **62**, 170-175 (2007).
1021
- 1022 73. Blechter B, *et al.* Sub-multiplicative interaction between polygenic risk score and household
1023 coal use in relation to lung adenocarcinoma among never-smoking women in Asia. *Environ Int*
1024 **147**, 105975 (2021).
1025
- 1026 74. Chen J, *et al.* Genomic landscape of lung adenocarcinoma in East Asians. *Nat Genet* **52**, 177-
1027 186 (2020).
1028
- 1029 75. Carrot-Zhang J, *et al.* Genetic Ancestry Contributes to Somatic Mutations in Lung Cancers
1030 from Admixed Latin American Populations. *Cancer Discov* **11**, 591-598 (2021).
1031
- 1032 76. Katki HA, Kovalchik SA, Berg CD, Cheung LC, Chaturvedi AK. Development and Validation
1033 of Risk Models to Select Ever-Smokers for CT Lung Cancer Screening. *JAMA* **315**, 2300-2311
1034 (2016).
1035
- 1036 77. Chien LH, *et al.* Predicting Lung Cancer Occurrence in Never-Smoking Females in Asia:
1037 TNSF-SQ, a Prediction Model. *Cancer Epidemiol Biomarkers Prev* **29**, 452-459 (2020).
1038
- 1039 78. Shiraishi K, *et al.* Association of variations in HLA class II and other loci with susceptibility to
1040 EGFR-mutated lung adenocarcinoma. *Nat Commun* **7**, 12451 (2016).
1041
- 1042 79. Dong J, *et al.* Association analyses identify multiple new lung cancer susceptibility loci and
1043 their interactions with smoking in the Chinese population. *Nat Genet* **44**, 895-899 (2012).
1044
- 1045 80. Wang L, *et al.* Genetically determined height was associated with lung cancer risk in East
1046 Asian population. *Cancer Med*, (2018).
1047
- 1048 81. Yang J, *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics
1049 identifies additional variants influencing complex traits. *Nat Genet* **44**, 369-375, S361-363
1050 (2012).
1051
- 1052 82. Finucane HK, *et al.* Partitioning heritability by functional annotation using genome-wide
1053 association summary statistics. *Nat Genet*, (2015).

- 1054
1055 83. Moyerbrailean GA, *et al.* High-throughput allele-specific expression across 250 environmental
1056 conditions. *Genome Res* **26**, 1627-1638 (2016).
1057
- 1058 84. Bosse Y, *et al.* Molecular signature of smoking in human lung tissues. *Cancer Res* **72**, 3753-
1059 3763 (2012).
1060
- 1061 85. Beane J, *et al.* Characterizing the impact of smoking and lung cancer on the airway
1062 transcriptome using RNA-Seq. *Cancer Prev Res (Phila)* **4**, 803-817 (2011).
1063
- 1064 86. Chang IS, *et al.* Genetic Modifiers of Progression-Free Survival in Never-Smoking Lung
1065 Adenocarcinoma Patients Treated with First-Line Tyrosine Kinase Inhibitors. *Am J Respir Crit*
1066 *Care Med* **195**, 663-673 (2017).
1067
- 1068 87. Gusev A, *et al.* Integrative approaches for large-scale transcriptome-wide association studies.
1069 *Nat Genet* **48**, 245-252 (2016).
1070
- 1071 88. Verbanck M, Chen CY, Neale B, Do R. Detection of widespread horizontal pleiotropy in causal
1072 relationships inferred from Mendelian randomization between complex traits and diseases. *Nat*
1073 *Genet* **50**, 693-698 (2018).
1074
- 1075 89. Mangino M, *et al.* Genome-wide meta-analysis points to CTC1 and ZNF676 as genes
1076 regulating telomere homeostasis in humans. *Hum Mol Genet* **21**, 5385-5394 (2012).
1077
1078
1079
1080

