

1 **Development and evaluation of a polygenic risk score for lung cancer in never-smoking**
2 **women: a large-scale prospective Chinese cohort study**

3
4 Xiaoxia Wei^{1†}, Dianjianyi Sun^{2,3†}, Jiabin Gao^{1†}, Jing Zhang^{1†}, Meng Zhu^{1,4}, Canqing Yu^{2,3},
5 Zhimin Ma¹, Yating Fu¹, Chen Ji¹, Pei Pei³, Ling Yang^{5,6}, Iona Y. Millwood^{5,6}, Robin
6 G.Walters^{5,6}, Yiping Chen^{5,6}, Huaidong Du^{5,6}, Guangfu Jin^{1,4}, Zhengming Chen⁶, Zhibin Hu^{1,7},
7 Liming Li^{2,3}, Hongbing Shen^{1,4,8}, Jun Lv^{2,3} ‡, Hongxia Ma^{1,4,8‡}

8
9 1 Department of Epidemiology and Biostatistics, International Joint Research Center on
10 Environment and Human Health, Center for Global Health, School of Public Health,
11 Nanjing Medical University, Nanjing, China

12 2 Department of Epidemiology and Biostatistics, School of Public Health, Peking
13 University Health Science Center, Beijing, China

14 3 Peking University Center for Public Health and Epidemic Preparedness & Response,
15 Beijing, China

16 4 Jiangsu Key Lab of Cancer Biomarkers, Prevention and Treatment, Collaborative
17 Innovation Center for Cancer Medicine, Nanjing Medical University, Nanjing, China

18 5 Medical Research Council Population Health Research Unit at the University of Oxford,
19 Oxford, United Kingdom

20 6 Clinical Trial Service Unit & Epidemiological Studies Unit (CTSU), Nuffield Department
21 of Population Health, University of Oxford, United Kingdom

22 7 State Key Laboratory of Reproductive Medicine, Center for Global Health, Nanjing
23 Medical University, Nanjing, China

24 8 Research Units of Cohort Study on Cardiovascular Diseases and Cancers, Chinese
25 Academy of Medical Sciences, Beijing 100730, China.

26
27 † Joint first authors, contributed equally.

28 ‡ Joint last authors, contributed equally.

29
30 **Correspondence to:**

31 Hongxia Ma, Department of Epidemiology, Center for Global Health, School of Public
32 Health, Nanjing Medical University, 101 Longmian Rd., Nanjing 211166, China. Email:
33 hongxiama@njmu.edu.cn; or Jun Lv, Department of Epidemiology and Biostatistics, Peking

34 University Health Science Center, 38 Xueyuan Road, Beijing 100191, China, Phone: 010-
35 82801620, Email: lvjun@bjmu.edu.cn.

36

37 **Abbreviations**

38 *AUC*: area under the curve

39 *BMI*: body mass index

40 *CI*s: confidence intervals

41 *CKB*: China Kadoorie Biobank

42 *COPD*: chronic obstructive pulmonary disease

43 *C+T*: clumping and thresholding

44 *DHS*: DNase I hypersensitivity site

45 *FLCCA*: Female Lung Cancer Consortium in Asia

46 *GSA*: Global Screening Array

47 *GWAS*s: Genome-wide association studies

48 *HR*s: hazard ratios

49 *LUAD*: lung adenocarcinoma

50 *LUSC*: lung squamous cell carcinoma

51 *NRI*: net reclassification improvement

52 *NSCLC*: non-small cell lung cancer

53 *OR*s: odds ratios

54 *PCA*: principal component analysis

55 *PRS*: polygenic risk score

56 *Q-Q*: quantile-quantile

57 *SE*s: standard errors

58

59 **Novelty and Impact Statement**

60 We newly developed PRS for lung cancer risk in never-smoking women, and further validated

61 the performance of PRS in an independent cohort study. The GWAS-derived PRS-21 may
62 provide complementary information to current lung cancer screening guidelines to better
63 identify high-risk groups, especially in never-smoking women.

64

65 **Abstract**

66 The proportion of lung cancer in never smokers is rising, especially among Asian women, but
67 there is no effective early detection tool. Here, we developed a polygenic risk score (PRS),
68 which may help to identify the population with higher risk of lung cancer in never-smoking
69 women. We first performed a large GWAS meta-analysis (8595 cases and 8275 controls) to
70 systematically identify the susceptibility loci for lung cancer in never-smoking Asian women,
71 and then generated a PRS using GWAS datasets. Furthermore, we evaluated the utility and
72 effectiveness of PRS in an independent Chinese prospective cohort comprising 55,266
73 individuals. The GWAS meta-analysis identified 8 known loci and 1 new locus (5q11.2) at the
74 genome-wide statistical significance level of $P < 5 \times 10^{-8}$. Based on the summary statistics of
75 GWAS, we derived a polygenic risk score including 21 variants (PRS-21) for lung cancer in
76 never-smoking women. Furthermore, PRS-21 had a hazard ratio per SD of 1.29 (95% CI =
77 1.18-1.41) in the prospective cohort. Compared with participants who had a low genetic risk,
78 those with an intermediate (HR=1.32, 95% CI: 1.00-1.72) and high (HR=2.09, 95% CI: 1.56-
79 2.80) genetic risk had a significantly higher risk of incident lung cancer. The addition of PRS-
80 21 to the conventional risk model yielded a modest significant improvement in AUC (0.697 to
81 0.711) and net reclassification improvement (24.2%). The GWAS-derived PRS-21
82 significantly improves the risk stratification and prediction accuracy for incident lung cancer
83 in never-smoking Asian women, demonstrating the potential for identification of high-risk
84 individuals and early screening.

