

## Little Evidence of Modified Genetic Effect of rs16969968 on Heavy Smoking Based on Age of Onset of Smoking

Christine Adjangba<sup>1‡</sup>, Richard Border, Ph.D.<sup>1,2</sup>, Pamela N. Romero Villela<sup>1,3</sup>, Marissa A. Ehringer, Ph.D.<sup>1,4</sup>, and Luke M. Evans, Ph.D.<sup>1,5\*</sup>

<sup>1</sup> Institute for Behavioral Genetics, University of Colorado, Boulder

<sup>2</sup> Department of Applied Mathematics, University of Colorado, Boulder

<sup>3</sup> Department of Psychology & Neuroscience, University of Colorado, Boulder

<sup>4</sup> Department of Integrative Physiology, University of Colorado, Boulder

<sup>5</sup> Department of Ecology & Evolutionary Biology, University of Colorado, Boulder

‡ Current address: Department of Biology, Duke University

\* Corresponding Author: Luke M. Evans, Ph.D., Institute for Behavioral Genetics, Boulder, CO 80303; [luke.m.evans@colorado.edu](mailto:luke.m.evans@colorado.edu)

**Word Count: 2,712**

## KEY POINTS

**Question:** Does the age of regular smoking initiation modify the effect of rs16969968, the strongest smoking-related association, as previously reported?

**Findings:** Limited support for an age of smoking $\times$ rs16969968 interaction. Our replication attempt of previous work failed. We found nominally (none genome-wide) significant differences in allele effect sizes, and only with certain combinations of phenotype scales, contrasting previous findings.

**Meaning:** Individual interactions at specific loci are likely very small, and detection will require biobanks many-fold larger than current samples. Furthermore, measurement and analysis scales can strongly impact whether or not an interaction is detected and resulting conclusions.

## ABSTRACT

**Importance:** Tobacco smoking is the leading cause of preventable death globally. Smoking quantity, measured in cigarettes per day (CPD), is influenced both by the age of onset of regular smoking (AOS) and by genetic factors, including a strong association with the non-synonymous single nucleotide polymorphism rs16969968. A previous study<sup>1</sup> reported an interaction between these two factors, whereby individuals carrying the risk allele for rs16969968 who started smoking earlier showed increased risk for heavy smoking compared to those who started later. This finding has yet to be replicated in a large, independent sample.

**Objective:** To directly replicate the previous finding, explore the influence of phenotypic and analysis scale on the results, and to assess statistical power to detect gene-moderator interactions.

**Design:** We performed a preregistered, direct replication attempt of the rs16969968 $\times$ AOS interaction on smoking quantity in 116,317 unrelated individuals from the UK Biobank. We fit statistical association models mirroring the original publication as well as formal interaction tests on multiple phenotypic and analytical scales.

**Setting:** We used self-report smoking behavior data from the UK Biobank.

**Participants:** 116,317 unrelated individuals of European ancestry with reported smoking.

**Exposures:** rs16969968 genotype (0, 1 or 2 risk alleles), and AOS, encoded as untransformed, binned, and binary early/late.

**Main Outcome and Measures:** Self-reported CPD measured untransformed,  $\log_{10}$ -transformed, binned and binary heavy/light smoking variables.

**Results:** We replicated the main effects of rs16969968 and AOS on CPD. We failed to replicate the interaction using previous methods. Nominal significance of the rs16969968×AOS interaction term depended strongly on the scales of the analysis and the phenotypes, as did associations stratified by early/late AOS. In no analysis did the interaction tests pass genome-wide correction ( $\alpha=5e-8$ ), and in all analyses, the estimated interaction effect size was smaller in magnitude than previous estimates.

**Conclusions and Relevance:** We failed to replicate the strong rs16969968×AOS interaction effect on smoking quantity previously reported. If such gene-moderator interactions influence complex traits, current biobanks lack the power to detect significant genome-wide associations given the minute effect sizes expected. Furthermore, many potential interaction effects are likely to depend on the scale of measurement employed.

## INTRODUCTION

Approximately 20% of deaths every year in the United States can be attributed to cigarette smoking, and smokers have life expectancies at least 10 years shorter than nonsmokers<sup>2</sup>. Furthermore, the rise in use among adolescents of various electronic cigarettes has emerged as a potentially dangerous trend about which little is known regarding long-term health and addiction consequences<sup>3</sup>. There is strong evidence from adoption, family, and twin studies that both genetic and environmental factors contribute to risk for smoking behaviors, with heritability estimates for nicotine dependence, ever becoming a regular smoker, and smoking quantity ranging between 33% and 71%<sup>4-8</sup>. Recently, genome-wide association studies (GWAS) have identified common variants associated with smoking behaviors<sup>1,9-13</sup>. In particular, the nicotinic acetylcholine receptor subunit genes *CHRNA5-CHRNA3-CHRNA4* on chromosome 15 have been identified in well-powered GWAS of smoking behaviors<sup>14,15</sup>. Within *CHRNA5*, which codes for the  $\alpha 5$  receptor subunit, the nonsynonymous G/A single nucleotide polymorphism (SNP) rs16969968 has been replicated through both large-scale GWAS<sup>1,9-13,16</sup> and functional assays<sup>17-20</sup> to influence smoking quantity, as measured by the number of cigarettes smoked per day, and nicotine dependence. The rs16969968-A risk allele has the largest estimated allelic effect on smoking quantity known to date<sup>13</sup>.

In addition to genetic risk factors for heavy smoking, earlier age at onset of regular smoking (AOS) is well-known to increase risk for later heavy use and nicotine dependence<sup>21,22</sup>. In light of previous findings, Hartz et al.<sup>1</sup> conducted a meta-analysis of 33,348 individuals across 43 European and American data sets to test whether genetic vulnerability to heavy smoking and nicotine dependence at rs16969968 depends on AOS, and found a strong, significant interaction between early AOS and the rs16969968-A allele (OR=1.16). However, despite the large rs16969968×AOS interaction effect size reported and animal model evidence to support the plausibility of such an interaction<sup>20</sup>, we are aware of no large-scale replication attempt in an independent sample. Here, we assessed whether there is an rs16969968×age of onset of smoking interaction in a well-powered (Fig. 1), independent sample, in an attempt to directly replicate the original study.

## METHODS

We preregistered our analyses through the Open Science Framework ([osf.io/ynh2j](https://osf.io/ynh2j)) after we had obtained the UK Biobank data, but before we had analyzed CPD or AOS.