1081 **Acknowledgements**

1082 This work utilized the computational resources of the NIH HPC Biowulf cluster. (<http://hpc.nih.gov>).

1083 Female Lung Cancer Consortium in Asia (NCI): This study was supported by a Grant-in-Aid for
1084 Scientific Research on Priority Areas from the Ministry of Education, Science, Sports, Culture and
1085 Technology of Japan, a Grant-in- Aid for the Third Term Comprehensive 10-Year Strategy for Cancer
1086 Control from the Ministry Health, Labor and Welfare of Japan, by Health and Labor Sciences
1087 Research Grants for Research on Applying Health Technology from the Ministry of Health, Labor and
1088 Welfare of Japan, by the National Cancer Center Research and Development Fund, the National
1089 Research Foundation of Korea (NRF) grant funded by the Korea government (MEST) (grant No.
1090 2011-0016106), a grant of the National Project for Personalized Genomic Medicine, Ministry for
1091 Health & Welfare, Republic of Korea (A111218-11-GM04), the Program for Changjiang Scholars and
1092 Innovative Research Team in University in China (IRT_14R40 to K.C.), the National Science &
1093 Technology Pillar Program (2011BAI09B00), MOE 111 Project (B13016), the National Natural
1094 Science Foundation of China (No. 30772531, and 81272618), Guangdong Provincial Key Laboratory
1095 of Lung Cancer Translational Medicine (No. 2012A061400006), Special Fund for Research in the
1096 Public Interest from the National Health and Family Planning Commission of PRC (No. 201402031),
1097 and the Ministry of Science and Technology, Taiwan (MOST 103-2325-B-400-023 & 104-2325-B-
1098 400-012). The Japan Lung Cancer Study (JLCS) was supported in part by the Practical Research for
1099 Innovative Cancer Control from Japan Agency for Medical Research and Development
1100 (15ck0106096h0002) and the Management Expenses Grants from the Government to the National
1101 Cancer Center (26-A-1) for Biobank. BioBank Japan was supported by the Ministry of Education,
1102 Culture, Sports, Sciences and Technology of the Japanese government. The Japan Public Health
1103 Center-based prospective Study (the JPHC Study) was supported by the National Cancer Center
1104 Research and Development Fund (23-A- 31[toku], 26-A-2, 29-A-4, and 2020-J-4) (since 2011) and a
1105 Grant-in-Aid for Cancer Research from the Ministry of Health, Labour and Welfare of Japan (from
1106 1989 to 2010). The Taiwan GELAC Study (Genetic Epidemiological Study for Lung
1107 AdenoCarcinoma) was supported by grants from the National Research Program on Genomic
1108 Medicine in Taiwan (DOH99-TD-G-111-028), the National Research Program for Biopharmaceuticals
1109 in Taiwan (MOHW 103-TDUPB-211-144003, MOST 103-2325-B-400-023) and the Bioinformatics
1110 Core Facility for Translational Medicine and Biotechnology Development (MOST 104-2319-B-400-
1111 002). This work was also supported by the Jinan Science Research Project Foundation (201102051),
1112 the National Key Scientific and Technological Project (2011ZX09307-001-04), the National Natural
1113 Science Foundation of China (No.81272293), the State Key Program of National Natural Science of
1114 China (81230067), the National Research Foundation of Korea (NRF) grant funded by the Korea
1115 government (MSIP) (No. NRF- 2014R1A2A2A05003665), Sookmyung Women’s University Research
1116 Grants, Korea (1-1603-2048), Agency for Science, Technology and Research (A*STAR), Singapore
1117 and the US National Institute of Health Grant (1U19CA148127-01). The overall GWAS project was
1118 supported by the intramural program of the US National Institutes of Health/National Cancer Institute.
1119 The following is a list of grants by study center: SKLCS (Y.T.K.)—National Research Foundation of
1120 Korea (NRF) grant funded by the Korea government (MEST) (2011-0016106). (J.C.) – This work was
1121 supported by a grant from the National R&D Program for Cancer Control, Ministry of Health
1122 & Welfare, Republic of Korea (grant no. 0720550-2). (J.S.S) – grant number is A010250. WLCS
1123 (T.W.)—National Key Basic Research and Development Program (2011CB503800). SLCS (B.Z.)—
1124 National Nature Science Foundation of China (81102194). Liaoning Provincial Department of
1125 Education (LS2010168). China Medical Board (00726). GDS (Y.L.W.)—Foundation of Guangdong