85 **Keywords:** lung cancer in never-smoking women, polygenic risk score, predictive markers

86 **Introduction**

87 Lung cancer is the second most common cancer worldwide, with 2.20 million newly
88 diagnosed cases and 1.80 million deaths estimated in 2020 (1). Although tobacco smoking is
89 known as the leading risk factor for lung cancer, the proportion of never-smokers with lung
90 cancer has been increasing over time, especially in Asian countries (2-4). Epidemiologic
91 studies indicate that lung cancer in never-smokers is mostly female, adenocarcinoma
92 histologic subtype and East Asian ethnicity (5). Additionally, some risk factors have been

93 investigated to explain the occurrence of lung cancer in never-smokers, such as passive
94 smoking, occupational exposure, indoor and outdoor pollution as well as genetic
95 susceptibility (6, 7). Due to the distinct etiological factors between lung cancer in never-
96 smokers and smokers, the challenges associated with prevention are heightened, thereby
97 emphasizing only screening never-smokers identified as at high enough risk to warrant
98 screening.

99 Polygenic risk score (PRS) is a useful tool to summarize an individual's genetic risk for a
100 complex disease, obtained by aggregating and quantifying the effect of common variants
101 associated with the disease (8, 9). In recent years, PRS has been widely used in cancer risk
102 prediction and is capable of identifying a larger fraction of the population with risk equivalent
103 to rare monogenic mutations (10, 11). For instance, PRS has been shown to effectively stratify
104 breast cancer risk in women and even refine risk estimates for *BRCA1* and *BRCA2* mutation
105 carriers (12). Genome-wide association studies (GWASs) have identified a number of genetic
106 variants associated with the risk of lung cancer in never-smoking women (13-16), which offer
107 an opportunity to generate a PRS to identify the subgroups of the population at the highest
108 risk for this disease. A previous study derived a PRS based on 10 significant SNPs from lung
109 cancer risk in never-smoking women in the Asian population (17); however, this PRS-10 was
110 derived from a limited number of SNPs, and its efficiency and significance were not validated
111 by prospective studies.

112 In this study, we performed a large genome-wide association meta-analysis with a total
113 of 8595 cases and 8275 controls to systematically identify genetic factors for the risk of lung
114 cancer in never-smoking Asian women, and then constructed and compared the risk
115 assessment effect of a set of PRSs. Furthermore, we selected PRS-21 to examine its prediction
116 value in a large prospective cohort of the Chinese population, a subset of women from the
117 China Kadoorie Biobank (CKB) and also evaluated the performance of the conventional risk
118 model when integrating PRS.

119 **Methods**

120 ***Study design and participants***

121 This study was performed in two stages, with a meta-analysis and a prospective cohort
122 study design (**Supplementary Figure 1**). In the first phase, we employed a case-control study
123 design and performed a meta-analysis of data from four independent GWAS datasets,
124 including a total of 8,595 cases and 8,275 controls from NJMU GWAS project (18) in never-
125 smoking women in Asia (513 cases and 962 controls), NJMU Global Screening Array (GSA)

126 Project(19) (3,657 cases and 3,271 controls), NJMU OncoArray project (19) (296 cases and
127 301 controls), and Female Lung Cancer Consortium in Asia (FLCCA) GWAS (14) (4,129
128 cases and 3,741 controls). All participants were restricted to never-smoking female lung
129 cancer cases and never-smoking female controls, and comprehensive details for individual
130 studies have been previously reported (14, 19). Using the estimated effect sizes of the
131 summary statistics from this GWAS meta-analysis, we derived a set of PRSs for lung cancer
132 in never-smoking women.

133 In the second stage of the study, we further evaluated the prediction performance of
134 PRSs in an independent cohort study of CKB. Details of the design and methods of the CKB
135 have been described previously (20). Briefly, 512,715 adults aged 30-79 years were recruited
136 during 2004-2008 across ten geographically diverse (five urban and five rural) regions in
137 China. At baseline, all participants completed a laptop-based questionnaire on demographic
138 characteristics, lifestyle factors, biomedical information, and provided physical measurements
139 and a blood sample. Among those, 55,266 never-smoking female individuals with both
140 genotypic and phenotypic data remained in the present study. The outcome of our study was
141 lung cancer (ICD-10: C34). Participants were assessed for lung cancer risk between baseline
142 and date of diagnosis, death, loss to follow-up, or study end date (December 31, 2016),
143 whichever came first.

144 ***Genotyping and quality control***

145 Multiple genotyping arrays for genome-wide genotyping were used in the four GWAS
146 and CKB study samples (**Supplementary Table 1**). Detailed descriptions of the quality
147 control procedures on genotyping data for each GWAS have been described previously (14,
148 19). Briefly, we performed quality control at both sample-level and variant-level. Qualified
149 genotypes were phased with SHAPEIT v2 (21, 22) and imputation was performed with
150 IMPUTE v2 (23) by using the reference of the 1000 Genomes Project Phase III database
151 (October 2014 release).

