1    Lung cancer, genetic predisposition and smoking: the Nordic Twin Study of Cancer

2

3    Jacob Hjelmborg*, Ph.D., Tellervo Korhonen*, Ph.D., Klaus Holst, Ph.D., Axel Skytthe,

4    Ph.D., Eero Pukkala, Ph.D., Julia Kutschke, Ph.D., Jennifer R. Harris, Ph.D., Lorelei A.

5    Mucci, Sc.D., Kaare Christensen, M.D., Ph.D., Kamila Czene, Ph.D., Hans-Olov Adami,

6    M.D., Ph.D., Thomas Scheike, Ph.D., and Jaakko Kaprio, M.D, Ph.D., *on behalf of the*

7    *Nordic Twin Study of Cancer (NorTwinCan) collaboration.*

8    *Drs. Hjelmborg and Korhonen contributed equally to this article.

9

10    **Author affiliations:**

11    Department of Epidemiology, Biostatistics and Biodemography, University of Southern

12    Denmark, Denmark (J.H., A.S., K.C.); The Danish Twin Registry, University of Southern

13    Denmark, Denmark (J.H., A.S., K.C.); Department of Biostatistics, University of

14    Copenhagen, Denmark (K.H., T.S.); Channing Division of Network Medicine, Brigham

15    and Women's Hospital, Harvard Medical School, Boston MA; (L.A.M); Centre for

16    Public Health Sciences, University of Iceland, Reykjavik, Iceland (L.A.M); Finnish

17    Cancer Registry, Institute for Statistical and Epidemiological Cancer Research, Helsinki,

18    Finland (E.P): School of Health Sciences, University of Tampere, Tampere, Finland

19    (E.P.); Division of Epidemiology, The Norwegian Institute of Public Health, Oslo,

20    Norway (J.Ku., J.R.H.); Department of Public Health, University of Helsinki, Finland

21    and Department of Health, National Institute for Health & Welfare, Helsinki, Finland

22    (J.Ka., T.K); Department of Medical Epidemiology and Biostatistics, Karolinska

23    Institutet, Stockholm, Sweden (K.Cz., H-O.A.); and Department of Epidemiology,

24 Harvard School of Public Health, Boston, MA, USA (L.A.M., H-O.A); Institute for

25 Molecular Medicine, University of Helsinki, Finland (J.Ka.); Institute of Public Health

26 and Clinical Nutrition, University of Eastern Finland (T.K.)

27

28

29 **Corresponding author:**

30 Jaakko Kaprio, Department of Public Health, PO Box 41, FI-00014 University of

31 Helsinki, Finland; Phone: +358 50 3715419; Email: Jaakko.kaprio@helsinki.fi

32

33 **Key messages**

34

35 What is the key question?

36

37 Is there a significant genetic component to the occurrence of lung cancer and is the
38 genetic influence modified by smoking and age?

39

40 What is the bottom line?

41

42 The interplay between genes and tobacco smoking in the etiology of lung cancer has
43 remained controversial, and we disentangle genetic and environmental causes in cancer
44 while taking smoking status into account.

45

46 Why read on?

47

48 Our study shows that tobacco exposure causes lung cancer even when adjusting for
49 genetic factors. Interactions between genes and environmental exposure in the
50 development of lung cancer are not supported from the largest twin cohort study with
51 longest follow-up ever. Familial effects have decreased influence with increasing age.

52

53

54 **Word count:** Abstract 250; text 3336

55

56 **Abstract**

57 **Background**

58 We aimed to disentangle genetic and environmental causes in lung cancer while

59 considering smoking status.

60 **Methods**

61 Four Nordic Twin Cohorts (43,512 monozygotic (MZ) and 71,895 same sex dizygotic

62 (DZ) twin individuals) had smoking data before cancer diagnosis. We used time-to-event

63 analyses accounting for censoring and competing risk of death to estimate incidence,

64 concordance risk and heritability of liability to develop lung cancer by smoking status.

65 **Results**

66 During a median of 28.5 years of follow-up we recorded 1,508 incident lung cancers. Of

67 the 30 MZ and 28 DZ pairs concordant for lung cancer, nearly all were current smokers at

68 baseline and only one concordant pair was seen among never smokers.  Among ever

69 smokers the case-wise concordance of lung cancer, that is the risk before a certain age

70 conditional on lung cancer in the co-twin before that age was significantly increased

71 compared with the cumulative incidence for both MZ and DZ pairs. This ratio, the

72 relative recurrence risk, significantly decreased by age for MZ, but was constant for DZ

73 pairs.  Heritability of lung cancer was 0.41 (95%CI 0.26–0.56) for currently smoking and

74 0.37 (95%CI 0.25–0.49) for ever smoking pairs. Among smoking discordant pairs, the

75 pairwise hazard ratio for lung cancer of the ever smoker twin compared to the never

76 smoker cot-win was 5.4 (95%CI 2.1–14.0) in MZ pairs and 5.0 (95%CI 3.2–7.9) in DZ

77 pairs.

78 **Conclusions**

79    The contribution of familial effects appears to decrease by age. The discordant pair

80    analysis confirms that smoking causes lung cancer.

81

82

83 **Introduction**

84      Smoking is the primary cause of lung cancer globally, though several other

85 environmental exposures play a role.[1] The estimated heritable genetic contribution to

86 variation in risk to lung cancer overall has been modest in family (heritability estimate of

87 0.08)[2] and twin (0.26[3] and 0.18[4]) studies. Genome-wide association (GWA) studies

88 further suggest that some gene loci are associated with lung cancer in both smokers and

89 non-smokers, while other variants, such as the functional D398N (rs16969968) variant in

90 CHRNA5, are associated with lung cancer only among smokers.[5,6] Thus, the heritability

91 of  lung cancer may vary as a function of smoking, but the differential effect of smoking

92 on genetic variation underlying development of lung cancer has not been quantified.

93      To this end, our aim is to estimate the heritability of liability to lung cancer based

94 on the largest twin cohort to date, the Nordic Twin Study of Cancer (NorTwinCan)[4],

95 which extends the Lichtenstein (2000)[3] study with longer follow-up and new birth

96 cohorts and refined methodology.  We sought to estimate the heritability in the liability to

97 lung cancer and whether it is modified by smoking or age.