1126 Science and Technology Department (2006B60101010, 2007A032000002, 2011A030400010).
1127 Guangzhou Science and Information Technology Bureau (2011Y2-00014). Chinese Lung Cancer
1128 Research Foundation, National Natural Science Foundation of China (81101549). Natural Science
1129 Foundation of Guangdong Province (S2011010000792). TLCS (K.C., B.Q)—Program for Changjiang
1130 Scholars and Innovative Research Team in University (PCSIRT), China (IRT1076). Tianjin Cancer
1131 Institute and Hospital. National Foundation for Cancer Research (US). FLCS (J.C.W., D.R., L.J.)—
1132 Ministry of Health (201002007). Ministry of Science and Technology (2011BAI09B00). National
1133 S&T Major Special Project (2011ZX09102-010-01). China National High-Tech Research and
1134 Development Program (2012AA02A517, 2012AA02A518). National Science Foundation of China
1135 (30890034). National Basic Research Program (2012CB944600). Scientific and Technological Support
1136 Plans from Jiangsu Province (BE2010715). NLCS (H.S.)—China National High-Tech Research and
1137 Development Program Grant (2009AA022705). Priority Academic Program Development of Jiangsu
1138 Higher Education Institution. National Key Basic Research Program Grant (2011CB503805). GEL-S
1139 (A.S.)—National Medical Research Council Singapore grant (NMRC/0897/2004, NMRC/1075/2006).
1140 (J.Liu)—Agency for Science, Technology and Research (A*STAR) of Singapore. GELAC (C.A.H.)—
1141 National Research Program on Genomic Medicine in Taiwan (DOH98-TDG-111-015). National
1142 Research Program for Biopharmaceuticals in Taiwan (DOH 100- TD-PB-111-TM013). National
1143 Science Council, Taiwan (NSC 100- 2319-B-400-001). YLCS (Q.L.)—Supported by the intramural
1144 pro- gram of U.S. National Institutes of Health, National Cancer Institute. SWHS (W.Z., W.H.C.,
1145 N.R.)—The work was supported by a grant from the National Institutes of Health (R37 CA70867,
1146 UM1 CA182910) and the National Cancer Institute intramural research program, including NCI
1147 Intramural Research Program contract (N02 CP1101066). JLCS (K.M., T.K.)—Grants-in-Aid from the
1148 Ministry of Health, Labor, and Welfare for Research on Applying Health Technology and for the 3rd-
1149 term Comprehensive 10-year Strategy for Cancer Control; by the National Cancer Center Research and
1150 Development Fund; by Grant-in-Aid for Scientific Research on Priority Areas and on Innovative Area
1151 from the Ministry of Education, Science, Sports, Culture and — Technology of Japan. (W.P.)—NCI
1152 R01-CA121210. HKS (J.W.)— General Research Fund of Research Grant Council, Hong Kong
1153 (781511M). The Environment and Genetics in Lung Cancer Etiology (EAGLE), Prostate, Lung,
1154 Colon, Ovary Screening Trial (PLCO), and Alpha-Tocopherol, Beta-Carotene Cancer Prevention
1155 (ATBC) studies were supported by the Intramural Research Program of the National Institutes of
1156 Health, National Cancer Institute (NCI), Division of Cancer Epidemiology and Genetics. ATBC was
1157 also supported by U.S. Public Health Service contracts (N01-CN-45165, N01-RC-45035, and N01-
1158 RC-37004) from the NCI. PLCO was also supported by individual contracts from the NCI to the
1159 University of Colorado Denver (NO1-CN-25514), Georgetown University (NO1-CN-25522), the
1160 Pacific Health Research Institute (NO1-CN-25515), the Henry Ford Health System (NO1-CN-25512),
1161 the University of Minnesota, (NO1-CN- 25513), Washington University (NO1-CN-25516), the
1162 University of Pittsburgh (NO1-CN-25511), the University of Utah (NO1-CN- 25524), the Marshfield
1163 Clinic Research Foundation (NO1-CN- 25518), the University of Alabama at Birmingham (NO1-CN-
1164 75022), Westat, Inc. (NO1-CN-25476), and the University of California, Los Angeles (NO1-CN-
1165 25404). The Carotene and Retinol Efficacy Trial (CARET) is funded by the National Cancer Institute,
1166 National Institutes of Health through grants U01-CA063673, UM1-CA167462, and U01-CA167462.
1167 The Cancer Prevention Study-II (CPS-II) Nutrition Cohort was supported by the American Cancer
1168 Society. The NIH Genes, Environment and Health Initiative (GEI) partly funded DNA extraction and
1169 statis- tical analyses (HG-06-033-NCI-01 and RO1HL091172-01), genotyping at the Johns Hopkins
1170 University Center for Inherited Disease Research. This research was supported by the National
1171 Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No.

