

1 SCAMPI: A scalable statistical framework for genome-wide interaction testing harnessing cross-
2 trait correlations

3

4 Shijia Bian, MS¹; Andrew J. Bass, PhD²; Yue Liu, PhD³; Aliza P. Wingo, MD, MS^{4,5}; Thomas
5 Wingo, MD³; David J. Cutler, PhD²; Michael P. Epstein, PhD^{2,*}

6

7 **Affiliations:**

8 ¹Department of Biostatistics and Bioinformatics, Emory University, Atlanta, GA, 30329, USA

9 ²Department of Human Genetics, School of Medicine, Emory University, Atlanta, GA, 30329,
10 USA

11 ³Department of Neurology, University of California, Davis, Sacramento, CA 95817, USA

12 ⁴Department of Psychiatry, University of California, Davis, Sacramento, CA 95817, USA

13 ⁵Division of Mental Health, VA Northern California Health Care System, CA 95655, USA

14

15

16

17

18

19

20

21

22

23 * Correspondence: mepste@emory.edu

24 **Abstract**

25 Family-based heritability estimates of complex traits are often considerably larger than their
26 single-nucleotide polymorphism (SNP) heritability estimates. This discrepancy may be due to non-
27 additive effects of genetic variation, including variation that interacts with other genes or
28 environmental factors to influence the trait. Variance-based procedures provide a computationally
29 efficient strategy to screen for SNPs with potential interaction effects without requiring the
30 specification of the interacting variable. While valuable, such variance-based tests consider only a
31 single trait and ignore likely pleiotropy among related traits that, if present, could improve power
32 to detect such interaction effects. To fill this gap, we propose SCAMPI (Scalable Cauchy
33 Aggregate test using Multiple Phenotypes to test Interactions), which screens for variants with
34 interaction effects across multiple traits. SCAMPI is motivated by the observation that SNPs with
35 pleiotropic interaction effects induce genotypic differences in the patterns of correlation among
36 traits. By studying such patterns across genotype categories among multiple traits, we show that
37 SCAMPI has improved performance over traditional univariate variance-based methods. Like
38 those traditional variance-based tests, SCAMPI permits the screening of interaction effects without
39 requiring the specification of the interaction variable and is further computationally scalable to
40 biobank data. We employed SCAMPI to screen for interacting SNPs associated with four lipid-
41 related traits in the UK Biobank and identified multiple gene regions missed by existing univariate
42 variance-based tests. SCAMPI is implemented in software for public use.

43

44

45 **Introduction**

46 Genome-wide association studies (GWAS) have successfully improved our understanding
47 of the role of common single-nucleotide polymorphisms (SNPs) on many complex human traits
48 and diseases. Researchers can further use SNP data from a GWAS study to estimate a trait's
49 narrow-sense heritability (proportion of trait variance due to additive genetic effects) using
50 statistical techniques like GCTA and LD Score Regression (LDSC).^{1; 2} Interestingly, SNP-based
51 heritability estimates of a complex trait are routinely smaller than the corresponding family-based
52 estimates of narrow-sense heritability based on kinship. For instance, studies have reported SNP-
53 based estimates of narrow-sense heritability for body mass index (BMI) to be 0.3, which is
54 considerably less than the narrow-sense heritability estimates of 0.47-0.90 for BMI reported in
55 twin studies.^{3; 4} For Alzheimer's Disease (AD), family-based heritability estimates of the disease
56 range from 0.60-0.80, whereas the latest population-based AD GWAS meta-analyses estimated
57 the narrow-sense heritability from SNP data to be between 0.06-0.41.⁵⁻¹³ Likewise, a GWAS
58 analysis of Amyotrophic Lateral Sclerosis (ALS) estimated SNP-based heritability of
59 approximately 0.21, which is significantly less than the estimates of 0.38-0.85 observed in twin
60 studies.¹⁴

61 The gap between family-based estimates of narrow-sense heritability and corresponding
62 SNP-based estimates may be due to several factors, including rare causal variation poorly tagged
63 by common SNPs as well as shared familial environmental effects ignored in traditional family-
64 based heritability estimation.^{15; 16} Here, we focus on another possible explanation for this gap - the
65 presence of non-additive effects (including higher-order genetic interactions) on complex traits
66 and diseases. As noted in the Supplemental Materials (S1), we can show that higher-order
67 interactions of a complex trait inflate narrow-sense heritability estimates more among close

68 relatives (traditionally used for family-based estimates of heritability) than distantly related
69 individuals (traditionally used to estimate GWAS heritability via LDSC/GCTA).¹⁷ Thus, higher-
70 order interactions can explain the discrepancy between family-based and SNP-based heritability
71 estimates observed for many complex human traits. This motivates the search for genetic variants
72 in large-scale genetic studies that demonstrate non-additive effects, including gene-gene and gene-
73 environment interactions.