98

99 **Methods**

100 Material

101      NorTwinCan includes population-based cohorts from the Danish, Finnish,

102 Norwegian, and Swedish twin registries.[7] Each twin has an individually unique national

103 registration number, allowing for linkage to the national cancer and mortality registries

104 with complete follow-up, drop-out being only due to death or emigration. Lung cancer

105 occurrence was obtained from the national cancer registries and computed from the

106 baseline when smoking status was determined until the end of follow-up (Table 1). In all

107 cohorts, zygosity - monozygotic (MZ) or dizygotic (DZ) - was determined at baseline by

108 validated questionnaire methodology, which classifies more than 95% of twin pairs

109 correctly.[3] Twins, who have not replied to the questionnaires, as well as a minority

110 providing inconsistent responses, are classified as unknown zygosity (UZ). The ethics

111 committees for each country approved the study.

112 Given the major role of smoking in the etiology of lung cancer, our analysis

113 includes twin individuals of known zygosity from the Danish, Finnish, Norwegian, and

114 Swedish registries, where data on smoking status was available prior to lung cancer

115 diagnosis. We excluded individuals from opposite-sex DZ pairs as data from them have

116 not been as comprehensively collected. For individuals who reported smoking behavior

117 on more than one questionnaire, we used the earlier information.

118 Characteristics of the four national twin cohorts included in the analyses are

119 summarized in Table 1. We classified the participants as never smokers, ever smokers

120 (former or current at time of questionnaire) and current smokers based on the survey

121 items used to assess smoking status. Smoking data in the Danish cohort came from the

122 eight questionnaire surveys conducted from 1959 to 2002.[8–10] In Finland smoking data

123 came primarily from the first questionnaire survey in 1975, but some twins who had not

124 replied in 1975 responded to a questionnaire survey in 1981.[11,12] In the Norwegian cohort

125 smoking data came from three questionnaire surveys in 1980–1982 & 1990–92 &

126 1998.[13,14] In the Swedish cohort smoking data came from questionnaire surveys in 1961,

127 1967, 1970, and 1973.[15,16]

128  We included individuals with histologically confirmed lung cancer.  Among those

129 with smoking data, we recorded a total of 1,508 incident lung cancers with a mean

130 follow-up time of 25.2 years (21.0 years in lung cancer patients).

131

132 Statistical analysis

133  After defining cohort-specific dates of entry and follow-up, we accounted for left-

134 truncation from variable initiation of cancer registration and right-censoring among those

135 censored at the end of follow-up, and lost to follow-up due to emigration (<2%). We

136 examined the individual risk of lung cancer diagnosis by age by estimating cumulative

137 lung cancer, incidence[17] and lifetime risk as the cumulative incidence (the probability of

138 lung cancer) by age 80 years. We modeled potential competing deaths[18,19] which allows

139 estimation of lung cancer risk in a twin given the occurrence of other disease in his/her

140 co-twin.  We obtained the case-wise concordances by age[18,19] (see supplementary

141 material for details) as well as relative recurrence risks in MZ and DZ pairs and the

142 multilocus index.[20,21]

143  We extended standard biometrical modelling methods to address issues of

144 censoring at follow-up[7,22] . Results would agree with those obtained from standard

145 models for twin data[18,23,24] if no censoring  were present. Quantitative models were

146 analyzed to estimate the magnitude of variation explained by genetic and environmental

147 influences[18] underlying the liability to develop lung cancer by smoking status. The

148 relative magnitude of genetic influences on variation in liability to lung cancer is thus

149 estimated among pairs in which neither had ever smoked, among pairs where both co-

150    twins are ever (former or current) smokers and among pairs in which both co-twins are

151    current smokers.

152         We use information on lung cancer incidence in MZ and DZ pairs to decompose

153    variation into additive genetic effects (A), dominant genetic effects (which represent

154    deviations of the heterozygote genotype from the mean of the homozygote genotype) (D),

155    common environmental effects (C), and individually unique environmental effects (E).

156    Within-pair covariance of liability is expressed as $\kappa$ var(A) + $\gamma$ var(D) +var(C), where $\kappa$ =

157    $\gamma$ = 1 for MZ pairs and $\kappa$ = 1/2 and $\gamma$ = 1/4 for DZ pairs.[18] We tested a series of models

158    sequentially to assess the significance of specific parameters. We estimated measurement

159    error in E which is the component of variance that does not contribute to within-pair

160    resemblance. Dominance effects are, typically, biologically implausible in the absence of

161    additive effects. The primary models are thus the ACE and ADE models, as well as their

162    sub-models AE, CE, and E. We assessed the fit of the sub-models by the Akaike

163    information criterion[22].

164         We tested for equal thresholds (i.e., normal quantiles of prevalence) between MZ

165    and DZ twins, which is equivalent to assuming that the risk of disease does not differ by

166    zygosity.  We tested for constant relative recurrence risk (RRR) over age by grouping

167    into five-year interval from age 65 to 90 years of age for MZ and DZ pairs. To correct for

168    possible bias due to censoring, individuals were assigned weights obtained by calculating

169    the inverse probability of being censored at time of follow-up[7,18,19,22]  Estimates have not

170    been adjusted for the effect of left-truncation that would cause an upwards bias, which is

171    not yet feasible for the approach.

172   For gene and smoking status interaction the magnitude on liability scale could not

173 be estimated due to having one concordant pair among all never-never and never-ever

174 smoking pairs. The presence of genetic interaction with smoking status was therefore

175 investigated by comparing observed concordance in strata of smoking status  to the

176 expected when assuming same variance components on the liability scale as in ever-ever

177 pairs but using smoking-status specific cumulative incidence by age as well as follow-up

178 time of the specific pairs in the cohort. This procedure leads to an approximate test,

179 which we later refer to as the binomial test, and takes into account the smoking-status

180 specific cumulative incidence by age, as well as follow-up time of the specific pairs in the

181 cohort and we then computed the probability that a randomly selected pairs were

182 concordant using the dependence parameters of the liability threshold model for the ever-

183 ever pairs.