152 ***Polygenic risk score***

153 Based on the summary statistics from the meta-analysis of GWAS, we derived the PRS
154 for lung cancer in never-smoking women, including two components: (i) derived PRSs using
155 a clumping and thresholding (C+T) approach, with a range of p-value (5×10^{-8} , 5×10^{-7} , 5×10^{-6} ,
156 5×10^{-5} , and 5×10^{-4}) and linkage disequilibrium-driven clumping procedure in PLINK version
157 1.90b (--clump), retaining the SNP with the smallest p-value excluding variants with $r^2 > 0.1$
158 in a 1000-kb window (**Supplementary Table 2-6**). The PRS with the highest area under the
159 curve (AUC) and the lowest *P* value in the cohort study was selected as the best-performing

160 PRS. The best-performing PRS was obtained at a P -value threshold of 5×10^{-7} , including 16
161 variants (PRS-16). (ii) additional SNPs that were previously reported in lung cancer GWAS
162 and replicated in our current study at P value < 0.00058 (after multiple comparisons at
163 $0.05/86$, 86 were the sum of 16 SNPs included in best-performing PRS and 71 reported lung
164 cancer susceptibility SNPs(19), including one duplicate SNP). When correlation exists (r^2
165 > 0.1), 5 reported lung cancer susceptibility SNPs (rs2293607, rs2293607, rs4236709,
166 rs10429489, rs1200399, and rs77468143) representing independent loci with the strongest
167 statistical significance were retained. The final PRS (PRS-21) was constructed by combining
168 above both components. The list of known lung cancer-related loci and the SNPs included in
169 PRS-21 were shown in **Supplementary Table 7** and **Supplementary Table 8**, respectively.

170 Polygenic risk scores were generated by multiplying the dosage of the effect allele for
171 each SNP by its respective weight (log OR). In addition, we compared the predictive
172 performance of our PRS with two known PRSs from Asian populations: PRS-10 for lung
173 cancer in never-smoking women (17) and PRS-19 for non-small cell lung cancer (NSCLC)
174 (19).

175 ***Statistical analysis***

176 For each GWAS dataset, a single-variant association analysis was performed using the
177 SNPTEST software (V.2.54) based on a log-additive model (24). Per-allele odds ratios (ORs)
178 and standard errors (SEs) from each GWAS dataset were combined by a fixed-effects meta-
179 analysis using the METAL software (25). The genome-wide significance threshold was set at
180 $P < 5.0 \times 10^{-8}$. For the CKB cohort study, we performed Cox proportional hazards regression
181 models to assess the association between PRS and risk of lung cancer in never-smoking
182 women to estimate hazard ratios (HRs) and 95% confidence intervals (CIs) with adjustment
183 for age, education, body mass index (BMI), chronic obstructive pulmonary disease (COPD),
184 personal history of cancer, family history of cancer, and the first ten principal components of
185 ancestry determined through principal component analysis. The assumption of proportional
186 hazards was assessed by testing the significance of Schoenfeld residuals.

187 A potential non-linear association between the PRS and the risk of lung cancer in never-
188 smoking women was assessed by fitting a restricted cubic spline (26). The participants were
189 categorized into low (quintile 1), intermediate (quintiles 2–4), and high (quintile 5) genetic
190 risk according to quintiles of the PRS, and those with the lowest PRS as the reference group.
191 Effect modification of the PRS by the epidemiologic risk factors was evaluated by fitting
192 additional interaction terms in the model. We also estimated the 5-year and lifetime absolute

193 risks (up to 80 years of age) of developing lung cancer in never-smoking women based on
194 Cox proportional hazards model.

195 In addition, we evaluated the incremental predictive ability of PRS for improving the
196 conventional risk model on incident lung cancer in never-smoking women. At first, we
197 selected risk factors based on prior evidence to construct a conventional risk model (27-29),
198 which included age, education, BMI, COPD history, personal history of cancer, and family
199 history of cancer. Then, AUC and net reclassification improvement (NRI) were used to
200 estimate the discriminatory accuracy when PRS was added to the conventional risk model.

201 To evaluate the robustness of our results, we carried out several further sensitivity
202 analyses: (1) excluding incident cases occurring during the first year of follow-up; (2)
203 reclassifying genetic risk levels based on quartiles (bottom, 2-3, and top quartiles defined as
204 low, intermediate, and high genetic risk, respectively) or tertiles (corresponding to low,
205 intermediate and high genetic risk) of PRSs; (3) using a 10-fold cross-validation approach to
206 avoid over-fitting of the risk model.

207 All *p*-values were two-sided and $p < 0.05$ was considered statistically significant. All
208 statistical analyses were performed with R software, version 3.5.1.

209 **Results**

210 **Characteristics of the subjects**

211 The GWAS meta-analysis included 8,595 never-smoking female lung cancer cases and
212 8,275 never-smoking controls from four independent GWAS datasets. The characteristics of
213 the participants are presented in **Supplementary Table 9**, in which 86.46% of cases were
214 classified as lung adenocarcinoma (LUAD) and 11.34% as lung squamous cell carcinoma
215 (LUSC). For the PRS application, a total of 55,266 never-smoking females with complete
216 genotypes and phenotypes from the CKB cohort were included in the analysis. During a
217 median follow-up of 10.53 years (IQR 9.58-11.38), 466 incident lung cancer cases occurred
218 (**Supplementary Table 10**).

219 **GWAS meta-analysis**

220 A quantile-quantile (Q-Q) showed minimal evidence of systematic inflation from
221 population stratification for the meta-analysis ($\lambda = 1.02$, **Supplementary Figure 2**). The
222 principal component analysis (PCA) revealed good clustering of all case and control samples
223 within the same cluster (**Supplementary Figure 3**). The GWAS meta-analysis identified 9
224 risk loci for lung cancer in never-smoking women at the genome-wide statistical significance
225 threshold ($P < 5 \times 10^{-8}$) (**Table 1 and Supplementary Figure 4**). Of those, we confirmed the

226 association of 8 previously reported loci (3q28, 5p15.33, 6p21.32, 6p21.32, 6p21.1, 6q22.1,
227 10q25.2, and 11q23.3) and identified one novel risk loci at 5q11.2 (lead variant: rs1498606, P
228 $= 4.79 \times 10^{-8}$). The lead SNP rs1498606 is located in the intronic region of the *PDE4D* gene
229 on chromosome 5q11.2 (**Figure 1 A**). **Figure 1 B** depicted rs1498606 associated with the risk
230 of lung cancer in never-smoking women in each GWAS study.