1172 2020R1A2C4002236). Genotyping of the samples in NJLCS was supported by the National Natural
1173 Science of China (81820108028).
1174
1175
1176 Female Lung Cancer Consortium in Asia (Tianjin): Tianjin Science and Technology Committee
1177 Foundation, 18YFZCSY00520.
1178
1179 Female Lung Cancer Consortium in Asia (Taiwan): The Ministry of Health and Welfare grants
1180 DOH97-TD-G-111-028 (ISC), DOH98-TD-G-111-017 (ISC), DOH99-TD-G-111-014 (ISC); DOH97-
1181 TD-G-111-026 (CAH), DOH98-TD-G-111-015 (CAH), DOH99-TD-G-111-028 (CAH); National
1182 Health Research Institutes grants NHRI-PH-110-GP-01, NHRI-PH-110-GP-03; and the Ministry of
1183 Science and Technology grants MOST108-2314-B-400-038(CAH), MOST109-2740-B-400-
1184 002(CAH).
1185
1186 The GWAS of lung cancer in European never smokers was supported by NIH R01 CA149462 (OYG).
1187
1188 OncoArray study in Europeans: The OncoArray data and analysis from INTEGRAL-ILCCO were
1189 supported by NIH U19 CA203654, and U19 CA148127. The data harmonization for ILCCO was
1190 supported by Canadian Institute for Health Research (CIHR) Canada Research Chair to R.J.H, and
1191 CIHR FDN 167273).
1192
1193 European never-smoking lung cancer study: CIA is a Research Scholar of the Cancer Prevention
1194 Institute of Texas (CPRIT) and supported by CPRIT grant RR170048.
1195
1196 Taiwan eQTL study: This study was supported by the Ministry of Health and Welfare grants DOH97-
1197 TD-G-111-028 (ISC), DOH98-TD-G-111-017 (ISC), DOH99-TD-G-111-014 (ISC); DOH97-TD-G-
1198 111-026 (CAH), DOH98-TD-G-111-015 (CAH), DOH99-TD-G-111-028 (CAH); National Health
1199 Research Institutes grants NHRI-PH-110-GP-01, NHRI-PH-110-GP-03; and the Ministry of Science
1200 and Technology grants MOST108-2314-B-400-038(CAH), MOST109-2740-B-400-002(CAH).
1201
1202 N.C. is supported by NIH grant 1R01HG010480. P.Y. is supported by Mayo Clinic Foundation
1203 Research Funds, NIH-CA77118 and CA80127. G.L. is supported is supported by the Alan Brown
1204 Chair and Lusi Wong Fund of the Princess Margaret Cancer Foundation. D.C.C is supported by
1205 U01CA209414. O.Y.G. is supported by NIH R01 CA231141.