184   Among pairs in which one twin was a smoker and the other was not, we computed

185 within pair hazard ratios for the association of smoking with lung cancer using a Cox

186 model with pair-specific baseline hazard functions. Given that MZ pairs share their

187 genomic sequence, an association of smoking with lung cancer risk within such pairs is

188 independent of genetic liability. This hypothesis has historically competed with the

189 hypothesis[25] of shared genes underlying both smoking and lung cancer. The statistical

190 program R was used for all analyses with the package *mets*.[26]

191

192

193

194    **Results**

195    Among those with smoking data, we recorded 1,508 incident lung cancers among

196    a total of 115,407 twin (43,512 MZ and 71,895 DZ) individuals.  Forty-seven percent

197    were never smokers (n=54238), 16% former smokers (n=18,231) and 37% current

198    smokers (n=42,938) at baseline. Figure 1 shows the cumulative incidence of lung cancer

199    by smoking status (never, former, current) and sex. The risk of lung cancer diagnosis

200    before 80 years of age is estimated at 0.6% (95% CI 0.5%–0.7%) among never smokers,

201    2.0% (1.7%–2.3%) among former and 5.7% (5.4%–6.0%) among current smokers

202    adjusting for censoring and competing risk of death.  The only sex difference is seen

203    among smokers. There was no difference in risk between MZ and DZ twin individuals.

204    The numbers of pairs concordant and discordant for lung cancer incidence are

205    presented in Table 2 for those with smoking data (n=50,595 pairs with smoking status on

206    both twins) overall and further classified by smoking status.

207    Among twin pairs where both are ever smokers, the risk of lung cancer in a twin

208    before a given age given that his or her co-twin also has lung cancer before that age, the

209    case-wise concordance by age is depicted in Figure 2 in both MZ and DZ pairs, as well as

210    the cumulative incidence of lung cancer by age in individuals.  The case-wise

211    concordance risk was larger in MZ twins than the individual cumulative incidence risk,

212    testing for a difference from the cumulative incidence across the five year age intervals

213    (chisq=22.1, df=6, p=0.001). For the DZ twins we found that the case-wise concordances

214    were borderline significantly different from the cumulative incidence (chisq=13.4, df=6,

215    p=0.04). The estimated case-wise concordance at 90 years of age was 0.20 (0.13-0.27) for

216    MZ pairs and 0.13 (0.08-0.17) for DZ pairs.

217       This excess risk of MZ and DZ pairs of the case-wise concordance relative to the

218      population based individual cumulative incidence of lung cancer, the relative recurrence

219      risk (also known as the lambda value) is depicted in Figure 3 and demonstrates the

220      presence of familial effects at all ages. The RRR is higher at younger ages, in fact the

221      lung cancer risk is increased 10.2 -fold (3.2-17.2) at 65 years of age and decreases

222      significantly to a 3.6 (2.3-4.9) -fold increase at 90 years of age if a MZ co-twin is

223      diagnosed (p-value = 0.04, test for trend). The RRR is suggested to be constant by age for

224      DZ twins (p-value = 0.25, test for trend) (Figure 3). (A table of relative risks by age-

225      group is provided in supplemental Table 1.) We tested if the absolute differences of the

226      MZ and DZ curves at each five-year interval from age 65 to age 90 years of age were

227      significantly different, which there was no sign of (p-value=0.21). Our results are thus

228      consistent with the hypothesis of rather strong familial influences that do not increase

229      across age. We hypothesize that the genetic part of the familial influence may become

230      weaker by age.

231       We then examined evidence for genetic factors in the liability to develop lung

232      cancer by smoking status. Among pairs in which neither had ever smoked (7,871 MZ

233      pairs and 10,768 DZ pairs), there was one lung cancer concordant MZ pair with 43 MZ

234      and 59 DZ lung cancer discordant pairs.  Heritability could not be estimated. However,

235      the dependence in the never-never and never-ever pairs was not significantly different

236      from the dependence among the ever-ever pairs (p=0.28, binomial test of observing more

237      than one concordant pair of lung cancer).

238       The overall estimate of familial aggregation (genetic variance and shared

239      environment component) for lung cancer liability is 44% with 38% (0.05- 0.72) of

240    variability attributed to genetic effects.  When adjusted for smoking status, effects of

241    country and sex, variability attributed to genetic effects was 34% (0.00-0.70) (Table 3). A

242    comparison of the MZ and DZ tetrachoric within-pair correlations in liability to develop

243    lung cancer (Table 3) adjusting for age, sex, country and smoking, and further adjustment

244    for censoring hypothesizing equal correlations, gave a p-value of 0.07 (Wald test).

245    Among the pairs where both twins are ever (current or former) smokers, the heritability

246    estimates ranged from 28% (0.00-0.66) to 37% (0.25-0.49), depending on the

247    assumptions of the genetic model (Table 4). A pure environmental model did not fit the

248    data. Among current smokers, the heritability was estimated at 29% (0.00-0.74) or 41%

249    (0.26-0.56), depending on genetic assumptions (Table 4).

250         Finally, for smoking discordant pairs, we examined whether smoking status was

251    associated with future lung cancer. In the ever smoking discordant pairs (3,274 MZ pairs

252    and 8,350 DZ pairs), 40 MZ pairs were discordant for lung cancer (Table 5). Of these 35

253    cases were among ever smokers (with their non-smoking co-twin being unaffected) and

254    only five in the never-smokers (while their smoking co-twin was unaffected), yielding a

255    paired analysis hazard ratio (HR) of 5.4. Results for DZ pairs and for current-smoking

256    *versus* never smoking discordant pairs are shown in Table 5. Most discordant pairs arose

257    from pairs in which the smoker still smoked at baseline. None of the smoking discordant

258    pairs were concordant for lung cancer.

259

260

**Discussion**

In the largest study of lung cancer in twins to date, we found that genetic effects account for a significant amount of the variation in the liability to develop lung cancer, and the magnitude of this estimate is independent of smoking status. The largest estimate of heritability in the liability to lung cancer was found in pairs where both were current smokers at baseline. Among twin pairs where both twins were never smokers, only one concordant lung cancer pair was seen and a formal estimate of heritability could not be derived. A test of gene by smoking interaction was not significant suggesting that the relative contribution of genetics does not vary by smoking status. Furthermore, testing suggests that the contribution of familial effects does not increase by age. Our pairwise analysis of smoking discordant pairs confirmed that smoking causes lung cancer independent of genetic liability either to smoking or to lung cancer.

Twin pairs discordant for both lung cancer and smoking status at baseline are informative for causal analyses. In the lung cancer and smoking doubly discordant pairs, the pairwise relative risk for lung cancer was 5.4 among ever smokers in MZ pairs. It is of historical interest that after the landmark papers of Doll and Hill[27] and Wynder and Graham[28] in the early 1950s, the causality of the relationship between smoking and lung cancer was soon challenged by the great statistician Ronald Fisher.[25] He pointed out the greater similarity of MZ vs. DZ pairs for smoking, and indicated genetics as a potential confounder.  MZ pairs discordant for smoking would help to resolve the issue of causality. Following up on prior twin studies of smoking discordant pairs,[29,30] we can now finally put this issue to rest, an issue debated for many years because of tobacco industry's prolonged refusal to acknowledge publicly that smoking causes lung cancer.