## *Study Population*

We used the UK Biobank, a large sample of approximately 500,000 participants with rich phenotype and genome-wide genotype data<sup>23</sup>. We included all participants with available genomic data who had reported CPD and AOS data. The participants were either current or former smokers aged 40 years or older. To avoid confounding influences of population stratification<sup>24</sup>, we performed analyses using individuals of European ancestry, the largest subsample within the UK Biobank, identified by those whose first scores on the first four principal components (PCs; UK Biobank data field ID 22009) fell within the range of the UK Biobank identified individuals of European ancestry (field 22006).

We only included unrelated individuals in analyses to avoid possible confounding due to shared environmental factors. Relatedness was estimated using MAF- and LD-pruned array markers (plink2<sup>25</sup> command: `--maf 0.01 --hwe 1e-8 --indep-pairwise 50 5 0.2`) after excluding those individuals with self-report and genetic sex mismatch (fields 31 and 22001), those with unusually high inbreeding coefficients ( $|F_{het}| > 0.2$ ), and those identified by the UK Biobank and Affymetrix as having poor-quality genomic data (fields 220010 and 22051). Unrelated individuals (estimated relatedness  $< 0.05$ ) were identified with GCTA<sup>26</sup> v1.91.3. After removing individuals with missing phenotype and covariate data (see below), a total of 116,317 unrelated individuals of European ancestry were included.

## *Variables*

Smoking quantity, as measured by CPD, was the primary dependent variable in analyses. Data on CPD (fields 2887, 3456, and 6183) were obtained from current or former smokers by asking the question “About how many cigarettes do/did you smoke on average each day?” These data were highly skewed; therefore, we added an analysis of  $\log_{10}$ -transformed CPD (Supplemental Fig. 1). Because of possible differences in the interaction on the additive compared to multiplicative scale<sup>27</sup> and observed evidence of scale dependence (see results below), we also added an analysis of heavy/light CPD analyzed on an additive analytical scale. These two additional analyses were the only deviations from our preregistered analyses. Final analyses used untransformed CPD,  $\log_{10}$ (CPD), heavy/light (analyzed on both multiplicative, i.e., logistic, and additive scales), and binned encodings. The dichotomous encoding defined smoking quantity as light smoking (CPD  $\leq 10$ ) versus heavy smoking (CPD  $> 20$ ), as in Hartz et al. The binned encoding defined smoking

quantity as a linear variable consisting of 0 (CPD  $\leq$  10), 1 (11-20 CPD), 2 (21-30 CPD), or 3 (CPD  $>$  30), also matching their secondary analysis.

Age of onset of regular smoking (AOS) was determined from fields 3426 and 2867, where participants were asked “How old were you when you first started smoking on most days?” AOS was analyzed based on a dichotomous encoding, a binned encoding, and the raw AOS data, replicating the methods of Hartz et al. (2012). The dichotomous encoding defined early as AOS  $\leq$  16 years and late as AOS  $>$  16 years. The binned encodings were 0 (AOS  $\leq$  15 years), 1 (AOS = 16 years), 2 (17-18 years AOS), or 3 (AOS  $>$  18 years).

Covariates included were sex (field 31), age at time of assessment (field 21003), age<sup>2</sup>, Townsend Deprivation Index (field 21003), educational attainment (“qualification”, categorical, field 6138), genotyping batch (field 22000), assessment center (field 54), and the first 10 genetic principal components as estimated with flashpca<sup>28</sup> applied to the MAF- and LD-pruned SNPs as described above. High collinearity of covariates within this sample resulted in a rank-deficient design matrix, which we addressed by performing a principal components analysis of the  $c=141$  fixed effects using the `prcomp` function in R v3.2.2<sup>29</sup>. We then estimated the rank of the resulting eigenvector matrix (rank  $r < c$ ) using the `matrix` R package<sup>30</sup> and included the first  $r=140$  principal components of the design matrix as covariates in all analyses.

### Statistical Analyses

For dichotomized light/heavy CPD, we performed logistic regression using `glm` (family='binomial') in R<sup>29</sup> to assess the multiplicative scale interaction. The model included the rs16969968 genotype (coded as 0, 1, or 2), AOS, and rs16969968×AOS. All genotype×covariate and AOS×covariate interactions were included within the models to appropriately control for confounding<sup>31</sup>. For continuous variables (binned and un- & log-transformed CPD) and the additive scale interaction model of the dichotomous heavy/light phenotype, we tested the same model using linear regression with the R `lm` function.

The above model varied from that tested by Hartz et al., who tested rs16969968 effects on smoking phenotypes stratified by AOS (early versus late), using logistic regression (i.e., multiplicative scale). To recapitulate their methods, we secondarily performed association tests of rs16969968 stratified by early versus late AOS using BOLT-LMM v2.3.2<sup>32</sup>, with 339,444 genome-wide SNPs (quality control as described above, but

without LD-pruning) to control for background polygenicity and cryptic relatedness. All covariates were included in the BOLT-LMM models, excluding interaction terms. Finally, to directly replicate previous methods, we performed AOS-stratified logistic regression of heavy/light CPD using only rs16969968 and sex as independent variables.

We also performed several power analyses, to determine the power to detect the previously reported effect size<sup>1</sup>, as well as to determine the sample size needed to achieve 80% power at specified effect sizes and  $\alpha$ . To estimate the power to detect the previously reported effect size in the UK Biobank sample under a multiplicative scale interaction model, we simulated 54,827 diploid genotypes and early/late AOS in R, with linear phenotypes simulated using the previously reported main effect sizes as,

$$lp = 0.33 + \log(1.28)g + \log(2.63)a + \log(OR_{AOS*rs16969968})ag \quad (1)$$

where genotypes,  $g$ , were simulated from a binomial distribution with MAF=0.34, the observed frequency of the A allele in the UK Biobank, early versus late AOS status,  $a$ , was randomly assigned to individuals. We varied the interaction effect size,  $\log(OR_{AOS*rs16969968})$  between 0.005 and 0.4, reflecting a range of plausible effect sizes and encompassing the previously reported interaction effect (OR=1.16). Binary phenotypes,  $y$ , were then simulated in R as,

$$y = \text{rbinom}(54827, 1, \exp(lp)/(1 + \exp(lp))). \quad (2)$$

For each simulated interaction effect size, we performed 1,000 replicate simulations, estimating the interaction effect using logistic regression as above, and recorded the number of observations with an interaction  $p$ -value below either nominal significance,  $\alpha=0.05$ , or genome-wide significance,  $\alpha=5e-8$ . We performed similar simulations with the main AOS and rs16969968 effect sizes estimated within the UK Biobank (see below). We varied the sample size from 1,000 to 2e6, varying interaction effect sizes (previously reported  $OR_{AOS*rs16969968}=1.16$  versus our estimate  $OR_{AOS*rs16969968}=1.03$ ), and nominal versus genome-wide significance thresholds ( $\alpha=0.05$  versus  $5e-8$ , respectively).