74 While studies have identified SNPs demonstrating interaction effects on complex traits,¹⁸⁻
75 ²³ genome-wide investigation of non-additive effects is inherently challenging.^{24; 25}
76 Comprehensive genome-wide testing of SNP-SNP (epistatic) interactions is computationally
77 intractable as 10 million SNPs can lead to approximately 5×10^{13} potential interaction tests. Even
78 if such analyses were tractable, the resulting multiple-testing adjustment cripples the power to
79 detect epistatic effects. Gene-environment interaction analyses require fewer tests and are more
80 computationally feasible, but measuring the right environmental determinants can be difficult and
81 is often unknown.²⁶⁻²⁹ To circumvent uncertainty about the right environmental factor yet still test
82 for evidence of interaction, Paré et al. proposed an efficient variance-based method for a
83 quantitative trait that screens for SNPs with possible interactive effects without requiring
84 specification of the interacting factor.³⁰ Recognizing that a SNP with an interaction effect on a trait
85 induces trait variance that differs by genotype (see Supplemental Figure S1), Paré screened for
86 SNPs with potential interaction effects by testing for equality of variances across genotype
87 categories using Levene's test.³¹ Researchers have successfully applied this type of variance-based
88 approach within the UK Biobank to identify genetic variants with interaction effects on obesity
89 phenotypes and cardiometabolic serum biomarkers.^{32; 33}

90 The variance-based test of Paré is a univariate test that considers whether a SNP has an
91 interactive effect with a single phenotype. However, biobanks routinely collect detailed
92 information on a large collection of related phenotypes with shared genetic effects. Many recent
93 methods of gene mapping illustrate the appeal of leveraging the ubiquitous phenomenon of
94 pleiotropy across related traits when present.³⁴⁻³⁷ Consequently, if pleiotropic genetic variants with
95 interactive effects exist, we expect a multi-trait statistical method that leverages this information
96 will have improved performance over existing univariate variance-based interaction procedures.
97 Bass et al. recently showed that a SNP with an interaction effect induces not only variance but also
98 covariance patterns between traits that differ by genotype (which we illustrate in Supplemental
99 Figure S2).³⁸ Based on this observation, the authors developed a kernel framework for interaction
100 testing that assessed where similarity in variance/covariance patterns among a group of modeled
101 traits correlated with genotypic similarity at a test SNP. While more powerful than standard
102 variance-based testing, the kernel framework of Bass lacks practical features for genetic analysis
103 such as the inability to identify the specific phenotypes (among those modeled) that demonstrate
104 interaction effects with the test SNP. Identifying these specific phenotypes are of substantial value
105 for further downstream analyses.

106 To this end, we propose here an efficient screening method SCAMPI (Scalable Cauchy
107 Aggregate test using Multiple Phenotypes to test Interactions) for identifying potential SNPs with
108 interaction effects using multiple phenotypes. SCAMPI fits simple regression models relating SNP
109 genotype to (standardized) cross products of all pairwise combinations of traits under
110 consideration and then aggregates the correlated p-values from these separate regression tests
111 together into an omnibus test using the Cauchy Combination Test.^{39; 40} Similar to variance-based
112 interaction tests, SCAMPI does not require specification of the factor that interacts with the SNP

113 of interest, thereby reducing the computational and testing burden and enabling the scaling of the
114 method to biobank-size datasets. Moreover, SCAMPI scales to handle many related phenotypes
115 and can identify the specific phenotype(s) that have interaction effects among those modeled.
116 Using simulations, we show that SCAMPI can detect interactions under various scenarios and has
117 improved performance over univariate variance-based interaction procedures. We also applied
118 SCAMPI to lipid panel data (an indicator of risk of heart disease and stroke) in the UK Biobank
119 (UKBB) and identified several genes with putative interaction effects that were missed by standard
120 univariate variance-based procedures. For public use, SCAMPI is implemented as an R package.