231 To annotate the underlying biological significance of novel loci in 5q11.2, we evaluated
232 the potential functions of the lead variant rs1498606 and their correlated variants ($r^2 \geq 0.60$)
233 and found that rs6450500 (with rs1498606 $r^2 = 1.0$) is more likely to be a functional SNP that
234 overlapped with DNase I hypersensitivity site (DHS) or promoter/enhancer histone marks (eg,
235 H3K4me1, H3K4me3, and H3K27ac) in lung tissues and transcription factor binding sites
236 (**Supplementary table 11 and supplementary figure 5**).

237 **PRS construction and evaluation**

238 We generated a set of PRSs using effect size derived from GWAS meta-analysis at
239 different significance thresholds. Then, we evaluated the PRSs in the CKB cohort and found
240 that all five GWAS-derived PRSs were significantly associated with the risk of lung cancer in
241 never-smoking women (**Table 2**). The best-performing PRS was obtained at a P value
242 threshold of 5×10^{-7} based on the C+T approach, resulting in a 16-SNP PRS (PRS-16), which
243 showed the strongest association ($HR_{per\ 1\ SD} = 1.28$, 95% CI: 1.17-1.40; AUC=0.573). When the
244 16-SNP PRS was supplemented with 71 reported lung cancer susceptibility SNPs, a new PRS
245 with 21 SNPs (PRS-21) was further derived and showed improved prediction accuracy (HR_{per}
246 $1\ SD} = 1.29$, 95% CI: 1.18-1.41; AUC=0.576), which was ultimately used to define genetic risk
247 for all analyses. Moreover, the PRS-21 had better performance prediction than the two
248 published PRSs.

249 The polygenic risk score derived from 21 SNPs could significantly predict risk of lung
250 cancer among never-smoking women in the CKB cohort. Incident lung cancer cases in never-
251 smoking women had a higher PRS-21 than those without incident lung cancer (**Figure 2 A**). A
252 linear and positive association between the PRS-21 and the risk of incident lung cancer in
253 never-smoking women using restricted cubic spline regression was found ($P_{non-linear} = 0.8762$;
254 **Figure 2 B**). We also observed a dose-response relationship between PRS quintiles and the
255 risk of lung cancer in never-smoking women ($P_{trend} = 5.90 \times 10^{-8}$) (**Figure 2 C and**
256 **Supplementary Table 12**). Additionally, we assessed the effect of the PRS by categorizing
257 the PRS distribution into low (bottom quintile), intermediate (quintiles 2-4) and high (top
258 quintile) genetic risk groups. Compared with participants who had a low genetic risk, those

259 with an intermediate (HR=1.32, 95% CI: 1.00-1.72) and high (HR=2.09, 95% CI: 1.56-2.80)
260 genetic risk had a significantly greater risk of incident lung cancer (**Figure 2 D and**
261 **Supplementary Table 13**). These associations did not materially change when excluding
262 incident cases that occurred in the first year of follow-up (**Supplementary table 14**).
263 Similarly, reclassification of genetic risk by dividing PRS into quartiles or tertiles also showed
264 a positive association with the risk of lung cancer in never-smoking women (**Supplementary**
265 **table 15**). Positive associations were also observed in stratified analyses including age,
266 education, BMI, COPD history, personal history of cancer, and family history of cancer, and
267 there was no evidence of interaction (**Supplementary table 16**).

268 **Estimated 5-year and lifetime absolute risks of lung cancer in never-smoking women**

269 Estimated 5-year and lifetime absolute risks for different PRS categories among never-
270 smoking women in the CKB cohort are shown in **Figure 3**. The 5-year absolute risk for
271 never-smoking women in the highest 1% of the PRS reached 2.48% (**Figure 3A**).
272 Furthermore, we assessed how the PRS affected lifetime trajectories of lung cancer in never-
273 smoking women at risk. The cumulative risks by age of 80 for lung cancer events were 2.26%
274 for individuals with low genetic risk, 2.95% for individuals with intermediate genetic risk,
275 and 4.71% among those with high genetic risk, respectively (**Figure 3B**).

276 **Incremental value of PRS in risk prediction**

277 We further evaluated the incremental predictive ability of PRS-21 for lung cancer in
278 never-smoking women. **Supplementary Table 17** lists the model parameters in the CKB
279 cohort. The addition of the PRS to this baseline model containing conventional risk factors
280 showed a statistically significant improvement in discrimination, with the AUC improved
281 from 0.697 to 0.711 (difference, 1.4%; $P<0.001$) (**Supplementary Figure 6**). Similar results
282 were observed in analyses using the 10-fold cross-validation approach (**Supplementary**
283 **Table 18**). Furthermore, there was a significant reclassification improvement, with an overall
284 improvement in the net risk stratification of 24.2% (13.7%-34.8%) (**Supplementary Table**
285 **19**).

286 **Discussion**

287 In this study, leveraging large GWAS data sets of women with Asian ancestry, we identified a
288 novel risk locus and newly developed PRS for lung cancer risk in never-smoking women, and
289 further validated the performance of PRS in an independent prospective cohort study. Our
290 results suggested that GWAS-derived PRS-21 is associated with the risk of lung cancer in
291 never-smoking women with a dose-response relationship and is powerful to identify high-risk

292 individuals. Furthermore, adding PRS to the conventional risk model can provide a modest
293 improvement in the prediction of lung cancer in never-smoking women.