## RESULTS

We observed significant main effects of the rs16969968 A allele and AOS on CPD (Figures 2A, 2B; Table 1). When estimated as predictors of heavy vs. light smoker status, the estimated genetic effect,  $OR_{rs16969968}=1.38$  ( $p=4.97e-48$ ), was similar to the previous estimate<sup>1</sup> of 1.28. The effect of early AOS,

$OR_{AOS}=1.58$  ( $p=1.59e-51$ ), was less than previously reported<sup>1</sup> ( $OR_{AOS}=2.63$ ). However, both main effects were associated with CPD in the expected direction, regardless of the CPD or AOS encoding, and represent strong evidence that both the rs16969968 A allele and early AOS are positively associated with heavier smoking.

Conversely, the interaction between rs16969968 genotype and AOS was only nominally significant ( $\alpha=0.05$ ) and only in some combinations of CPD and AOS encoding (Figure 2C; Table 1). Specifically, when treating both CPD and AOS as binary phenotypes the logistic model interaction was not significant ( $OR_{rs16969968 \times AOS}=1.03$ ,  $p=0.36$ ) and the effect was notably lower than the previously reported estimate of 1.16. Interestingly, the interaction effect was nominally significant ( $p<0.05$ ) for the binned CPD phenotype and dichotomized AOS, and when heavy/light CPD was analyzed on the additive scale, but not when the CPD phenotype was either heavy vs. light analyzed on the multiplicative scale model or when CPD was log-transformed. Across all tests and all CPD and AOS encodings, interaction effects did not reach genome-wide significance ( $p>0.02$ ).

Associations of rs16969968 stratified by AOS also produced mixed results. 95% confidence intervals ( $\alpha=0.05$ ) were non-overlapping for all CPD encodings except  $\log_{10}(\text{CPD})$  phenotypes (Fig. 3, Table S1). When examining heavy versus light CPD,  $OR_{\text{Early}}/OR_{\text{Late}}=1.076$ , roughly half of that previously reported (Table S1). All confidence intervals overlapped when using a genome-wide significance threshold ( $\alpha=5e-8$ ). The direct replication test using Hartz et al.<sup>1</sup> methods with only rs16969968 and sex as independent variables found no evidence of different allelic effects between early and late smokers ( $p=0.16$ ; Table S2).

Our power analyses yielded two main results. First, our sample was well powered (>99%) to detect an interaction effect of the size previously reported at nominal significance (Figs. 1,4), though not at genome-wide significance (power ~20%), even with over 54,000 subjects. Second, a sample roughly seven times larger than that analyzed here would be required to detect an interaction effect of  $OR_{rs16969968 \times AOS}=1.03$ , as estimated within our sample, with 80% power at  $\alpha=0.05$ . Achieving 80% power at a genome-wide threshold ( $\alpha=5e-8$ ) would require a sample of approximately 2,000,000 participants (Fig. 4).

## DISCUSSION

We replicated the substantial main effects of rs16969968 and early age of onset of smoking on CPD, across all phenotypic and analytical scales (Table 1). Estimates were in the same direction and of roughly the



same magnitude as previously reported<sup>1</sup>.

Conversely, we found limited evidence of an rs16969968×AOS interaction effect. Formal interaction model results were mixed and depended heavily on measurement scale and phenotype encoding. Notably, our attempt to directly replicate the methods of Hartz et al. failed to identify a significant difference in the rs16969968-A allele effect on heavy smoking between early and late AOS ( $p=0.16$ ; Table S2). This is in contrast to the results from stratified linear mixed model analyses, where the genetic effect in early AOS individuals was 1.076-fold higher than in late AOS individuals, perhaps reflecting greater statistical power of linear mixed models<sup>33</sup>, as well as more control of potential confounding variables, such as genetic ancestry and geographic variation throughout the UK. However, this effect was also scale-dependent, and did not persist when CPD was log-transformed, limiting our confidence in the interaction. Across multiple analytic frameworks and phenotype encodings, the majority of our results were incongruent with an interaction between rs16969968 and AOS.

#### *Magnitude of Effects and Power*

No interaction test, and no comparison of stratified estimates, reached genome-wide significance ( $\alpha=5e-8$ ) despite the comparatively large sample size of our study. With genome-wide genotyping arrays and imputation commonly applied<sup>34</sup>, and as genome-wide interaction associations and heritability studies have become more frequent<sup>35-41</sup>, focusing on genome-wide significance thresholds is paramount to avoid false positives, even in situations where there are *a priori* hypotheses of interaction, as in rs16969968×AOS.

In all tests related to interaction effects and stratified associations, the estimated interaction effect sizes were much smaller than previously reported<sup>1</sup> ( $OR_{AOS \times rs16969968}=1.16$ ). We estimated the effect to be only 1.03 (or 1.07 in the stratified associations). It is important to recognize that both main effects were strong, significant, and in the expected direction, reflecting the strongest single-locus genetic effect on CPD<sup>13</sup> and a strong, consistent risk factor of heavy smoking (early age of initiation). This suggests that if such interactions were to exist, their effect would be much less than previously expected and would contribute only minimally to phenotypic variance.