284        Smoking is the most important cause of lung cancer. Taking smoking into

285    account permits us to test for the dependence of genetic effects on smoking status. The

286    overall estimate of familial aggregation (genetic variance and shared environment

287    component) for lung cancer liability is 44%, with most variability attributed to genetic

288    effects (38%), higher but still consistent with the estimate 26% (95%CI 0%–49%) by

289    Lichtenstein et al.[3] also unadjusted for smoking and for censoring, but based on a smaller

290    number of affected pairs. We recently reported on the heritability for liability to lung

291    cancer in the entire NorTwinCan data, with an overall estimate of familial aggregation of

292    42%.[4] The present analysis extends these estimates by accounting for the effect of

293    smoking status prior to disease occurrence and examines heritability among the smoking

294    pairs.

295        In our analysis, adjustment for smoking eliminates the estimates for shared

296    environmental effects. Shared environmental effects (i.e. exposure to smokers in the

297    childhood home, and among peers in adolescence) are of importance for the initiation of

298    smoking[31] so it is not surprising that adjustment for smoking controls for this source of

299    variation. The highest estimates of heritability and recurrence risks were seen among

300    current smoking pairs. Among never smokers, we cannot estimate the heritability of lung

301    cancer.

302        Prior family[2] and twin[3,4] studies of lung cancer have demonstrated familial

303    aggregation and provided very modest estimates for the role of genes. The Swedish

304    multi-generational register family study[2] estimated the heritability of lung cancer to be

305    8% (95% CI 5%–9%), without information on smoking in the families.   The American

306    World War II veterans' study [32] followed 12,938 male twin pairs for 44 years for

307     mortality. Among pairs with at least one lung cancer death, only 10 of 269 MZ pairs and

308     21 of 373 DZ pairs were concordant, and no heritability estimate was provided. Smoking

309     information was not used in the analysis, but smoking-related cancers showed less MZ –

310     DZ differences in similarity than other cancers. Despite the large number of pairs in our

311     present study, the final number of concordant pairs with smoking information was

312     limited. Thus, we could not examine heritability of lung cancer risk in relation to time

313     trends in lung cancer or histological subtypes of lung cancer. Nor did we have

314     information on smoking amount, duration or changes in smoking status comprehensively

315     and comparably assessed in all the twin cohorts.

316         Since detailed smoking information was not available, it should be acknowledged

317     as a potential limitation that there might be residual confounding that remains in the

318     estimates of heritability estimation. Because MZ twins, who are smokers, are also more

319     similar than DZ pairs in age of smoking initiation, amount smoked and duration of

320     smoking[31], the heritability of lung cancer among smokers may still contain residuals

321     effects of genetics on smoking, and thus on lung cancer risk.

322         The overall genetic contribution to lung cancer as a function of smoking status is

323     relevant for gene discovery.  Since 2007, 21 lung cancer genome-wide analysis (GWA)

324     and genome-wide meta-analysis studies[33] ([www.genome.gov/gwastudies](http://www.genome.gov/gwastudies)) have found the

325     strongest association to the CHRNA5 functional D398N (rs16969968) variant. The

326     functional changes[34,35] in nicotinic acetylcholine receptor activity are linked to increased

327     risk for nicotine dependence, higher amount smoked[36-39] and higher cotinine levels.[40,41]

328     Thus, those with a risk allele smoke more, are more tobacco-dependent and are less likely

329     to quit, and  therefore at higher risk of developing lung cancer. However, D398N is not a

330    risk factor for lung cancer in non-smokers, based on a GWA meta-analysis of 14,900

331    lung cancer cases and 29,485 controls[6] and among 56,037 individuals from the HUNT

332    population study in Norway.[5] This variant requires exposure to smoking to affect lung

333    cancer risk and thus contributes to the heritability seen among current smokers. In

334    contrast to D398N, associations with other loci found to be significant for lung cancer

335    such as those in 5p15 (TERT and CLPTM1L genes) and 6p21 (BAG6/BAT3) are found

336    also in non-smokers.[33,6] The existence of a modest familial liability to lung cancer

337    independent of smoking status was also observed in the analysis of Utah genealogical

338    data.[42] An increased risk of lung cancer was seen even in distant relatives; the high

339    proportion of non-smoking lung cancer cases (31%) and a large proportion of missing

340    data on smoking status (which was assessed through the death certificate and not

341    prospectively) calls for replication in other populations. A recent large meta-analysis

342    yielded an array-based heritability estimate for lung cancer of 21% (95% CI 14-27%).[43]

343    This is somewhat smaller than our overall twin estimates suggesting that much of the

344    genetic liability to lung cancer is attributable to common variants, but other genetic

345    effects may exist. The same study estimated that 24% of the heritability of lung cancer is

346    accounted for by genetic determinants of smoking behavior.

347        In conclusion, our study extends earlier studies to examine the heritability in

348    liability to lung cancer by smoking status and age. We find no formal evidence for a gene

349    by environmental exposure interaction in lung cancer; more detailed environmental

350    exposures and larger sample sizes may be required. We hypothesize that a genetic part of

351    the rather strong familial influence demonstrated may become weaker by age. Studies of

352    genetic factors and hence molecular mechanisms in cancer would benefit by carefully

353    taking into account known environmental risk factors and identifying the population

354    groups at highest genetic risk using environmental stratification.  However, the discordant

355    pair analysis conclusively demonstrates that tobacco exposure causes lung cancer even

356    when adjusting for genetic factors.

357

**Contributions**

Jacob Hjelmborg (J.H.) designed the study, contributed to developing the statistical methodology, conducted the data analysis, interpreted the data, and wrote the methods section of the manuscript.

Tellervo Korhonen (T.K.) contributed to the design and wrote the manuscript together with J.H. and J.K.

(Drs. Hjelmborg and Korhonen contributed equally to this article.)

Klaus Holst made central contributions to developing the statistical methodology, took part in conducting the statistical analysis as well as in revising the manuscript.

Axel Skytthe was responsible for quality assurance of the combined data set, conducted the data analysis, reviewed and commented the manuscript.

Eero Pukkala contributed to quality assurance of the combined data set reviewed, commented and edited the manuscript.