1206
1207 **Author contributions**

1208 Organized and designed the study: Q.L., J.S., N.R., J.C., S.J.C., N. Chatterjee, K.S., K.M., T.K.,
1209 M.T.L., R.J.H., C.I.A., O.Y.G., H.S., I-S.C., C.A.H. Conducted and supervised new genotyping for the
1210 project: K.S., K.M., T.K. Contributed to the design and execution of statistical analyses: J.S., J.C.,
1211 S.J.C., N.R., Q.L., L.S., B.D.R., S.L, R.J.H., C.I.A., O.Y.G., N.C., I-S.H., K.M.F., K.S., K.M., T.K.
1212 Wrote the first draft: J.S., J.C., N.R., Q.L., K.S., K.M., T-Y.C., J.D., R.J.H., K.C., N.C., O.Y.G.,
1213 C.A.H., S.J.C., C.I.A., H.S., T.K. Conducted epidemiology studies and contributed samples to GWAS
1214 and/or conducted initial genotyping. Q.L., M.T.L., B.A.B., W.H., N.E.C., BT.J., M.Song, H.P., D.A.,
1215 C.C.C., L.B., M.Y., A.H., B.H., J.Liu, B.Zhu, S.I.B., C.K., K.Wyatt, S.A.L., A.Chao, J.F.F.J., S.J.C.,

1216 N.R., Z.Wang, C.L., J.C., C.W., W.T., D.Lin, S.J.A., X.C.Z., J.S., Y.L.W., M.P.W., L.P.C., J.C.M.H.,
1217 V.H.F.L., Z.H., K.M., J.Y.P., Jia.Liu, H.S.J., J.E.C., Y.Y.C., H.N.K., M.H.S., S.S.K., Y.C.K., I.J.O.,
1218 S.W.S., H.I.Y., Y.T.K., Y.C.H., J.H.K., Y.H.K., J.S.S., Y.J.J., K.H.P., C.H.K., J.S.K., I.K.P., B.S.,
1219 Jie.L, Z.W., S.C., J.Y., J.C.W., Y.Y., Y.B.X., Y.T.G., D.L., J.Y.Y.W., H.C., L.J., J.Z., G.J., K.F., Z.Y.,
1220 B.Z., W.W., P.G., Q.H., X.L., Y.R., A.S., Y.L., Y.C., W.Y.L., W.Z., X.O.S., Q.C., G.Y., B.Q., T.W.,
1221 H.G., L.L., P.X., F.W., G.W., J.X., J.L., R.C.H.V., B.B., H.D.III.H., J.Wang, A.D.L.S., J.K.C.C.,
1222 V.L.S., K.C., H.Z., H.D., C.A.H., T.Y.C., L.H.C., IS.C., C.Y.C., S.S.J., CH.C., G.C.C., C.F.H., Y.H.T.,
1223 W.C.W., K.Y.C., M.S.H., W.C.S., Y.M.C., C.L.W., K.C.C., C.J.Y., H.H.H., F.Y.T., H.C.L., C.J.C., P.C.Y.,
1224 K.Shiraishi, T.K., H.K., S.M., H.H., K.Goto, Y.Ohe, S.W., Y.Yatabe, M.T., R.Hamamoto,
1225 A.Takahashi, Y.Momozawa, M.Kubo, Y.K., Y.D., Y.Miyagi, H.N., T.Y., N.S., M.I., M.H., Y.N.,
1226 K.Takeuchi, K.W., K.Matsuda, Y.Murakami, K.S., K.T., Y.O., M.S., H.Suzuki, A.G., Y.M., T.H.,
1227 M.K., K.O., H.S., J.D., H.M., M.Z., R.J.H., S.L., A.T., C.C., S.E.B., M.Johansson, A.R., H.Bö.,
1228 H.E.W., D.C., G.R., S.A., P.B., J.M.K., J.K.F., S.S.S., L.L.M., O.M., H.Bö., G.L., A.A., L.A.K., S.ZN.,
1229 K.G., M.J., A.C., J.M.Y., P.L., M.B.S., M.C.A., C.I.A., A.G.S., R.H., M.R.S., O.Y.G., I.P.G., X.W.,
1230 P.Y. All authors contributed to the writing and final review of the manuscript.
1231