The discrepancy between our results and those reported by Hartz et al.<sup>1</sup> could reflect differences between the study populations and models used for analyses. The study by Hartz et al.<sup>1</sup> exemplified a

tremendous effort to collect the largest available sample size at the time. They were able to do so by meta-analyzing multiple individual studies together, an effort that must be recognized and applauded. One possible outcome of this approach is heterogeneity of effect estimates, which they found and noted. Our analysis focused on a single, relatively homogeneous dataset instead of many studies, removing potential heterogeneity that could have influenced the previous results. However, although Hartz et al. included some European datasets, consistent cultural differences may exist between American and UK samples, such as general attitudes towards smoking, which cannot be ruled out as leading to different findings. Additional methodological differences include testing a full statistical interaction model with complete covariate $\times$ AOS and covariate $\times$ genotype terms and using a linear mixed model in our stratified analyses, neither of which were previously employed. Mixed model approaches generally improve power<sup>33</sup>, and including the covariate interaction terms should lead to unbiased estimates of the rs16969968 $\times$ AOS interaction<sup>31</sup>. On the other hand, comparing estimates across different subsamples, as in stratified linear mixed model analysis, introduces an additional potential source of confounding. However, the respective strengths and weaknesses of these methods cannot account for our failure to directly replicate the original finding; our stratified association tests with only sex as a covariate (mirroring the approach of Hartz et al.) failed to identify significant differences in allelic effect sizes between early and late AOS individuals ( $p=0.16$ ; Table S2), despite being well-powered to do so.

Regardless, with respect to particular phenotype encodings and analyses (e.g., stratified analyses of heavy vs. light smoker status, with linear mixed models), we did find nominally significant small differences in allelic effect size estimates between early- and late-onset smokers. These findings are thus potentially congruent with a small interaction between rs1696968 and AOS. If there is a true rs16969968 $\times$ AOS interaction of roughly the magnitude we estimated ( $OR=1.03$ ), it would require seven times the sample size to detect it with nominal significance, and roughly 20 times the sample size to detect it with a genome-wide significance threshold (Fig. 4). We must therefore conclude that such interactions are likely of very small effect at individual loci, and will be very difficult to identify, even with the largest available biobanks.

### *Conclusions*

We found limited support for the rs16969968 $\times$ AOS interaction. To the extent that AOS might moderate the

effect of rs16969968, we estimate this effect to be far smaller than previously reported. We suggest that even larger sample sizes will be required to identify, with genome-wide significance, interactions at individual loci given the expected magnitude of the interaction effects. On the other hand, our unambiguous replications of the main effects of both rs16969968 and AOS on smoking quantity support epidemiological evidence that individuals who begin regularly smoking at a young age are at a higher risk for nicotine dependence later in life<sup>21,22</sup>. This provides further evidence in support of public health interventions for adolescent smoking that could help reduce tobacco use, which would in turn lower the number of tobacco-related deaths and illnesses.

## ARTICLE INFORMATION

**Corresponding Author:** Luke M. Evans, Ph.D., Institute for Behavioral Genetics, University of Colorado Boulder, 1430 30th St., Boulder, CO 80303; ([luke.m.evans@colorado.edu](mailto:luke.m.evans@colorado.edu))

### Author Contributions:

*Concept and design:* Evans, Ehringer, Adjangba

*Statistical analysis:* Adjangba, Evans, Border, Romero

*Drafting of the manuscript:* Adjangba, Evans

*Critical revision of the manuscript for important intellectual content:* All authors

*Supervision:* Evans

**Conflict of Interest Disclosures:** The authors declare no conflict of interest.

**Funding/Support:** Ms. Adjangba was supported by the Summer Multicultural Access to Research Training Program at the University of Colorado. Drs. Evans and Border are supported by National Institute of Mental Health R01 MH100141-06 (PI: Matthew C. Keller) and Dr. Evans is supported by National Institute on Drug Abuse R01 DA044283-01A1 (PI: Scott I. Vrieze).

**Role of the Funder/Sponsor:** The funding sources had no role in the design and conduct of the study; collection, management, analysis, and interpretation of the data; preparation, review, or approval of the manuscript; and decision to submit the manuscript for publication.

**Additional Contributions:** We thank Peter Visscher, Ph.D. (University of Queensland) and Matthew C. Keller, Ph.D. (University of Colorado) for helpful discussion.