121

122 **Materials and Methods**

123 Motivation: We first show that a SNP with a pleiotropic interaction effect yields trait
124 correlation patterns that differ by genotype category. We could analogously show that a SNP with
125 a pleiotropic interaction trait effect influences the covariance patterns between traits but chose to
126 focus on correlation due to the scale-free nature of the latter measure. For subject i , define G_i as
127 the subject's genotype at a test SNP and define W_i as some factor (either genetic or environmental)
128 that interacts with the SNP to influence multiple traits. Suppose subject i possesses two correlated
129 traits $Y_{i,1}$ and $Y_{i,2}$ that are generated under the relationships:

$$130 \quad Y_{i,1} = \alpha_1 + \beta_1 G_i + \gamma_1 W_i + \delta_1 G_i W_i + \epsilon_{i,1}; \quad Y_{i,2} = \alpha_2 + \beta_2 G_i + \gamma_2 W_i + \delta_2 G_i W_i + \epsilon_{i,2}.$$

131 Here, $\beta_j, \gamma_j, \delta_j$ denote the main effect of genotype, the main effect of the factor, and two-way
132 interaction effect between genotype and factor, respectively, on trait j ($j = 1, 2$). We further
133 assume each of the error terms $\epsilon_{i,1}$ and $\epsilon_{i,2}$ has a standard normal distribution $\epsilon_{i,1}, \epsilon_{i,2} \sim N(0,1)$.
134 Without loss of generality, further assume W_i is distributed as $W_i \sim N(0,1)$ and is independent of
135 G_i .

136 Based on the trait models listed above, Paré previously showed that that variance of Y_1 (Y_2)
137 differs by G when the genotype has an interaction effect on trait 1 (trait 2), respectively.³⁰
138 Additionally, when pleiotropic interaction effects exist, we can show the correlation of traits 1 and
139 2 also differ by genotype. In Supplementary Materials S2, we derive the correlation between $Y_{i,1}$
140 and $Y_{i,2}$ conditional on genotype G_i as:

$$141 \quad \text{cor}(Y_{i,1}, Y_{i,2} | G_i = g_i) = \frac{\gamma_1 \gamma_2 + (\delta_1 \gamma_2 + \delta_2 \gamma_1) g_i + \delta_1 \delta_2 g_i^2}{\sqrt{(\gamma_1 + \delta_1 g_i)^2 + 1} \times \sqrt{(\gamma_2 + \delta_2 g_i)^2 + 1}} \quad (1)$$

142 Equation (1) shows that the correlation between two traits differs by genotype when either a) the
143 genotype interacts with the factor on both phenotypes or b) the genotype interacts with the factor
144 on at least one of the phenotypes, provided the factor has a main effect on the other phenotype.
145 We can see that if the SNP has no interaction effect on either phenotype ($\delta_1 = \delta_2 = 0$), the
146 phenotypic correlation will not differ by genotype even when main effects for the factor exist ($\gamma_1 \neq$
147 $0, \gamma_2 \neq 0$).

148 The above result suggests an efficient strategy for screening SNPs with potential
149 interaction effects. Instead of performing traditional interaction analyses, which mandates defining
150 potential interacting factors W_i , we can instead screen for SNPs with interaction effects without
151 having to specify W_i by examining whether the correlation between traits changes as a function of
152 the linear and quadratic effects of genotype. Such modeling provides a workaround in situations
153 where interacting covariates are uncollected or inaccurately recorded. The screening procedure
154 further provides an efficient alternative strategy for genome-wide epistatic analysis in that it does
155 not require direct modeling of the interacting genetic factor, which substantially reduces the
156 number of tests to be considered. If we are analyzing M SNPs, SCAMPI requires only M tests

157 whereas comprehensive epistatic analysis requires $\binom{M}{2}$ tests. Thus, when $M = 200\text{K}$ ($M = 2M$),