Julia Kutschke (nee Isaeva) helped to prepare the Norwegian data.

Jennifer R. Harris helped in the drafting and providing critical comments to manuscript.

Lorelei A. Mucci reviewed, commented and edited the manuscript.

Kaare Christensen reviewed, commented and edited the manuscript.

Hans-Olov Adami was involved in initiating, designing and funding the study as well as in interpreting the results and editing the manuscript.

Thomas Scheike contributed to statistics and took part in revising the manuscript.

Jaakko Kaprio (J.K.) designed the study, contributed to data interpretation, and wrote the manuscript together with J.H. and T.K.

392    **Conflict of interest statement:**

393    Tellervo Korhonen and Jaakko Kaprio have consulted for Pfizer on nicotine dependence

394    from 2012 to 2015. Other authors declare no conflict of interest.

395

396

397     **References**

398     1.      Boffetta P, Trichopoulos D. Cancer of the lung, larynx, and pleura. In (Adami H,

399     Hunter DJ, Trichopoulos D, eds). *Textbook of cancer epidemiology*. Oxford; New York:

400     Oxford University Press, 2008.

401     2.      Czene K, Lichtenstein P, Hemminki K. Environmental and heritable causes of

402     cancer among 9.6 million individuals in the Swedish Family-Cancer Database. *Int J*

403     *Cancer* 2002;99:260–6.

404     3.      Lichtenstein P, Holm NV, Verkasalo PK, et al. Environmental and heritable

405     factors in the causation of cancer--analyses of cohorts of twins from Sweden, Denmark,

406     and Finland. *N Engl J Med* 2000;343:78–85.

407     4. Mucci LA,  Hjelmborg JB, Harris JR, et al. Familial Risk and Heritability of Cancer

408     Among Twins in Nordic Countries. *JAMA.* 2016;315(1):68-76.

409     5. Gabrielsen ME, Romundstad P, Langhammer A, Krokan HE, Skorpen F. Association

410     between a 15q25 gene variant, nicotine-related habits, lung cancer and COPD among

411     56,307 individuals from the HUNT study in Norway. *Eur J Hum Genet* 2013; 21:1293–

412     1299.

413     6.      Timofeeva MN, Hung RJ, Rafnar T, et al. Influence of common genetic variation

414     on lung cancer risk: meta-analysis of 14 900 cases and 29 485 controls. *Hum Mol Genet*

415     2012;21:4980–95.

416     7.      Hjelmborg JB, Scheike T, Holst K, et al. The Heritability of Prostate Cancer in

417     the Nordic Twin Study of Cancer. *Cancer Epidemiol Biomarkers Prev* 2014;23:2303–10.

418  8.      Wienke A, Herskind AM, Christensen K, Skytthe A, Yashin AI. The heritability

419  of CHD mortality in Danish twins after controlling for smoking and BMI. *Twin Res Hum*

420  *Genet* 2005;8:53–9.

421  9.      Osler M, McGue M, Christensen K. Socioeconomic position and twins' health: a

422  life-course analysis of 1266 pairs of middle-aged Danish twins. *Int J Epidemiol*

423  2007;36:77–83.

424  10.     Johnson W, Kyvik KO, Mortensen EL, Skytthe A, Batty GD, Deary IJ. Does

425  education confer a culture of healthy behavior? Smoking and drinking patterns in Danish

426  twins. *Am J Epidemiol* 2011;173:55–63.

427  11.     Kaprio J, Koskenvuo M. A prospective study of psychological and socioeconomic

428  characteristics, health behavior and morbidity in cigarette smokers prior to quitting

429  compared to persistent smokers and non-smokers. *J Clin Epidemiol* 1988;41:139–50.

430  12.     Kaprio J, Koskenvuo M. Genetic and environmental factors in complex diseases:

431  the older Finnish Twin Cohort. *Twin Res* 2002;5:358–65.

432  13.     Nilsen TS, Brandt I, Magnus P, Harris JR. The Norwegian Twin Registry. *Twin*

433  *Res Hum Genet*. 2012;15:775–80.

434  14.     Harris JR, Magnus P, Tambs K. The Norwegian Institute of Public Health twin

435  program of research: an update. *Twin Res Hum Genet.* 2006; 9:858–64.

436  15.     Lichtenstein P, De Faire U, Floderus B, Svartengren M, Svedberg P, Pedersen

437  NL. The Swedish Twin Registry: a unique resource for clinical, epidemiological and

438  genetic studies. *J Intern Med* 2002;252:184–205.

439  16.     Pedersen NL, Lichtenstein P, Svedberg P. The Swedish Twin Registry in the third

440  millennium. *Twin Res* 2002;5:427–32.

441    17.    Allignol A, Schumacher M, Beyersmann J. Empirical Transition Matrix of Multi-

442    State Models: The etm Package. *J Stat Softw* 2011;38.

443    18.    Scheike TH, Holst KK, Hjelmborg JB. Estimating twin concordance for bivariate

444    competing risks twin data. *Stat Med* 2014;33:1193–204.

445    19.    Scheike TH, Holst KK, Hjelmborg JB. Estimating heritability for cause specific

446    mortality based on twin studies. *Lifetime Data Anal* 2014;20:210–33.

447    20.    Risch N. Linkage strategies for genetically complex traits. I. Multilocus models.

448    *Am J Hum Genet* 1990;46:222–8.

449    21.    Risch N. The genetic epidemiology of cancer: interpreting family and twin studies

450    and their implications for molecular genetic approaches. *Cancer Epidemiol Biomarkers*

451    *Prev* 2001;10:733–741.

452    22.    Holst KK, Scheike T, Hjelmborg JB. The liability threshold model for censored

453    twin data [published online ahead of print January 2015]. *Computational Statistics &*

454    *Data Analysis* 2015; doi: 10.1016/j.csda.2015.01.014.

455    23.    Neale MC, Cardon LR, North Atlantic Treaty Organization. Scientific Affairs

456    Division. *Methodology for genetic studies of twins and families*. Dordrecht ; Boston:

457    Kluwer Academic Publishers, 1992.

458    24.    Sham P. *Statistics in human genetics.* London; New York: Arnold; John Wiley &

459    Sons, Inc., 1998.

460    25.    Fisher RA. Cancer and smoking. *Nature* 1958;182:596.

461    26.    Holst K, Scheike, T. H. mets: Analysis of Multivariate Event Times, R package

462    version 0.2.8.1, http://lava.r-forge.r-project.org/

463    27.    Doll R, Hill AB. A study of the aetiology of carcinoma of the lung. *Br Med J*

464    1952; 2:1271–1286.