294 Most of the existing screening criteria for lung cancer, both the NLST (27) and USPSTF
295 (30), focus on heavy smokers, while the proportion of non-smoking lung cancer is increasing,
296 and there has been no clear guideline for screening criteria. The PRS can identify individuals
297 with relatively high risk and provide valuable information for individual risk assessments to
298 guide individualized screening strategies. A previous study derived a PRS for lung cancer in
299 never-smoking Asian women based on 10 SNP and examined the interaction between the PRS
300 and household coal use(17); however, the performance of PRS lacked validation in an
301 independent cohort study and thus its prediction value could not be examined. Besides, our
302 previous study used 19 SNPs to construct a PRS for NSCLC and showed predictive effects in
303 the CKB cohort comprising 95,408 individuals (19). In the present study, we generated a PRS
304 with 21 SNPs based on a large-sample, multi-center GWAS dataset of women with Asian
305 ancestry, and proved that the GWAS-derived PRS significantly predicted the risk of lung
306 cancer in never-smoking women in an independent nationwide prospective cohort study, and
307 had better performance prediction than the two known PRSs. Moreover, the PRS-21 displayed
308 the ability to substantially stratify lung cancer in never-smoking women risk trajectories
309 within each genetic risk category. Participants in the top 20% of PRS-21 compared with those
310 in the bottom 20% had a 2.09-fold risk of developing lung cancer. Furthermore, our study
311 revealed that individuals within the top 1% of the PRS exhibit a five-year lung cancer risk of
312 2.48%, a significant enough value to warrant the implementation of lung cancer screening.
313 (27). Consequently, the utilization of genetic risk score may aid in formulating a tailored
314 screening strategy for never-smoking women, taking into account their unique risk factors.
315 These findings suggest that PRS-21 may provide complementary information to current lung
316 cancer screening guidelines to better identify high-risk groups, especially in never-smoking
317 women.

318 The incorporation of PRS and nongenetic risk factors have been reported to improve the
319 risk discrimination of diseases such as coronary heart disease (8), breast (31), gastric (11) and
320 esophageal cancer (32). In a recent analysis based on the UK Biobank prospective cohort,
321 adding the 128-SNP European PRS to the epidemiologic risk factors model yielded an
322 increment of AUC of about 1.7% in never-smoking lung cancer events (33). In our analyses,
323 we also demonstrated that PRS-21 was a significant predictor of lung cancer in never-
324 smoking Asian women events, with a modest improvement in prediction accuracy when the

325 PRS was added to a conventional model (difference AUC, 1.4%; $P < 0.001$). Further studies to
326 validate a comprehensive model for individualized risk prediction based on the combined
327 effects of genetic and conventional risk factors will be needed.

328 In addition, it should be noted that the new locus at 5q11.2 identified by our GWAS
329 meta-analysis is associated with the risk of breast cancer (34) and epithelial ovarian cancer
330 (35), but the first time involved in the risk of lung cancer in never-smoking women. The lead
331 SNP rs1498606 is located in the intronic region of the *PDE4D* gene on chromosome 5q11.2,
332 and functional annotation showed that rs6450500 (with rs1498606 $r^2 = 1.0$) at 5q11.2 localizes
333 in a region with histone modification signatures, suggesting that the locus has transcriptional
334 enhancer and promoter activity. *PDE4D*, a tumor-promoting factor (36), plays a role in cancer
335 pathophysiology by altering tumor cell proliferation (37, 38), tumor cell differentiation (39),
336 and loss of E-cadherin (40). However, there is currently a lack of research on the relationship
337 between target genes and the pathogenesis of lung cancer in never-smoking women.
338 Therefore, further studies to elucidate the underlying biological mechanisms are warranted.

339 Our study has several important strengths. We developed PRS based on a large-sample,
340 multi-center GWAS database of lung cancer in never-smoking women, which provided robust
341 effect size estimates for SNPs. Furthermore, the validation samples from the CKB cohort
342 were utterly independent of the GWAS meta-analysis, which avoided overfitting of the PRS
343 effect and confirmed the predictive value of the PRS. However, we also acknowledged several
344 limitations of this study. First, missing genotype data were imputed in GWAS using different
345 genotyping arrays, which may lead to partial deviation in estimates of SNP effects. Second,
346 we only assessed the overall risk of lung cancer, but the risk estimates might differ by tumor
347 stage or subtype; however, this information was not available in the CKB cohorts, making
348 further analysis limited. Third, the construction and evaluation of the conventional risk model
349 were constructed in the CKB cohort, which may lead to overfitting to some extent. However,
350 this would not affect the discriminative power of assessing the addition of PRS to the
351 conventional risk model.

352 In summary, we developed a GWAS-derived PRS and validated its effective prediction
353 value in a large cohort for the risk of lung cancer in never-smoking women. These findings
354 suggest that the PRS-21 may be useful for identifying high-risk individuals and early
355 screening of lung cancer in never-smoking women.

356 **Author Contributions**

357 XW, DS, JG, JL and HM contributed to the study design and sample collection, and

358 supervised the whole project. XW, JZ, MZ, and HM contributed to the data interpretation,
359 data analysis, and writing of the manuscript. CY, ZM, YF, CJ, PP, LY, IM, RW, YC, HD, GJ,
360 ZC, ZH, LL, and HS contributed to the study design, sample collection, data interpretation of
361 the present analysis. All of the authors reviewed or revised the manuscript. The work reported
362 in the paper has been performed by the authors, unless clearly specified in the text.

363 **Acknowledgments**

364 We thank all the study participants and research staff for their contributions and commitment
365 to the present study.