158 SCAMPI reduces the number of tests required by approximately 5 (6) orders of magnitude.

159 Rather than model trait correlation as a function of linear and quadratic effects of genotype
160 mentioned above, we note that we can alternatively parameterize this relationship using a general
161 genotype model that allows for separate effects of each genotype relative to a baseline category.
162 That is, for some outcome Y^* , the coefficient estimates of $\hat{\alpha}_1$, $\hat{\alpha}_2$ and $\hat{\alpha}_3$ in the regression model
163 $Y^* = \alpha_1 + \alpha_2 G + \alpha_3 G^2 + \epsilon^*$ can be directly mapped to coefficient estimates $\hat{\beta}_1$, $\hat{\beta}_2$ and $\hat{\beta}_3$ in a
164 model $Y^* = \beta_1 + \beta_2 G_1 + \beta_3 G_2 + \epsilon$, where G_1 and G_2 are genotype indicators for those with 1 and
165 2 copies of the reference allele, respectively (those with 0 copies are treated as baseline). Given
166 the familiarity of this general genotype model in GWAS,⁴¹⁻⁴⁴ we chose to use this alternative
167 parameterization in our method moving forward.

168 *Notation and Trait Standardization:* Assume a sample of N unrelated subjects that possess
169 J continuous phenotypes. Let $\mathbf{Y}_j = (Y_{1,j}, Y_{2,j}, \dots, Y_{N,j})^T$ denote the $N \times 1$ vector of observations for
170 trait j ($j = 1, \dots, J$). Define $\mathbf{G} = (G_1, G_2, \dots, G_N)^T$ as an $N \times 1$ vector of genotypes for one test
171 SNP, where G_i represents $[0, 1, 2]$ copies of the minor allele that subject i possesses at the site. As
172 noted in the previous section, we are interested in applying a general genotype model for
173 interaction testing as it naturally captures the linear and quadratic effects of genotype shown in
174 equation (1). Consequently, further define $G_i^{(1)} = I[G_i = 1]$ and $G_i^{(2)} = I[G_i = 2]$ as indicator
175 variables for genotype categories 1 and 2, respectively (we treat genotype category 0 as baseline).
176 Finally, let \mathbf{Z} be an $N \times K$ matrix of confounding variables. These confounding variables can be a
177 mixture of continuous or categorical features. Common confounder examples include age,
178 biological sex, batch ID, and principal components of ancestry to deal with population
179 stratification.

180 Our goal is to detect a SNP with an interaction effect that yields correlation patterns that
 181 differ by genotype. Such trait pattern differences can erroneously arise if the main effect of the
 182 genotype, as well as main and variance effects of confounders (such as population structure), are
 183 unaccounted for prior to analysis.^{45; 46} To avoid this issue, we first standardize and adjust each
 184 $Y_j (j = 1, \dots, J)$ prior to analysis using a double generalized linear model (DGLM) that corrects for
 185 the mean effects of the test SNP and confounders, as well as the potential variance effects of
 186 confounders.^{47,48}

187 DGLM is composed of two sub-models, where the first sub-model controls population
 188 mean, and the second sub-model controls population variance. For our work, the first sub-model
 189 adjusts Y_j for the mean effects of $\mathbf{G}^{(1)}$ & $\mathbf{G}^{(2)}$ and confounders \mathbf{Z} using the following framework:

$$190 \quad Y_j = [\mathbf{1} \quad \mathbf{G}^{(1)} \quad \mathbf{G}^{(2)} \quad \mathbf{Z}] \begin{bmatrix} \beta_{j,0} \\ \beta_{j,G^{(1)}} \\ \beta_{j,G^{(2)}} \\ \boldsymbol{\beta}_{j,Z} \end{bmatrix} + \boldsymbol{\varepsilon}_j$$

191 where $\beta_{j,0}$ is the intercept associated with the j^{th} trait. $\beta_{j,G^{(1)}}$ and $\beta_{j,G^{(2)}}$ are the regression
 192 coefficient for $\mathbf{G}^{(1)}$ and $\mathbf{G}^{(2)}$ respectively, and $\boldsymbol{\beta}_{j,Z}$ is a $K \times 1$ vector of regression coefficients for
 193 confounders \mathbf{Z} . Finally, $\boldsymbol{\varepsilon}_j$ is a $N \times 1$ vector of residual errors that follow

$$194 \quad \boldsymbol{\varepsilon}_j \sim MVN \left(\boldsymbol{\mu}_{\boldsymbol{\varepsilon}_j} = \mathbf{0}_{N \times 1}, \boldsymbol{\Sigma}_{\boldsymbol{\varepsilon}_j} = \begin{bmatrix} \sigma_{\boldsymbol{\varepsilon}_{1,j}}^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_{\boldsymbol{\varepsilon}_{N,j}}^2 \end{bmatrix} \right) \quad (2)$$