1232 **Competing interests**

1233 The authors declare no competing interests.

1234

1235

1236 Table 1. Demographic characteristics of the subjects in the discovery and the replication datasets for a GWAS of
 1237 lung adenocarcinoma in East Asians

1238

	Discovery ^a		Replication ^b		Combined	
	Cases	Controls	Cases	Controls	Cases	Controls
Male	4,021 (34%)	11,609 (38%)	5,650 (57%)	62,596 (52%)	9,671(45%)	74,205 (49%)
Female	7,732 (66%)	18,953 (62%)	4,255 (43%)	57,518 (48%)	11,987 (55%)	76,471 (51%)
Individuals with smoking history	3,751 (32%)	9,780 (32%)	6,108 (62%)	58,430 (49%)	9,859 (46%)	68,210 (45%)
Never-smokers	8,002 (68%)	20,782 (68%)	3,797 (38%)	61,684 (51%)	11,799 (54%)	82,466 (55%)
Total	11,753	30,562	9,905	120,114	21,658	150,676

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

^a The discovery dataset includes 4,438 cases and 4,544 controls from the FLCCA study, 1,923 cases and 3,544 controls from the NJLCS study, 3,921 cases and 19,910 controls from the NCC study and 1,471 cases and 2,564 controls from the ACC study. ^b The replication dataset consists of new candidate variant genotyping conducted in Japanese study LUAD subjects by the NCC study center and controls from the BioBank Japan. More details can be found in Supplementary Table 1 and Methods.

Table 2: Novel genetic variants associated with lung adenocarcinoma in East Asians. All p-values are nominal and two-sided.

Chr	BP	SNP	Genes	Discovery				Replication		Combined	
				Eff/Ref	EAF	OR (95% CI)	P	OR (95% CI)	P	OR (95% CI)	P
3	138570011	rs137884934	<i>PIK3CB</i>	T/C	0.09	0.81(0.74,0.89)	6.33×10 ⁻⁶	0.80(0.76,0.85)	1.88×10 ⁻¹⁵	0.80(0.77,0.84)	6.21×10 ⁻²⁰
2	25757709	rs682888	<i>DTNB</i>	C/T	0.47	0.89(0.86,0.93)	4.94×10 ⁻¹⁰	0.91(0.88,0.94)	1.57×10 ⁻¹⁰	0.90(0.88,0.92)	5.96×10 ⁻¹⁹
11	61581656	rs174559	<i>FADS1</i>	A/G	0.39	0.91(0.88,0.94)	6.10×10 ⁻⁷	0.91(0.89,0.94)	6.22×10 ⁻⁹	0.91(0.89,0.93)	1.93×10 ⁻¹⁴
15	49757466	rs71467682 ^a	<i>FGF7, SECISBP2L</i>	G/A	0.31	0.91(0.87,0.95)	2.46×10 ⁻⁶	0.90(0.88,0.93)	2.30×10 ⁻⁹	0.91(0.88,0.93)	2.81×10 ⁻¹⁴
10	126324209	rs10901793	<i>FAM53B, METTL10</i>	A/G	0.30	1.10(1.06,1.14)	3.14×10 ⁻⁷	1.07(1.04,1.10)	1.03×10 ⁻⁵	1.08(1.06,1.11)	3.04×10 ⁻¹¹
7	124373384	rs4268071 ^b	<i>GPR37</i>	T/G	0.04	1.39(1.25,1.54)	7.27×10 ⁻¹⁰	NA	NA	1.39(1.25,1.54)	7.27×10 ⁻¹⁰
6	53389995	rs531557	<i>GCLC</i>	T/A	0.60	0.90(0.87,0.94)	7.73×10 ⁻⁷	0.94(0.91,0.97)	8.49×10 ⁻⁵	0.93(0.90,0.95)	9.25×10 ⁻¹⁰
19	725066	rs116863980	<i>PALM</i>	A/G	0.06	1.31(1.16,1.47)	7.94×10 ⁻⁶	1.17(1.09,1.26)	2.50×10 ⁻⁵	1.21(1.14,1.29)	2.63×10 ⁻⁹
15	56454223	rs764014	<i>RFX7</i>	G/A	0.47	0.91(0.88,0.95)	5.75×10 ⁻⁷	0.95(0.92,0.98)	7.36×10 ⁻⁴	0.94(0.91,0.96)	7.73×10 ⁻⁹
4	44174404	rs117715768	<i>KCTD8</i>	T/C	0.06	1.24(1.14,1.34)	4.48×10 ⁻⁷	1.10(1.04,1.17)	1.28×10 ⁻³	1.15(1.09,1.21)	2.45×10 ⁻⁸
4	157894892	rs1373058	<i>PDGFC</i>	A/T	0.57	1.10(1.05,1.15)	8.55×10 ⁻⁶	1.06(1.03,1.09)	3.60×10 ⁻⁴	1.07(1.05,1.10)	3.86×10 ⁻⁸