## REFERENCES

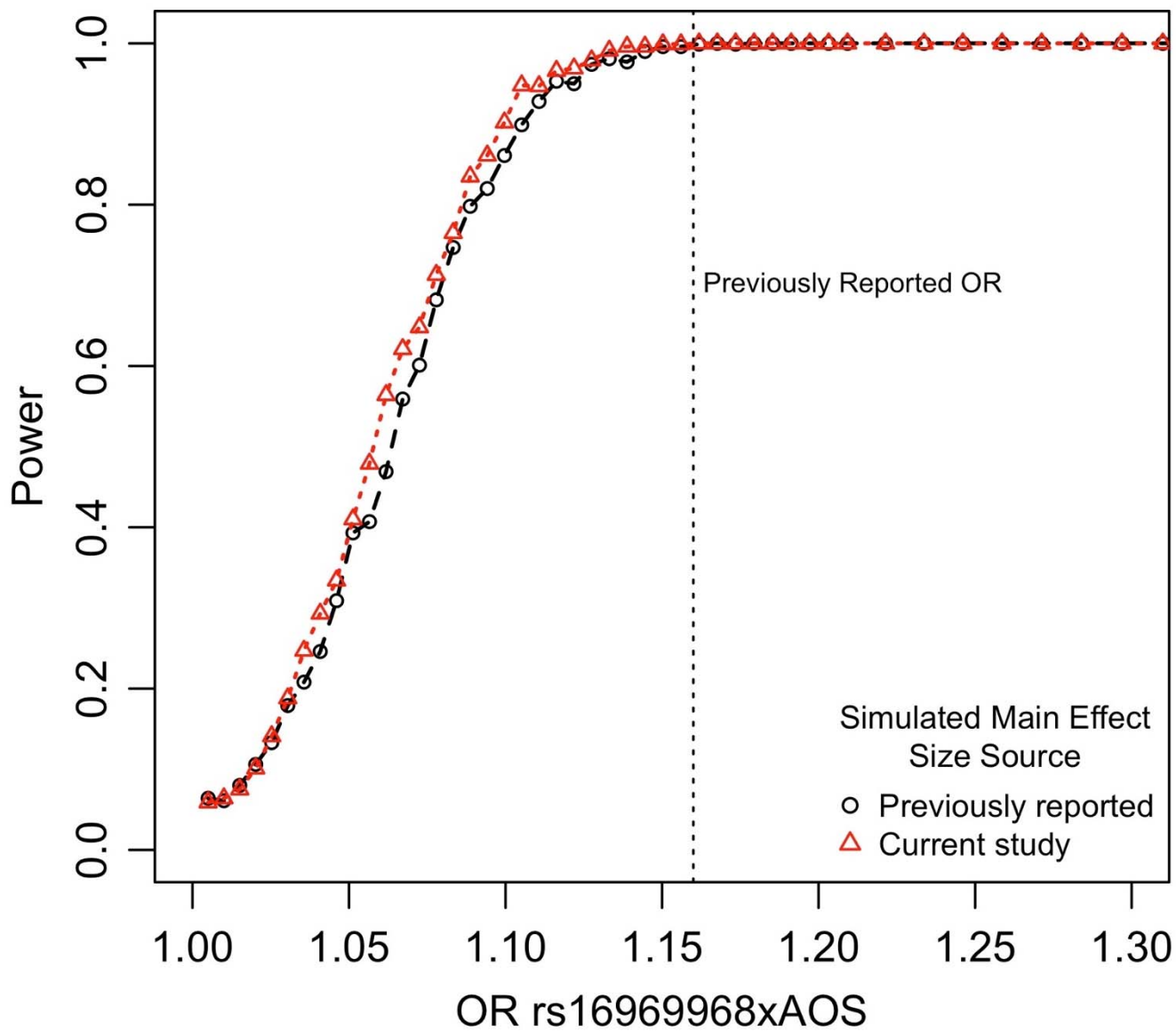
1. Hartz SM, Short SE, Saccone NL, et al. Increased genetic vulnerability to smoking at CHRNA5 in early-onset smokers. *Arch Gen Psychiatry*. 2012;69(8):854-860.
2. US Department of Health and Human Services. Health Consequences of Smoking—50 Years of Progress A Report of the Surgeon General. *Report of the Surgeon general*. 2014;1081.
3. Fadus MC, Smith TT, Squeglia LM. The rise of e-cigarettes, pod mod devices, and JUUL among youth: Factors influencing use, health implications, and downstream effects. *Drug Alcohol Depend*. 2019;201:85-93.
4. Haberstick BC, Ehringer MA, Lessem JM, Hopfer CJ, Hewitt JK. Dizziness and the genetic influences on subjective experiences to initial cigarette use. *Addiction*. 2011;106(2):391-399.
5. Haberstick BC, Zeiger JS, Corley RP, et al. Common and drug-specific genetic influences on subjective effects to alcohol, tobacco and marijuana use. *Addiction*. 2011;106(1):215-224.
6. Kaprio J. Genetic epidemiology of smoking behavior and nicotine dependence. *COPD*. 2009;6(4):304-306.
7. Rose R.J., Broms U., Korhonen T., Dick D.M., J. K. Genetics of Smoking Behavior. In: YK K, ed. *Handbook of Behavior Genetics*. New York, NY: Springer; 2009.
8. Kendler KS, Schmitt E, Aggen SH, Prescott CA, Virginia V. Genetic and Environmental Influences on Alcohol, Caffeine, Cannabis, and Nicotine Use From Early Adolscence to Middle Adulthood. *Arch Gen Psychiatry*. 2008;65:674-682.
9. Hancock DB, Guo Y, Reginsson GW, et al. Genome-wide association study across European and African American ancestries identifies a SNP in DNMT3B contributing to nicotine dependence. *Mol Psychiatry*. 2017.
10. Hancock DB, Reginsson GW, Gaddis NC, et al. Genome-wide meta-analysis reveals common splice site acceptor variant in CHRNA4 associated with nicotine dependence. *Transl Psychiatry*. 2015;5:e651.
11. Saccone NL, Emery LS, Sofer T, et al. Genome-Wide Association Study of Heavy Smoking and Daily/Nondaily Smoking in the Hispanic Community Health Study/Study of Latinos (HCHS/SOL). *Nicotine Tob Res*. 2018;20(4):448-457.

12. Wen L, Yang Z, Cui W, Li MD. Crucial roles of the CHRN3-CHRNA6 gene cluster on chromosome 8 in nicotine dependence: update and subjects for future research. *Transl Psychiatry*. 2016;6(6):e843.
13. Liu M, Jiang Y, Wedow R, et al. Association studies of up to 1.2 million individuals yield new insights into the genetic etiology of tobacco and alcohol use. *Nat Genet*. 2019.
14. Bierut LJ, Stitzel JA, Wang JC, et al. Variants in nicotinic receptors and risk for nicotine dependence. *Am J Psychiatry*. 2008;165(9):1163-1171.
15. Berrettini W, Yuan X, Tozzi F, et al. Alpha-5/alpha-3 nicotinic receptor subunit alleles increase risk for heavy smoking. *Mol Psychiatry*. 2008;13(4):368-373.
16. Tobacco, Genetics C. Genome-wide meta-analyses identify multiple loci associated with smoking behavior. *Nat Genet*. 2010;42(5):441-447.
17. Bailey CD, Tian MK, Kang L, O'Reilly R, Lambe EK. Chrna5 genotype determines the long-lasting effects of developmental in vivo nicotine exposure on prefrontal attention circuitry. *Neuropharmacology*. 2014;77:145-155.
18. Kuryatov A, Berrettini W, Lindstrom J. Acetylcholine receptor (AChR) alpha5 subunit variant associated with risk for nicotine dependence and lung cancer reduces (alpha4beta2)(2)alpha5 AChR function. *Mol Pharmacol*. 2011;79(1):119-125.
19. George AA, Lucero LM, Damaj MI, Lukas RJ, Chen X, Whiteaker P. Function of human alpha3beta4alpha5 nicotinic acetylcholine receptors is reduced by the alpha5(D398N) variant. *J Biol Chem*. 2012;287(30):25151-25162.
20. O'Neill HC, Wageman CR, Sherman SE, Grady SR, Marks MJ, Stitzel JA. The interaction of the Chrna5 D398N variant with developmental nicotine exposure. *Genes Brain Behav*. 2018;17(7):e12474.
21. Lydon DM, Wilson SJ, Child A, Geier CF. Adolescent brain maturation and smoking: what we know and where we're headed. *Neurosci Biobehav Rev*. 2014;45:323-342.
22. Kendler KS, Myers J, Damaj MI, Chen X. Early smoking onset and risk for subsequent nicotine dependence: a monozygotic co-twin control study. *Am J Psychiatry*. 2013;170(4):408-413.
23. Bycroft C, Freeman C, Petkova D, et al. The UK Biobank resource with deep phenotyping and genomic data. *Nature*. 2018;562(7726):203-209.

24. Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet.* 2006;38(8):904-909.
25. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience.* 2015;4:7.
26. Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide complex trait analysis. *Am J Hum Genet.* 2011;88(1):76-82.
27. VanderWeele TJ, Knol MJ. A Tutorial on Interaction. *Epidemiologic Methods.* 2014;3(1):33-72.
28. Abraham G, Inouye M. Fast principal component analysis of large-scale genome-wide data. *PLoS One.* 2014;9(4):e93766.
29. *R: A language and environment for statistical computing.* [computer program]. Vienna, Austria: R Foundation for Statistical Computing; 2015.
30. *Matrix: Sparse and dense matrix classes and methods. R package version 1.2-2.* [computer program]. 2015.
31. Keller MC. Gene x environment interaction studies have not properly controlled for potential confounders: the problem and the (simple) solution. *Biol Psychiatry.* 2014;75(1):18-24.
32. Loh PR, Kichaev G, Gazal S, Schoech AP, Price AL. Mixed-model association for biobank-scale datasets. *Nat Genet.* 2018;50(7):906-908.
33. Yang J, Zaitlen NA, Goddard ME, Visscher PM, Price AL. Advantages and pitfalls in the application of mixed-model association methods. *Nat Genet.* 2014;46(2):100-106.
34. McCarthy S, Das S, Kretschmar W, et al. A reference panel of 64,976 haplotypes for genotype imputation. *Nat Genet.* 2016;48(10):1279-1283.
35. Rawlik K, Canela-Xandri O, Tenesa A. Evidence for sex-specific genetic architectures across a spectrum of human complex traits. *Genome Biol.* 2016;17(1):166.
36. Young AI, Wauthier FL, Donnelly P. Identifying loci affecting trait variability and detecting interactions in genome-wide association studies. *Nat Genet.* 2018;50(11):1608-1614.
37. Dahl A, Nguyen K, Cai N, Gandal MJ, Flint J, Zaitlen N. A Robust Method Uncovers Significant Context-Specific Heritability in Diverse Complex Traits. *Am J Hum Genet.* 2020;106(1):71-91.

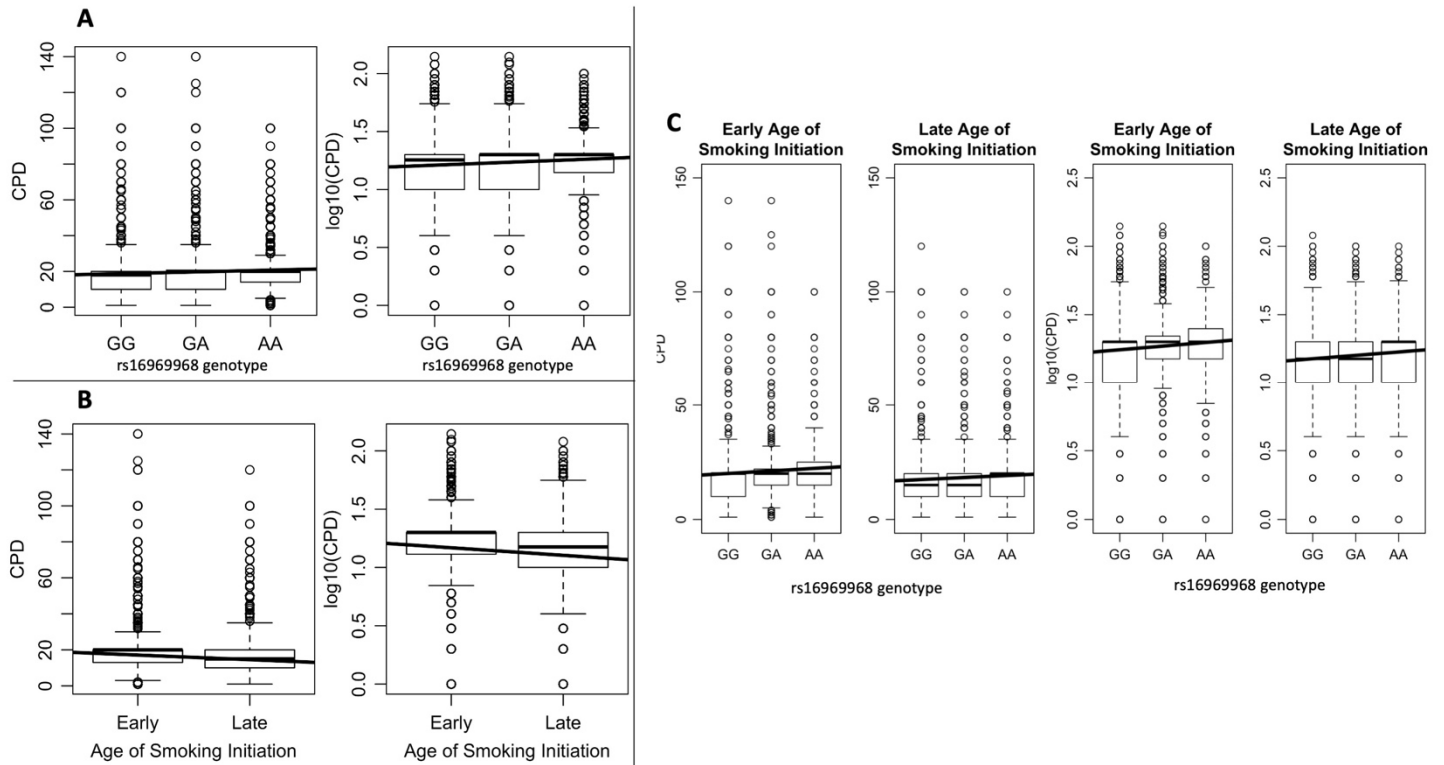
38. Peterson RE, Cai N, Dahl AW, et al. Molecular Genetic Analysis Subdivided by Adversity Exposure Suggests Etiologic Heterogeneity in Major Depression. *Am J Psychiatry*. 2018;175(6):545-554.
39. Arnau-Soler A, Adams MJ, Generation S, Major Depressive Disorder Working Group of the Psychiatric Genomics C, Hayward C, Thomson PA. Genome-wide interaction study of a proxy for stress-sensitivity and its prediction of major depressive disorder. *PLoS One*. 2018;13(12):e0209160.
40. Nivard MG, Middeldorp CM, Lubke G, et al. Detection of gene–environment interaction in pedigree data using genome-wide genotypes. *European Journal of Human Genetics*. 2016;24(12):1803-1809.
41. Robinson MR, English G, Moser G, et al. Genotype-covariate interaction effects and the heritability of adult body mass index. *Nat Genet*. 2017;49(8):1174-1181.