195 The second sub-model of the DGLM then models $\boldsymbol{\varepsilon}_j$ in (2) as a function of confounders \mathbf{Z} using
 196 the following framework using the log link function:

$$197 \quad \log \left(E \left(\text{diag} \left(\boldsymbol{\Sigma}_{\boldsymbol{\varepsilon}_j} \right) \right) \right) = \boldsymbol{\varphi}_j = [\mathbf{1} \quad \mathbf{Z}] \begin{bmatrix} \gamma_0 \\ \boldsymbol{\gamma}_Z \end{bmatrix}$$

198 where $\boldsymbol{\varphi}_j$ is the $N \times 1$ column vector representing the expected residual variance of the j^{th}
 199 observed trait. Here, γ_0 is the intercept while $\boldsymbol{\gamma}_Z$ represents the $K \times 1$ column vector of
 200 confounder effects on the variance. The error distribution to be used in the two sub-models is
 201 Gaussian.

202 We fit the above DGLM using the R package “dglm”. Let $\tilde{\mathbf{Y}}_j$ ($j = 1, \dots, J$) denote the
 203 adjusted and standardized form for trait j produced from the DGLM model fit. We subsequently
 204 use $\tilde{\mathbf{Y}}_j$ ($j = 1, \dots, J$) to construct appropriate measures for our downstream screening analyses for
 205 interaction effects.

206 Analysis Strategy: For $J = 2$ traits, we show in Supplemental Materials (S3) that we can
 207 approximate the sample Pearson correlation coefficient of traits \mathbf{Y}_1 and \mathbf{Y}_2 as the average of the
 208 $N \times 1$ vector of cross products of the traits after standardization, $\tilde{\mathbf{Y}}_1$ and $\tilde{\mathbf{Y}}_2$. That is, we estimate
 209 the Pearson correlation between \mathbf{Y}_1 and \mathbf{Y}_2 as the sample average of

$$210 \quad \tilde{\mathbf{Y}}_1 \odot \tilde{\mathbf{Y}}_2 = (\tilde{Y}_{1,1} \cdot \tilde{Y}_{1,2}, \tilde{Y}_{2,1} \cdot \tilde{Y}_{2,2}, \dots, \tilde{Y}_{N,1} \cdot \tilde{Y}_{N,2})^T$$

211 where \odot denotes the row-wise product operator of two vectors. Similarly, we can estimate the
 212 variance of \mathbf{Y}_1 and \mathbf{Y}_2 by $\tilde{\mathbf{Y}}_1 \odot \tilde{\mathbf{Y}}_1$ and $\tilde{\mathbf{Y}}_2 \odot \tilde{\mathbf{Y}}_2$, respectively.

213 Using these estimates, we construct a screening procedure to identify a SNP with an
 214 interaction effect on trait \mathbf{Y}_1 and/or \mathbf{Y}_2 by assessing whether SNP genotype \mathbf{G} is associated with
 215 either $\tilde{\mathbf{Y}}_1 \odot \tilde{\mathbf{Y}}_1$, $\tilde{\mathbf{Y}}_2 \odot \tilde{\mathbf{Y}}_2$, or $\tilde{\mathbf{Y}}_1 \odot \tilde{\mathbf{Y}}_2$. Examination of the relationship of \mathbf{G} with $\tilde{\mathbf{Y}}_1 \odot \tilde{\mathbf{Y}}_1$ (or
 216 $\tilde{\mathbf{Y}}_2 \odot \tilde{\mathbf{Y}}_2$) is similar to assessing whether trait variance differs by genotype (which Paré³⁰
 217 investigated using Levene’s test) while the study of \mathbf{G} with $\tilde{\mathbf{Y}}_1 \odot \tilde{\mathbf{Y}}_2$ leverages additional
 218 information on interactions based on differences in trait correlations. To implement our procedure,
 219 we fit 3 separate linear regression models; each model treating one of $\tilde{\mathbf{Y}}_1 \odot \tilde{\mathbf{Y}}_1$, $\tilde{\mathbf{Y}}_2 \odot \tilde{\mathbf{Y}}_2$, or
 220 $\tilde{\mathbf{Y}}_1 \odot \tilde{\mathbf{Y}}_2$ as outcome with SNP genotype $\mathbf{G}^{(1)}$ and $\mathbf{G}^{(2)}$ as predictors. Each regression models

221 produces a p-value based on a two-degree-of-freedom test. Since the resulting 3 p-values from
222 these regression tests are correlated, we can then combine them into an omnibus p-value (described
223 in the next section) to assess whether the SNP has an interaction effect on at least one of the two
224 traits under study.