465    28.    Wynder EL, Graham EA. Tobacco smoking as a possible etiologic factor in

466    bronchiogenic carcinoma: a study of 684 proved cases. *J Am Med Assoc.* 1950; 143:329–

467    336.

468    29.    Floderus B, Cederlof R, Friberg L. Smoking and mortality: a 21-year follow-up

469    based on the Swedish Twin Registry. *Int J Epidemiol* 1988;17:332–340.

470    30.    Kaprio J,  Koskenvuo M. Cigarette smoking as a cause of lung cancer and

471    coronary heart disease. A study of smoking-discordant twin pairs. *Acta Genet Med*

472    *Gemellol* (Roma) 1990;39:25–34.

473    31.    Rose RJ, Broms U, Korhonen T, Dick DM, Kaprio J. Genetics of Smoking

474    behavior. In: Kim YK, ed. *Handbook of Behavior Genetics*. New York: Springer,

475    2009:411–432.

476    32.    Braun MM, Caporaso NE, Page WF, Hoover RN. A cohort study of twins and

477    cancer. *Cancer Epidemiol Biomarkers Prev* 1995;4:469–73.

478    33.    Yang IA, Holloway JW, Fong KM. Genetic susceptibility to lung cancer and co-

479    morbidities. *J Thorac Dis* 2013;5:S454-62.

480    34.    Bierut LJ, Stitzel JA, Wang JC, et al. Variants in nicotinic receptors and risk for

481    nicotine dependence. *Am J Psychiatry* 2008;165:1163–1171.

482    35.    Fowler CD, Lu Q, Johnson PM, Marks MJ, Kenny PJ. Habenular alpha5 nicotinic

483    receptor subunit signalling controls nicotine intake. *Nature* 2011;471:597–601.

484    36.    Thorgeirsson TE, Geller F, Sulem P, et al. A variant associated with nicotine

485    dependence, lung cancer and peripheral arterial disease. *Nature* 2008;452:638–642.

486 37.     Thorgeirsson TE, Gudbjartsson DF, Surakka I, et al. Sequence variants at

487 CHRNB3-CHRNA6 and CYP2A6 affect smoking behavior. *Nat Genet* 2010;42:448–53.

488 38.     Liu JZ, Tozzi F, Waterworth DM, et al. Meta-analysis and imputation refines the

489 association of 15q25 with smoking quantity. *Nat Genet* 2010;42:436–440.

490 39.     Tobacco and Genetics Consortium. Genome-wide meta-analyses identify multiple

491 loci associated with smoking behavior. *Nat Genet* 2010;42:441–447.

492 40.     Keskitalo K, Broms U, Heliovaara M, et al. Association of serum cotinine level

493 with a cluster of three nicotinic acetylcholine receptor genes

494 (CHRNA3/CHRNA5/CHRNB4) on chromosome 15. *Hum Mol Genet* 2009;18:4007–40.

495 41.     Munafo MR, Timofeeva MN, Morris RW, et al. Association between genetic

496 variants on chromosome 15q25 locus and objective measures of tobacco exposure. *J Natl*

497 *Cancer Inst* 2012;104:740–748.

498 42. Carr SR, Akerley W, Hashibe M, Cannon-Albright LA. Evidence for a genetical

499 contribution to non-smoking-related lung cancer. *Thorax* 2015; doi:10.1136/thoraxjnl-

500 2014-206584.

501 43.     Sampson JN, Wheeler WA, Yeager M et al. Analysis of Heritability and Shared

502 Heritability Based on Genome-Wide Association Studies for Thirteen Cancer Types.

503 *J Natl Cancer Inst*. 2015 Oct 12;107(12):djv279. doi: 10.1093/jnci/djv279. Print 2015

504 Dec.

505     **Table 1**. Characteristics of the twin cohorts included in the analyses by zygosity and sex

506     (individuals with smoking data), NorTwinCan

| Sex and zygosity of twin individuals | Denmark | Finland | Norway | Sweden | Total |
|---|---|---|---|---|---|
| Males | | | | | |
| MZ | 5,309 | 3,421 | 2,532 | 8,525 | 19,787 |
| DZ | 8,263 | 8,035 | 3,313 | 14,262 | 33,873 |
| UZ | 480 | 1,247 | - | 1,131 | 2,858 |
| All males | 14,052 | 12,703 | 5,845 | 23,918 | 56,519 |
| Females | | | | | |
| MZ | 6,570 | 3,940 | 3,074 | 10,141 | 23,725 |
| DZ | 9,525 | 8,092 | 3,788 | 16,617 | 38,022 |
| UZ | 473 | 1,049 | - | 996 | 2,518 |
| All females | 16,568 | 13,081 | 6,862 | 27,754 | 64,265 |
| **Birth cohort included** | 1870–1982 | 1880–1957 | 1915–1960 | 1886–1958 | |
| **1st Year of assessment of smoking and start of lung cancer occurrence follow-up** | 1959 | 1975 | 1980 | 1961 | |
| **End of follow-up for lung cancer occurrence** | 2010 | 2011 | 2009 | 2010 | |
| **Number of incident lung cancers** | 354 | 341 | 152 | 661 | 1508 |
| **Mean age at baseline (years)** | 49.0 | 36.2 | 38.3 | 38.9 | |
| **Mean follow-up time (years)** | 10.2* | 30.1 | 24.6 | 32.1 | |

507
508     Note: **The 5,376 twins with unknown zygosity are included in the table but are**
509     **excluded from pairwise analysis.**

510
511     *In Denmark, smoking data came from eight surveys conducted from 1959 to 2002.
512

513  **Table 2**. The numbers of pairs concordant and discordant for lung cancer at the end of follow-up by baseline pairwise smoking status
514  and zygosity.

515

| | Pairwise lung cancer status | | | | | |
|---|---|---|---|---|---|---|
| | **Monozygotic** | | | **Dizygotic** | | |
| **Baseline pairwise smoking status** | Number of Concordant Pairs | | Number of Discordant Pairs | Number of Concordant Pairs | | Number of Discordant Pairs |
| **Concordant pairs for smoking** | Neither affected | Both affected | One twin in the pair affected | Neither affected | Both affected | One twin in the pair affected |
| Never / Never | 7827 | 1 | 43 | 10709 | 0 | 59 |
| Ever / Ever | 7942 | 29 | 332 | 11474 | 28 | 527 |
| Current / Current# | 4741 | 24 | 241 | 6341 | 24 | 356 |
| **Discordant pairs for smoking** | | | | | | |
| Never / Ever | 3234 | 0 | 40 | 8177 | 0 | 173 |
| Never / Current## | 1982 | 0 | 35 | 5511 | 0 | 144 |

516

517  # Current/current pairs are a subset of ever/ever pairs

518  ## Never/current pairs are a subset of the never/ever pairs.