366 **Funding**

367 This work was supported by National Natural Science Foundation of China (81922061,
368 81973123, 81803306); Natural Science Foundation of Jiangsu Province (BK20180675);
369 CAMS Innovation Fund for Medical Sciences (2019RU038); and National Science
370 Foundation for Post-doctoral Scientists of China (2018M640466).

371 **Conflict of interest**

372 No competing interests to declare.

373 **Ethics Statement**

374 Written informed consent was obtained from all participants in this study. The institutional
375 review boards of Nanjing Medical University approved the study protocol. The China
376 Kadoorie Biobank was approved by the Ethical Review Committee of the Oxford Tropical
377 Research Ethics Committee, University of Oxford, and the Chinese Center for Disease
378 Control and Prevention.

379 **Data Availability Statement**

380 The GWAS summary statistics are available through the NHGRI-EBI GWAS Catalog under
381 study accession number GCST90277434. The acquisition policies and procedures of China
382 Kadoorie Biobank are available at www.ckbiobank.org. Further information is available from
383 the corresponding author upon request.

384 **References**

- 385 1. H. Sung, J. Ferlay, R. L. Siegel *et al.*, Global Cancer Statistics 2020: GLOBOCAN
386 Estimates of Incidence and Mortality Worldwide for 36 Cancers in 185 Countries. *CA:
387 a cancer journal for clinicians*. 2021; **71**: 209-249.
- 388 2. R. Zheng, S. Zhang, H. Zeng *et al.*, Cancer incidence and mortality in China, 2016.
389 *Journal of the National Cancer Center*. 2022; **2**: 1-9.
- 390 3. C. K. Toh, W. S. Ong, W. T. Lim *et al.*, A Decade of Never-smokers Among Lung
391 Cancer Patients-Increasing Trend and Improved Survival. *Clinical lung cancer*. 2018;
392 **19**: e539-e550.
- 393 4. H. A. Wakelee, E. T. Chang, S. L. Gomez *et al.*, Lung cancer incidence in never
394 smokers. *Journal of clinical oncology : official journal of the American Society of
395 Clinical Oncology*. 2007; **25**: 472-478.
- 396 5. D. Planchard, B. Besse, Lung cancer in never-smokers. *The European respiratory
397 journal*. 2015; **45**: 1214-1217.
- 398 6. J. Subramanian, R. Govindan, Lung cancer in never smokers: a review. *Journal of
399 clinical oncology : official journal of the American Society of Clinical Oncology*.
400 2007; **25**: 561-570.
- 401 7. S. Sun, J. H. Schiller, A. F. Gazdar, Lung cancer in never smokers--a different disease.
402 *Nature reviews. Cancer*. 2007; **7**: 778-790.
- 403 8. X. Lu, Z. Liu, Q. Cui *et al.*, A polygenic risk score improves risk stratification of
404 coronary artery disease: a large-scale prospective Chinese cohort study. *European
405 heart journal*. 2022; **43**: 1702-1711.
- 406 9. N. Mavaddat, K. Michailidou, J. Dennis *et al.*, Polygenic Risk Scores for Prediction of
407 Breast Cancer and Breast Cancer Subtypes. *American journal of human genetics*.
408 2019; **104**: 21-34.
- 409 10. A. V. Khera, M. Chaffin, K. G. Aragam *et al.*, Genome-wide polygenic scores for
410 common diseases identify individuals with risk equivalent to monogenic mutations.
411 *Nature genetics*. 2018; **50**: 1219-1224.
- 412 11. G. Jin, J. Lv, M. Yang *et al.*, Genetic risk, incident gastric cancer, and healthy lifestyle:
413 a meta-analysis of genome-wide association studies and prospective cohort study. *The
414 Lancet. Oncology*. 2020; **21**: 1378-1386.
- 415 12. N. Zeinomar, W. K. Chung, Cases in Precision Medicine: The Role of Polygenic Risk
416 Scores in Breast Cancer Risk Assessment. *Annals of internal medicine*. 2021; **174**:
417 408-412.

- 418 13. W. J. Seow, K. Matsuo, C. A. Hsiung *et al.*, Association between GWAS-identified
419 lung adenocarcinoma susceptibility loci and EGFR mutations in never-smoking Asian
420 women, and comparison with findings from Western populations. *Human molecular*
421 *genetics*. 2017; **26**: 454-465.
- 422 14. Q. Lan, C. A. Hsiung, K. Matsuo *et al.*, Genome-wide association analysis identifies
423 new lung cancer susceptibility loci in never-smoking women in Asia. *Nature genetics*.
424 2012; **44**: 1330-1335.
- 425 15. K. Shiraishi, H. Kunitoh, Y. Daigo *et al.*, A genome-wide association study identifies
426 two new susceptibility loci for lung adenocarcinoma in the Japanese population.
427 *Nature genetics*. 2012; **44**: 900-903.
- 428 16. C. A. Hsiung, Q. Lan, Y. C. Hong *et al.*, The 5p15.33 locus is associated with risk of
429 lung adenocarcinoma in never-smoking females in Asia. *PLoS Genet*. 2010; **6**.
- 430 17. B. Blechter, J. Y. Y. Wong, C. Agnes Hsiung *et al.*, Sub-multiplicative interaction
431 between polygenic risk score and household coal use in relation to lung
432 adenocarcinoma among never-smoking women in Asia. *Environment international*.
433 2021; **147**: 105975.
- 434 18. Z. Hu, C. Wu, Y. Shi *et al.*, A genome-wide association study identifies two new lung
435 cancer susceptibility loci at 13q12.12 and 22q12.2 in Han Chinese. *Nature genetics*.
436 2011; **43**: 792-796.
- 437 19. J. Dai, J. Lv, M. Zhu *et al.*, Identification of risk loci and a polygenic risk score for
438 lung cancer: a large-scale prospective cohort study in Chinese populations. *The*
439 *Lancet. Respiratory medicine*. 2019; **7**: 881-891.
- 440 20. Z. Chen, J. Chen, R. Collins *et al.*, China Kadoorie Biobank of 0.5 million people:
441 survey methods, baseline characteristics and long-term follow-up. *Int J Epidemiol*.
442 2011; **40**: 1652-1666.
- 443 21. O. Delaneau, J. F. Zagury, J. Marchini, Improved whole-chromosome phasing for
444 disease and population genetic studies. *Nature methods*. 2013; **10**: 5-6.
- 445 22. O. Delaneau, J. Marchini, J. F. Zagury, A linear complexity phasing method for
446 thousands of genomes. *Nature methods*. 2011; **9**: 179-181.
- 447 23. B. N. Howie, P. Donnelly, J. Marchini, A flexible and accurate genotype imputation
448 method for the next generation of genome-wide association studies. *PLoS genetics*.
449 2009; **5**: e1000529.
- 450 24. J. Marchini, B. Howie, S. Myers *et al.*, A new multipoint method for genome-wide
451 association studies by imputation of genotypes. *Nature genetics*. 2007; **39**: 906-913.