225 The above example considered two traits under study. However, the strategy easily extends
226 to the study of $J > 2$ correlated traits as well. Assuming J traits, we fit J regression models that
227 regress $\tilde{Y}_j \odot \tilde{Y}_j$ on $\mathbf{G}^{(1)}$ and $\mathbf{G}^{(2)}$ ($j = 1, \dots, J$) and further fit $\binom{J}{2}$ additional regression models
228 that regress $\tilde{Y}_j \odot \tilde{Y}_l$ ($j, l = 1, \dots, J; j \neq l$) on $\mathbf{G}^{(1)}$ and $\mathbf{G}^{(2)}$. We then can combine the $J_C =$
229 $(J + \binom{J}{2})$ p-values from these tests together to assess whether the SNP has an interaction effect on
230 at least one of the J traits under study.

231 Cauchy Combination Test (CCT): After obtaining the J_C p-values above, we create a final
232 omnibus test for whether the test SNP has an interactive effect on any of the traits under
233 consideration using the Cauchy Combination Test (CCT),^{39; 40} which is a popular technique for
234 aggregating many potentially dependent tests of high dimension together into an omnibus
235 framework. CCT has provable type I error rate control for genome wide significance thresholds
236 even when p-values are dependent. CCT is especially useful when an SNP signal is sparse and
237 only affects a subset of the traits under consideration. The test statistics of CCT is a weighted sum
238 of the Cauchy transformation of individual p-values in SCAMPI. Let p_r to denote the dependent
239 individual p-value from the r^{th} regression test ($r = 1, 2, \dots, J_C$). The CCT statistic is defined as

240
$$T = \sum_{r=1}^{J_C} \frac{1}{J_C} \tan\{(0.5 - p_r)\pi\} \quad (3)$$

241 Under the null hypothesis of no SNP interactive effect with any of the traits under consideration,
242 T in (3) follows a standard Cauchy Distribution, i.e., $T \sim \text{Cauchy}(X_0 = 0, \gamma = J_C)$. This derived p-
243 value is the SCAMPI p-value at the given genotype \mathbf{G} .

244 Overview of the SCAMPI Framework: Our SCAMPI framework aggregates the regression tests
245 outlined earlier with the CCT to produce an omnibus p-value for testing whether the SNP has an
246 interactive effect with at least one of the traits under study. SCAMPI, which is implemented in a
247 public R package of the same name, requires the following inputs:

- 248 a. Multiple target traits are denoted as \mathbf{Y} . Should these traits not follow a normal distribution,
249 users can apply a rank-based Inverse Normal Transformation to normalize the traits, if
250 desired.
- 251 b. The confounding variables, represented by \mathbf{Z} ;
- 252 c. One test SNP, represented by \mathbf{G} and coded as $\mathbf{G}^{(1)}$ and $\mathbf{G}^{(2)}$.

253 SCAMPI then follows the workflow depicted in Figure 1.

254 Application to UK Biobank Data: We applied SCAMPI to identify SNPs with potential
255 interaction effects on lipid measures within the UK Biobank (application ID 42223). We focused
256 attention on four lipid-related measures: high-density lipoprotein cholesterol (HDL-C), low-
257 density lipoproteins cholesterol (LDL-C), triglycerides (TGs), and Body Mass Index (BMI). Both
258 the sample and SNP QC procedures are in accordance with Marderstein et al.⁴⁹ Similar QC
259 procedures were also carried out in multiple studies.⁵⁰ From the cohort, we excluded individuals
260 who either (1) had missing heterozygosity information, (2) were outliers in terms of heterozygosity
261 or had missing genotype rates greater than 0.02, (3) had over 10 putative third-degree relatives in
262 the kinship table, (4) were omitted from the kinship inference procedure, or (5) were either self-
263 reported as anything other than ‘White British’ or did not show similar genetic ancestry to this

264 group based on a principal components analysis of the genotypes. After performing this quality
265 control, 337,422 independent subjects remained ($N_{\text{Female}}= 181,203$; $N_{\text{Male}}= 156,219$). Moreover,
266 the UKB employed two genotyping arrays. In this post-QC sample, we have the UK Biobank
267 Axiom array ($N_{\text{UKBB}}= 300,345$) and the UK BiLEVE array ($N_{\text{UKBL}}= 37,077$). For the SNP QC,
268 genotypes were discarded if they had an INFO score < 0.8 , MAF < 0.05 and HWE p-value $< 10^{-10}$.
269 After SNP QC procedures, 288,910 SNPs were retained. Finally, 277,653 SNPs were included
270 for analysis using SCAMPI after applying a 10% missing rate threshold.