519

520

521 **Table 3.** Heritability estimates for lung cancer in the NorTwinCan cohort among those in the present analysis with smoking data, with

522 and without adjustment for smoking status (n=1508 cases). All estimates adjusted for country and sex.

523

| Number of complete MZ/DZ pairs | Casewise concordance rates 95% Confidence Intervals | | Adjustment for smoking | Variance component estimates 95% Confidence Intervals | | |
|---|---|---|---|---|---|---|
| | MZ | DZ | | A | C | E |
| 5299 9359 | 0.22 0.15 to 0.29 | 0.13 0.09 to 0.17 | No | 0.38 0.05 to 0.72 | 0.06 0.00 to 0.31 | 0.55 0.43 to 0.68 |
| | | | Yes | 0.34 0.00 to 0.70 | 0.02 0.00 to 0.29 | 0.64 0.50 to 0.78 |

524

525 Note: Variance components are: A: additive genetic effects, C: common environmental effects, and E: individually unique

526 environmental effects estimated from biometrical twin model taking into account censoring (see methods in the online supplement).

527
528

529 **Table 4.** Pairwise correlations in liability, heritability estimates and model fit parameters for liability to incident lung cancer among
530 ever smoking and current smoking concordant twin pairs from the NorTwinCan study. Estimates of genetic (A), shared environmental
531 (C), and unshared environmental (E) variance are presented for the ACE, AE, and CE models.
532

| Model | Correlation (95% CI) | | A Estimate (95%CI) | C Estimate (95%CI) | E Estimate (95%CI) | AIC | p-value |
|---|---|---|---|---|---|---|---|
| | MZ | DZ | | | | | |
| Ever smokers | | | | | | | |
| ACE | | | 0.28 (0.0–0.66) | 0.07 (0.0–0.36) | 0.65 (0.50–0.79) | 38759.12 | 0.01[1] |
| AE | 0.35 (0.21–0.49) | 0.21 (0.09–0.33) | 0.37 (0.25–0.49) | 0 - | 0.63 (0.51–0.75) | 38757.92 | 0.35 |
| CE | | | 0 | 0.28 (0.19–0.37) | 0.72 (0.63–0.81) | 38764.19 | 0 |
| Current smokers | | | | | | | |
| ACE | | | 0.29 (0.0–0.74) | 0.10 (0.0–0.44) | 0.62 (0.44–0.79) | 30484.27 | 0.12[1] |
| AE | 0.39 (0.20–0.55) | 0.24 (0.10–0.38) | 0.41 (0.26–0.56) | 0 - | 0.59 (0.44–0.74) | 30483.46 | 0.27 |
| CE | | | 0 | 0.31 (0.20–0.42) | 0.69 (0.58–0.80) | 30488.49 | 0.01 |

533 [1]Compared to saturated model, the other models are compared to ACE model.

534 [2] 95%CI for C effect here could not be estimated reliably

535     **Table 5**. Lung cancer in twin pairs discordant for smoking at baseline by zygosity and smoking status
536

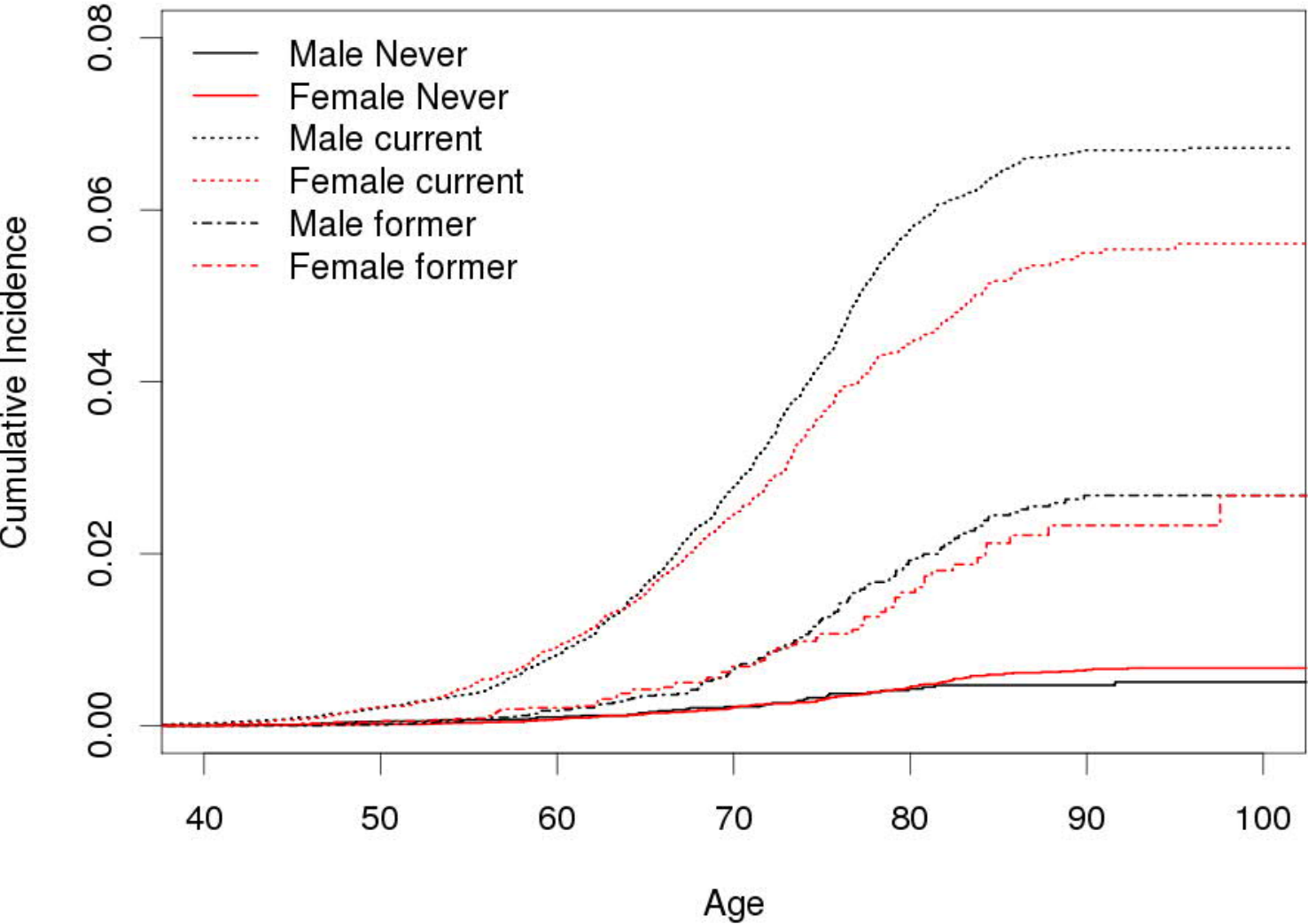| Smoking discordance | Zygosity | Pairs in which smoker had lung cancer and the non-smoking cotwin did not | Pairs in which non- smoker had lung cancer and the smoking cotwin did not | Hazard ratios (95% CI) and p-value |
|---|---|---|---|---|
| Ever/never | MZ | 35 | 5 | 5.4 ( 2.1–14.0); p=0.0005 |
| | DZ | 145 | 28 | 5.0 ( 3.2–7.9); p=1.4e-12 |
| Current/never | MZ | 31 | 4 | 6.0 (2.1-17.3) p=0.001 |
| | DZ | 124 | 20 | 5.9 (3.5-9.8) p=1.4e-11 |