^a: rs71467682 is in weak LD with rs77468143 ($R^2 = 0.27$ in EA) that was previously reported to be associated with LUAD in EUR populations¹⁶.

^b: Replication data not available.

Table 3: Conditional and joint analyses identified independently associated risk SNPs for lung adenocarcinoma at two existing loci in East Asians. All p-values are nominal and two-sided.

Chr	BP	SNP	Gene	GWAS analysis ^a				Conditional analysis ^b		Joint analysis ^c	
				Eff/Ref	EAF	OR (95% CI)	P	OR (95% CI)	P	OR (95% CI)	P
5	1280477	rs13167280	<i>TERT</i>	A/G	0.22	1.47(1.37,1.57)	6.99×10 ⁻³⁰	1.33(1.24,1.42)	8.36×10 ⁻¹⁷	1.29(1.20,1.38)	4.07×10 ⁻¹³
5	1286516	rs2736100		A/G	0.56	0.75(0.72,0.77)	7.92×10 ⁻⁵⁸			0.80(0.77,0.83)	9.83×10 ⁻³²
5	1290319	rs62332591		G/T	0.52	0.79(0.75,0.83)	3.53×10 ⁻²³	0.87(0.83,0.91)	2.95×10 ⁻⁹	0.87(0.83,0.92)	3.21×10 ⁻⁸
6	41483390	rs9367106	<i>FOXP4</i>	C/G	0.32	1.20(1.15,1.26)	1.06×10 ⁻¹⁴			1.19(1.14,1.25)	2.39×10 ⁻¹³
6	41483960	rs12664490		T/C	0.16	0.80(0.75,0.85)	5.52×10 ⁻¹²	0.81(0.76,0.86)	1.34×10 ⁻¹⁰	0.81(0.76,0.86)	1.24×10 ⁻¹⁰

^a: Data from single-variant analysis in GWAS.

^b: Conditional analysis using GCTA, conditioning on the lead variant in each locus.

^c: Joint analysis using GCTA including the lead variant and the significant variants in conditional analysis.