271 We first adjusted the four lipid-related traits for confounders, including the first six genetic
272 principal components, biological sex, age, age squared (age^2), and the type of genotyping array,
273 before applying SCAMPI to these traits. Notably, the first six principal components effectively
274 captured population structure at subcontinental geographic scales.^{51,52} Of the initial set of 337,422
275 independent subjects, 288,709 possessed complete information on all traits and confounders and
276 were considered moving forward. We first transformed the four traits using the inverse normal
277 transformation (INT) to align the traits, which is a common practice to ensure the residual of traits
278 is normally distributed in a regression model such as DGLM.⁵³⁻⁵⁶ The distribution of the four traits,
279 both pre and post-INT, can be found in Supplemental Figure S3 (a) - (d). Correlation between post-
280 INT traits was 0.1246 for HDL-C and LDL-C, -0.4938 for HDL-C and TG, -0.3809 for HDL-C
281 and BMI, 0.2797 for LDL-C and TG, 0.0394 for LDL-C and BMI, and 0.3708 for TG and BMI.

282 Simulations: We conducted comprehensive simulations to evaluate the type-I error rate of
283 SCAMPI under a variety of scenarios. For each scenario, we simulated a sample size of 300,000
284 to reflect biobank-scale datasets. Each scenario is analyzed based on 100,000 simulations. We
285 assumed $J = 2, 4, 8$ traits and simulated the trait values for the i^{th} individual based on the
286 multivariate normal distribution illustrated below:

$$287 \begin{pmatrix} Y_{i,1} \\ \vdots \\ Y_{i,J} \end{pmatrix} \sim MVN \left(\boldsymbol{\mu} = \begin{bmatrix} \alpha_1 + \beta_1 G_i + \gamma_1 W_i \\ \vdots \\ \alpha_J + \beta_J G_i + \gamma_J W_i \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & \dots & \sigma^2 \\ \vdots & \ddots & \vdots \\ \sigma^2 & \dots & 1 \end{bmatrix} \right) \quad (4)$$

288 For predictors, we generated the test SNP genotype G_i under Hardy-Weinberg equilibrium,
289 assuming the SNP had a minor-allele frequency of either 0.05 or 0.25. We further generated a
290 factor W_i that followed a standard normal distribution. For the choice of parameters in the
291 equation, we simulated the intercept α_j from $N(0, 5)$, the genotype main effect β_j from
292 $Unif(0, 0.2)$, and the factor main effect γ_j from $Unif(0, 0.3)$ ($j = 1, \dots, J$). In the covariance
293 matrix Σ in (4), the off-diagonal covariance elements are assigned as σ^2 . We performed different
294 simulations assuming $\sigma^2 = 0.01$ (negligibly correlated traits), 0.25 (moderately correlated traits),
295 and 0.5 (strongly correlated traits). For $J = 4$ traits, we conducted additional simulations where
296 we considered a specific covariance matrix that mirrored the observed covariance structure of the
297 lipid-related traits that we studied in the UKBB dataset. Finally, we conducted additional type-I
298 error simulations based directly on our UKBB sample. Specifically, we randomly permuted the
299 UKBB phenotype data (consisting of our four trait outcomes and confounding variables) across
300 subjects and then re-ran SCAMPI on the genome-wide data. We repeated the permutation process
301 four times, which resulted in a total of >1M SCAMPI p-values under the null hypothesis.