- 452 25. C. J. Willer, Y. Li, G. R. Abecasis, METAL: fast and efficient meta-analysis of
453 genomewide association scans. *Bioinformatics (Oxford, England)*. 2010; **26**: 2190-
454 2191.
- 455 26. H. Heinzl, A. Kaider, Gaining more flexibility in Cox proportional hazards regression
456 models with cubic spline functions. *Computer methods and programs in biomedicine*.
457 1997; **54**: 201-208.
- 458 27. M. C. Tammemägi, T. R. Church, W. G. Hocking *et al.*, Evaluation of the lung cancer
459 risks at which to screen ever- and never-smokers: screening rules applied to the PLCO
460 and NLST cohorts. *PLoS medicine*. 2014; **11**: e1001764.
- 461 28. L. H. Chien, C. H. Chen, T. Y. Chen *et al.*, Predicting Lung Cancer Occurrence in
462 Never-Smoking Females in Asia: TNSF-SQ, a Prediction Model. *Cancer
463 epidemiology, biomarkers & prevention : a publication of the American Association
464 for Cancer Research, cosponsored by the American Society of Preventive Oncology*.
465 2020; **29**: 452-459.
- 466 29. F. Wang, F. Tan, S. Shen *et al.*, Risk-stratified Approach for Never- and Ever-Smokers
467 in Lung Cancer Screening: A Prospective Cohort Study in China. *American journal of
468 respiratory and critical care medicine*. 2023; **207**: 77-88.
- 469 30. V. A. Moyer, Screening for lung cancer: U.S. Preventive Services Task Force
470 recommendation statement. *Annals of internal medicine*. 2014; **160**: 330-338.
- 471 31. A. Lee, N. Mavaddat, A. N. Wilcox *et al.*, BOADICEA: a comprehensive breast
472 cancer risk prediction model incorporating genetic and nongenetic risk factors.
473 *Genetics in medicine : official journal of the American College of Medical Genetics*.
474 2019; **21**: 1708-1718.
- 475 32. J. Dong, M. F. Buas, P. Gharahkhani *et al.*, Determining Risk of Barrett's Esophagus
476 and Esophageal Adenocarcinoma Based on Epidemiologic Factors and
477 Genetic Variants. *Gastroenterology*. 2018; **154**: 1273-1281.e1273.
- 478 33. R. J. Hung, M. T. Warkentin, Y. Brhane *et al.*, Assessing Lung Cancer Absolute Risk
479 Trajectory Based on a Polygenic Risk Model. *Cancer research*. 2021; **81**: 1607-1615.
- 480 34. M. Xu, Y. Xu, M. Chen *et al.*, Association study confirms two susceptibility loci for
481 breast cancer in Chinese Han women. *Breast cancer research and treatment*. 2016;
482 **159**: 433-442.
- 483 35. J. B. Permuth, A. Pirie, Y. Ann Chen *et al.*, Exome genotyping arrays to identify rare
484 and low frequency variants associated with epithelial ovarian cancer risk. *Human
485 molecular genetics*. 2016; **25**: 3600-3612.

- 486 36. D. C. Lin, L. Xu, L. W. Ding *et al.*, Genomic and functional characterizations of
487 phosphodiesterase subtype 4D in human cancers. *Proceedings of the National*
488 *Academy of Sciences of the United States of America*. 2013; **110**: 6109-6114.
- 489 37. N. He, N. Kim, M. Song *et al.*, Integrated analysis of transcriptomes of cancer cell
490 lines and patient samples reveals STK11/LKB1-driven regulation of cAMP
491 phosphodiesterase-4D. *Molecular cancer therapeutics*. 2014; **13**: 2463-2473.
- 492 38. S. S. Pullamsetti, G. A. Banat, A. Schmall *et al.*, Phosphodiesterase-4 promotes
493 proliferation and angiogenesis of lung cancer by crosstalk with HIF. *Oncogene*. 2013;
494 **32**: 1121-1134.
- 495 39. F. Baty, D. Klingbiel, F. Zappa *et al.*, High-throughput alternative splicing detection
496 using dually constrained correspondence analysis (DCCA). *Journal of biomedical*
497 *informatics*. 2015; **58**: 175-185.
- 498 40. E. Kolosionek, R. Savai, H. A. Ghofrani *et al.*, Expression and activity of
499 phosphodiesterase isoforms during epithelial mesenchymal transition: the role of
500 phosphodiesterase 4. *Molecular biology of the cell*. 2009; **20**: 4751-4765.