**Figure 1.** Power ( $\alpha=0.05$ ) to detect the rs16969968xAOS interaction effect, given the main effect estimates previously reported or estimated in the current study.

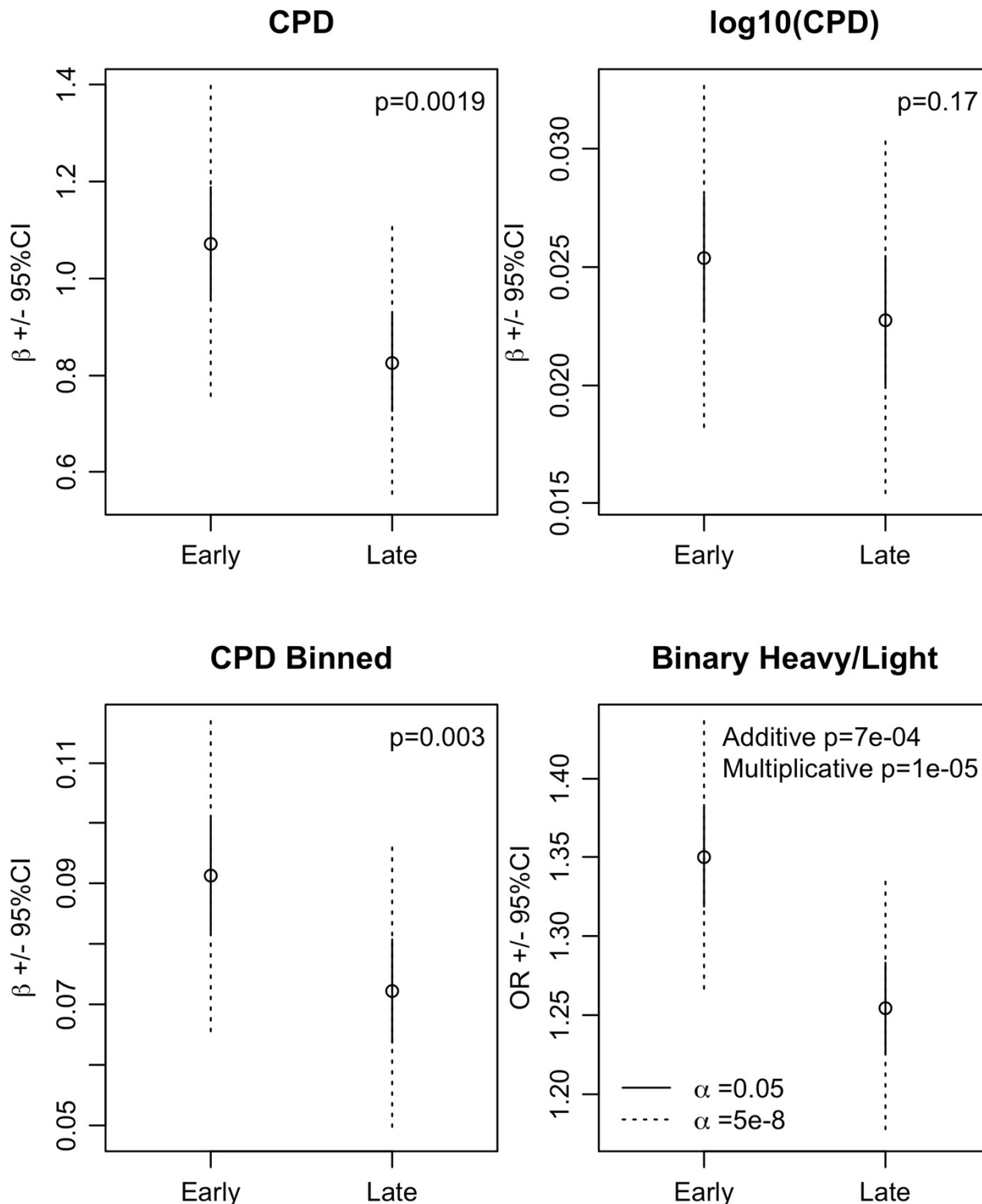




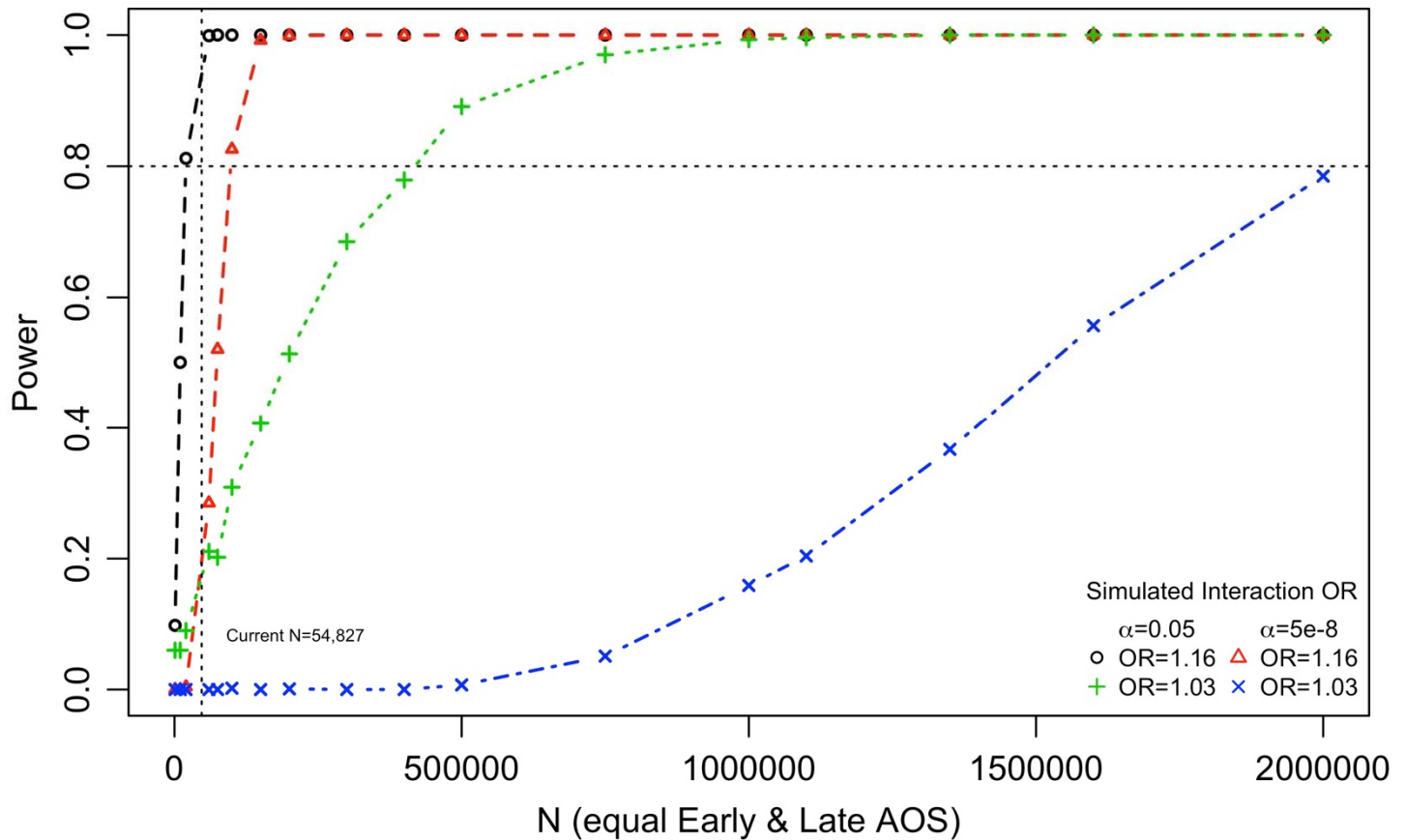
**Figure 2.** Main effects of rs16969968 genotype (A) and AOS (B) on cigarettes per day (CPD or  $\log_{10}(\text{CPD})$ ), and rs16969968 effect on CPD and  $\log_{10}(\text{CPD})$  as a function of early or late age of initiation (C).