Legend

Fig. 1. Manhattan plot for GWAS meta-analysis of lung adenocarcinoma in East Asians. The x-axis represents chromosomal location, and the y-axis represents $-\log_{10}(\text{p-value})$. All p-values were two-sided and not adjusted for multiple testing. The red horizontal line denotes the p-value threshold for declaring genome-wide significance at 5×10^{-8} . For each box, red text represents a novel variant (12 novel variants, including the lead variants from 10 novel loci, rs12664490 by conditional analysis at 6p21.1, a locus previously reported in East Asians, and rs71467682 at 15q21.2, a locus previously reported in Europeans); black text represents a previously reported association (16 variants in total, including three independently associated variants in 5p15.33 locus). For each locus, a green circle represents the top p-value from the discovery samples, a red diamond represents the p-value combining the discovery and the replication data, a black square represents the p-value combining our discovery data and Chinese samples in Dai *et al.*²⁴ (for three variants identified in a cross-ancestry analysis of East Asians and Europeans in Dai *et al.*²⁴, Supplementary Table 3). In summary, 28 variants at 25 loci achieved genome-wide significance, including 16 previously reported variants and 12 novel variants.

Fig. 2. Colocalization of lung adenocarcinoma GWAS signal from the new locus on Chr11 with *FADS1* eQTL signal. Colocalization analysis was performed using HyPrColoc with summary statistics from Taiwanese lung eQTL data (for *FADS1* gene, Panel A) and those of EA GWAS discovery set (Panel B). LD R^2 (1000 Genomes, EA) of each SNP with the GWAS lead SNP, rs174559 (red circle), is color-coded as shown in the top band. Colocalization posterior probability (PP) is shown next to the candidate SNP, rs174559. Note that the p-value of rs174559 in GWAS was based on the discovery data and did not include the Japanese replication data. All eQTL p-values were two-sided and not adjusted for multiple testing.

Fig. 3. Comparing odds ratios (ORs) of lung adenocarcinoma susceptibility variants between East Asian (EA) and European (EUR) populations. Here, the effect allele was defined as the minor allele in EA. Each error bar represents the 95% confidence interval of the OR (the center). A: Susceptibility variants previously discovered (at genome-wide significance) in both EA and EUR populations. B: Variants previously identified by multiple-ancestry meta-analysis of Chinese and EUR populations; C: Variants were identified by multiple-ancestry meta-analysis combining EA samples in our study and EUR samples in ILLCO. D: Variants identified only in EA populations. E: Novel variants identified in the current study; F: Variants identified only in EUR populations. Variants are labeled with *, **, *** and **** corresponding to $0.01 \leq p_{\text{het}} < 0.05$, $0.001 \leq p_{\text{het}} < 0.01$, $0.0001 \leq p_{\text{het}} < 0.001$ and $p_{\text{het}} < 0.0001$, respectively; here, p_{het} (t-statistic, two-sided) is the p-value for testing the heterogeneity of effect sizes between EA and EUR populations. Sample sizes for EUR populations in all panels: 11,273 cases and 55,483 controls. Sample sizes for EA populations: 11,753 cases and 30,562 controls for panels A, B, C, D, and F; 21,658 cases and 150,676 controls for panel E.

Fig. 4. A polygenic risk score (PRS) is more strongly associated with risk of lung adenocarcinoma in never-smokers than in individuals with a history of smoking ($P = 0.0058$). The PRS was defined based on 25 independent variants that achieved genome-wide significance in EA with weights derived from the meta-analysis of the current study (Supplementary Table 4). The odds ratios (ORs) and the standard errors of the 12 novel variants were based on 21,658 cases and 150,676 controls. The ORs and the standard errors of the other 13 variants were based on 11,753 cases and 30,562 controls. The figure shows the ORs and their 95% confidence intervals comparing each quintile group to the middle quintile for individuals with a history of smoking (blue) and never-smokers (red).

Fig. 5. The expected area under the receiver operating characteristic curve (AUC) of a polygenic risk score (PRS) built based on a GWAS of specified sample sizes for lung adenocarcinoma in never-smoking East Asians. For “1 million controls”, the x-coordinate represents the number of cases, assuming the study has 1 million controls. For “Equal number of cases and controls”, the x-coordinate represents the numbers of cases, assuming the same number of cases and controls.