501

502 **Table 1. Association of GWAS identified SNPs and risk of lung cancer in never-smoking women**

Locus	Variants	Position ^a	Gene	EA	OA	EAF	OR (95%CI)	P value
New loci								
5q11.2	rs1498606	58337549	PDE4D	C	T	0.29	1.15 (1.09-1.20)	4.79E-08
Known loci								
3q28	rs11928222	189350265	TP63	G	T	0.43	1.21 (1.15-1.26)	3.09E-16
5p15.33	rs7726159	1282319	CMT1M/TERT	A	C	0.38	1.42 (1.35-1.49)	1.10E-46
6p21.32	rs9275164	32652929	HLA-DQB1	C	T	0.34	1.20 (1.14-1.26)	8.43E-13
6p21.1	rs2496646	41483367	FOXP4	T	C	0.35	1.18 (1.12-1.24)	3.21E-10
6q22.1	rs9374663	117782634	ROS1/DCBLD1	A	G	0.50	1.15 (1.10-1.20)	1.56E-09
10q25.2	rs11196063	114460013	VTG1A	C	A	0.28	1.27 (1.21-1.33)	1.89E-21
11q23.3	rs55768116	118108331	AMICA1	C	A	0.45	1.15 (1.09-1.20)	3.92E-09
17q24.2	rs7284527 8	65877076	BPTF	T	C	0.71	1.20 (1.13-1.26)	9.09E-11

503 *Abbreviations:* EA, effect allele; OA, other allele; EAF, effect allele frequency; OR, odds (log-additive)
504 ratio; 95% CI, 95% confidence interval.

505 ^aChromosome position, hg19/GRCh37 build.

506

507 **Table 2. Association of different polygenic risk scores with risk of lung cancer in never-smoking**
 508 **women in the CKB cohort**

PRS development methods	N SNPs in the PRS	HR (95% CI) ^a	P value	AUC
Clumping and thresholding approach^b				
5×10 ⁻⁸	12	1.24 (1.13-1.35)	2.75E-06	0.565
5×10 ⁻⁷	16	1.28 (1.17-1.40)	9.37E-08	0.573
5×10 ⁻⁶	44	1.25 (1.14-1.36)	1.60E-06	0.568
5×10 ⁻⁵	202	1.16 (1.06-1.26)	1.77E-03	0.541
5×10 ⁻⁴	1448	1.08 (0.98-1.18)	1.16E-01	0.524
Integrate PRS				
PRS-21	21	1.29 (1.18-1.41)	4.40E-08	0.576
Published PRS				
PRS-10(17)	10	1.21 (1.11-1.33)	3.77E-05	0.557
PRS-19(19)	19	1.20 (1.09-1.09)	1.31E-04	0.547

509 *Abbreviations:* CI, confidence interval; HR, hazard ratio; PRS, polygenic risk score; SNP, single nucleotide
 510 polymorphism.

511 ^aAdjusted for age, education, BMI, COPD history, personal history of cancer, family history of cancer, and
 512 the first ten principal components of ancestry. HR denotes the increase of hazard risk of lung cancer in
 513 never-smoking women per SD for PRS

514 ^bPolygenic risk score (PRS) derived from different significance thresholds in the GWAS meta-analysis of
 515 lung cancer in never-smoking women

516 **Figure 1. Regional plot and forest plot for rs1498606.** (A) Regional plot of chromosome 5q11.2 locus.
517 The index SNP rs1498606 is colored purple, with other SNPs colored according to the degree of linkage
518 disequilibrium r^2 with the index SNP. The r^2 values are calculated based on the 1000 Genome data. Each
519 point represents a variant with chromosomal position on the x axis and the $-\log(P \text{ value})$ on the y axis. (B)
520 Forest plot for rs1498606 association with risk of lung cancer in never-smoking women in each GWAS
521 study.

522 **Figure 2. The distribution and relationship of polygenic risk score (PRS) with incident lung cancer in**
523 **never-smoking women in the CKB cohort.** (A) Distribution of PRS in participants affected with or
524 without lung cancer in never-smoking women in the CKB cohort; (B) Linear relationship between PRS and
525 risk of lung cancer in never-smoking women was assessed using a restricted cubic spline analysis in the
526 CKB cohort; (C) Participants in the CKB cohort were divided into five equal groups according to their
527 polygenic risk score, and the HRs for each group were compared with those in quintile 1 of the polygenic
528 risk score; error bars show the 95% CIs; (D) The standardized cumulative incidence of lung cancer in
529 never-smoking women in low (quintile 1), intermediate (quintiles 2–4), and high (quintile 5) genetic risk
530 groups in the CKB cohort during median 10.5 years of follow-up; shaded regions represent the 95% CIs,
531 and hazard ratios (HRs) were estimated with adjustment for age, education, BMI, COPD history, personal
532 history of cancer, family history of cancer, and the first ten principal components of ancestry.

533 **Figure 3. Absolute risk estimates of lung cancer in never-smoking women in the CKB cohort.** (A)
534 Five-year absolute risk; (B) Cumulative absolute risk until age 80. Cox proportional hazard model was used
535 to estimate the HRs (95% CIs) and the cumulative risk adjusted for age, education, BMI, COPD history,
536 personal history of cancer, family history of cancer, and the first ten principal components of ancestry.
537 Genetic risk categories: low (quintile 1), intermediate (quintiles 2-4), and high (quintile 5) risk according to
538 quintiles of PRS.

539

540

541

542

543

544