**Figure 3.** Allelic effect size estimates,  $\beta$ , of the rs16969968 risk allele, A, estimated in stratified association analyses by early and late age of initiation, with CIs indicated using either  $\alpha=0.05$  (solid line) or genome-wide Bonferroni-corrected  $\alpha=5e-8$  (dashed line). For the heavy vs. light smoker phenotype, allelic effects ( $\beta$ ) were transformed to OR using the BOLT-LMM<sup>32</sup>-suggested transformation, of  $e^{\beta / (u(1-u))}$ , where  $u=0.42$  is the proportion of cases. Effect size difference Z-test p-values are shown for each comparison. See Table S1 for estimates and statistics.



**Figure 4.** Power to detect the statistical interaction effect across a range of sample sizes for the interaction effect size estimated previously (OR=1.16) and from the current study (OR=1.03), applying either a nominal  $\alpha=0.05$  or genome-wide Bonferroni-corrected  $\alpha=5e-8$ .



**Table 1.** Estimated main and interaction effects for rs16969968, age of smoking initiation (AOS), and their interaction. Shown are estimates for each encoding of CPD and AOS.

CPD Encoding	N	AOS Encoding	Intercept		rs16969968 (risk allele=A)				AOS				rs16969968 x AOS			
			Estimate	SE	$\beta^a$	SE	Stat <sup>b</sup>	p	$\beta^a$	SE	Stat <sup>b</sup>	p	$\beta^a$	SE	Stat <sup>b</sup>	p
CPD (raw)	116317	raw	22.631	0.218	1.343	0.229	5.86	4.7E-09	-0.287	0.013	-22.82	5.5E-115	-0.023	0.013	-1.76	7.9E-02
		binned	18.753	0.063	1.030	0.066	15.60	8.3E-55	-0.784	0.037	-21.24	6.1E-100	-0.063	0.039	-1.62	1.0E-01
		Early/Late	16.988	0.059	0.848	0.061	13.99	1.9E-44	1.497	0.084	17.89	1.9E-71	0.194	0.088	2.20	2.8E-02
log <sub>10</sub> (CPD)	116317	raw	1.312	0.005	0.027	0.006	4.86	1.2E-06	-0.007	0.000	-24.06	1.2E-127	0.000	0.000	-0.56	5.8E-01
		binned	1.212	0.002	0.025	0.002	15.20	3.8E-52	-0.020	0.001	-22.62	4.7E-113	0.000	0.001	-0.29	7.7E-01
		Early/Late	1.166	0.001	0.023	0.001	15.31	7.9E-53	0.038	0.002	18.39	1.9E-75	0.003	0.002	1.37	1.7E-01
binned	116317	raw	1.315	0.018	0.108	0.019	5.79	7.2E-09	-0.022	0.001	-21.54	9.6E-103	-0.002	0.001	-1.50	1.3E-01
		binned	1.017	0.005	0.088	0.005	16.41	2.0E-60	-0.061	0.003	-20.36	5.8E-92	-0.006	0.003	-1.83	6.7E-02
		Early/Late	0.879	0.005	0.073	0.005	14.78	2.1E-49	0.118	0.007	17.31	4.7E-67	0.015	0.007	2.12	3.4E-02
Heavy/Light (Multiplicative Scale)	54827	raw	0.914	0.086	0.344	0.087	3.94	8.2E-05	-0.086	0.005	-16.93	2.7E-64	0.000	0.005	-0.05	9.6E-01
		binned	-0.244	0.022	0.345	0.023	14.96	1.3E-50	-0.235	0.014	-17.29	5.5E-67	-0.004	0.014	-0.26	7.9E-01
		Early/Late	-0.773	0.022	0.326	0.022	14.57	4.5E-48	0.457	0.030	15.10	1.6E-51	0.028	0.031	0.91	3.6E-01
Heavy/Light (Additive Scale)	54827	raw	0.618	0.014	0.098	0.015	6.49	8.4E-11	-0.013	0.001	-16.30	1.3E-59	-0.002	0.001	-2.29	2.2E-02
		binned	0.444	0.004	0.071	0.004	16.10	3.6E-58	-0.043	0.002	-17.33	4.3E-67	-0.005	0.003	-1.97	4.9E-02
		Early/Late	0.346	0.004	0.058	0.004	14.09	5.2E-45	0.085	0.006	15.06	3.9E-51	0.014	0.006	2.31	2.1E-02

<sup>a</sup>  $\beta$  refers to the regression slope. For CPD coded as heavy/light,  $\exp(\beta)$  is the Odds Ratio (OR) when analyzed on the multiplicative scale.

<sup>b</sup> Stat refers to either the t- or z-statistic for the linear or logistic regression, respectively.