**Figure legends**

**Figure 1.** Cumulative incidence of lung cancer by smoking status (never, former, current) and sex (male, female). Cumulative incidence curves are adjusted for censoring, delayed entry to cancer registration, and competing risk of death. (Continuous lines are for never smokers, dashed lines for former smokers and dotted lines for current smokers; black for males and red for females).
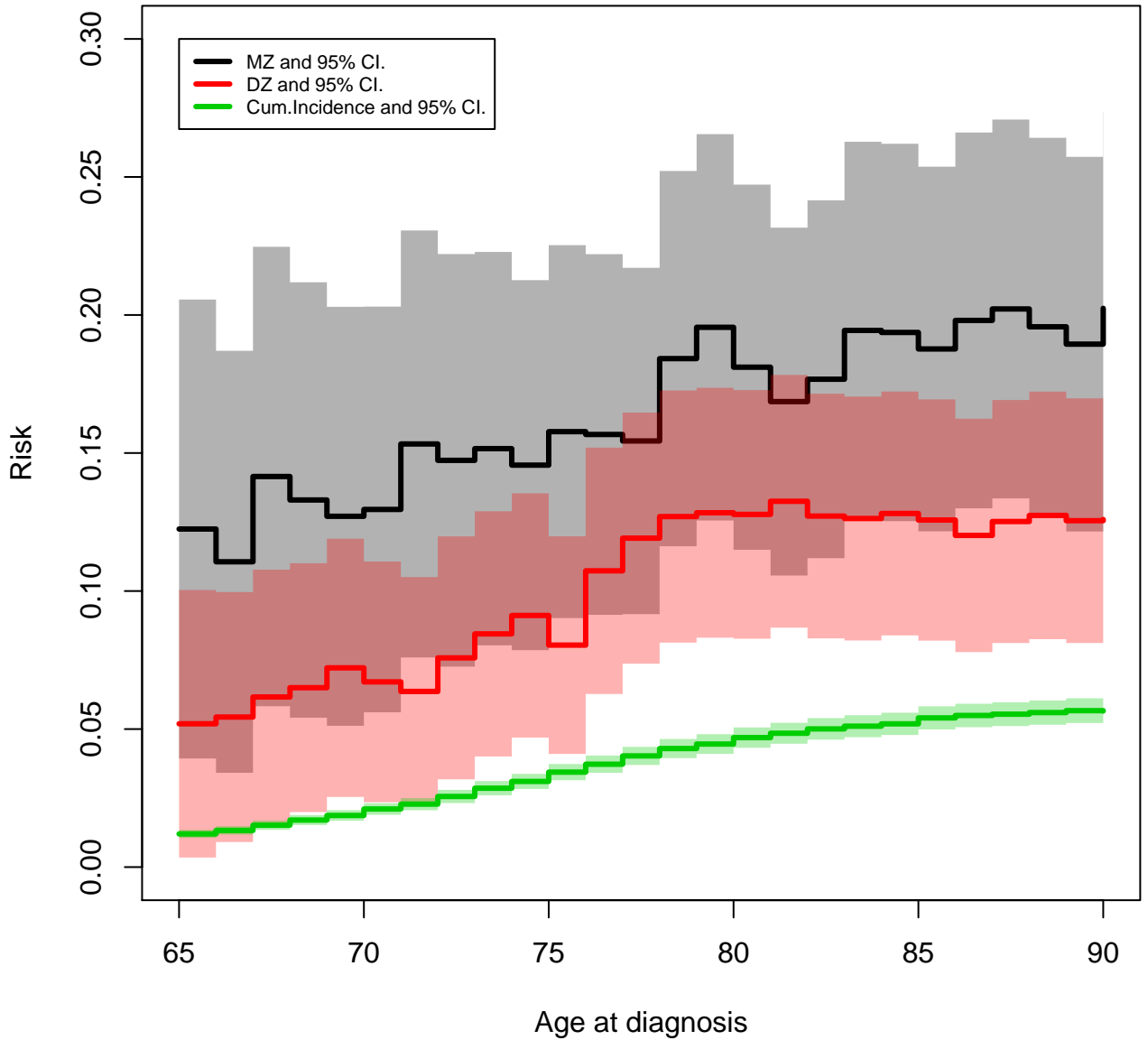
**Figure 2.** Case-wise concordance risk of lung cancer in MZ and DZ pairs compared to population risk among ever smokers, by age at diagnosis.

**Figure 3.** Relative recurrence risk ratio of lung cancer in MZ and DZ pairs compared to population risk among ever smokers, by age at diagnosis.

**Cumulative incidence by sex and smoking status**

Legend:
- Male Never
- Female Never
- Male current
- Female current
- Male former
- Female former

Y-axis: Cumulative Incidence (0.00, 0.02, 0.04, 0.06, 0.08)

X-axis: Age (40, 50, 60, 70, 80, 90, 100)

**Ever Smokers**

**Ever Smokers**