302 For power simulations, we implemented a similar simulation design as for our type-I error
303 simulations but introduced additional parameters to model the effect of the interaction between

304 SNP and the factor on the simulated traits. Specifically, we generated J traits based on the
305 multivariate normal distribution as presented in Eq. (5):

$$306 \quad \begin{pmatrix} Y_{i,1} \\ \vdots \\ Y_{i,J} \end{pmatrix} \sim MVN \left(\boldsymbol{\mu} = \begin{bmatrix} \alpha_1 + \beta_1 G_i + \gamma_1 W_i + \delta_1 G_i W_i \\ \vdots \\ \alpha_J + \beta_J G_i + \gamma_J W_i + \delta_J G_i W_i \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 1 & \dots & \sigma^2 \\ \vdots & \ddots & \vdots \\ \sigma^2 & \dots & 1 \end{bmatrix} \right)$$

307 (5)

308 δ_j ($j = 1, \dots, J$) in equation (5) represents the interaction effect of the SNP and factor on trait j .
309 For a given simulation scenario, we vary the percentage of traits that possess such an interaction
310 (i.e. the sparsity of the interaction signal) among the values 25%, 50%, 75%, and 100%. For those
311 traits with an interaction effect, we vary the value of δ_j across a range of values from 0.01 to 0.50
312 to study how the power trends change as δ_j increases for each scenario. The settings for the number
313 of traits, MAF, α_j , β_j align with those in the Type I error simulations. However, γ_j is held at fixed
314 values for all traits instead of being simulated from a distribution. Without loss of generality, this
315 approach eliminates the potential for power fluctuations arising from the randomness in γ_j . We
316 simulated the results for various combinations under different parameter sets. To illustrate the
317 overall pattern of the power simulation, we selected the simulation with $\gamma = 0.05$ and 0.25 , and
318 $\sigma^2 = 0.1, 0.3$ and 0.5 . For each simulation scenario, we assumed a sample size of 20K and
319 generated 10K replicates for inference.

320 We chose to benchmark SCAMPI against an enhanced multi-phenotype version of
321 Levene's test that was originally restricted to a single phenotype.³⁰ This enhanced version is termed
322 as the multivariate Levene's test in our context. The multivariate Levene's test applies Levene's
323 test (described in Supplemental Materials S4) to each trait separately, resulting in J p-values. These
324 J p-values are then aggregated together into an omnibus test using the CCT methodology detailed
325 in the prior section (see Supplemental Figure S4 for an outline of the framework). While this

326 benchmark examines how variances vary by genotype across different traits, it does not consider
327 difference in correlation patterns among traits that SCAMPI integrates within its framework.

328 **Results**

329 Simulation Studies: Table 1 provides empirical type 1 error rates for SCAMPI summarized
330 at a nominal rate α of 10^{-2} and 10^{-3} across varying numbers of phenotypes, MAF and Σ when
331 γ is simulated from *Unif*(0, 0.3). As described in Supplemental S5, we focused primarily on
332 studying the empirical type-I error rate at 10^{-3} based on the number of simulations performed and
333 observed that SCAMPI was well calibrated at such a threshold. To examine whether SCAMPI was
334 well calibrated at more stringent thresholds, we studied type-I error rates based on permutation of
335 the UKBB data, which yielded $> 1M$ tests under the null hypothesis. For these null simulations,
336 we observed the type I error rates of SCAMPI to be 1.08×10^{-2} , 1.06×10^{-3} and 9.81×10^{-5}
337 at α of 10^{-2} , 10^{-3} and 10^{-4} , respectively. SCAMPI p-values generally followed the same pattern
338 as p-values of other statistical methodology that employs CCT.^{39; 57-59}

339 We assessed the power of SCAMPI in different scenarios. Figure 2 provides representative
340 power results at a genome-wide significance threshold of 6.25×10^{-8} (based on a multiple-testing
341 correction for the total number of $\sim 800,000$ SNPs in the UKBB) assuming $J=4$ traits and a
342 correlation matrix that mirrored the observed correlation structure of the lipid-related traits that we
343 studied in the UKBB dataset. Figure 2 is comprised of four sub-figures, with each sub-figure
344 presenting simulation results and assuming a different level of sparsity for the interaction effect
345 among the traits modeled. For example, Figure 2a assumes the test SNP has an interaction effect
346 with only one of the four traits, while Figure 2d assumes the test SNP has an interaction effect on
347 all four traits. Within each sub-figure, the yellow solid line represents the power of SCAMPI while
348 the dashed green line represents the power of Multivariate Levene's test. Within each sub-figure,

349 results show, as expected, that the power of both SCAMPI and Multivariate Levene's test increases
350 as the magnitude of the interaction effect increases. Further, the power of each method increases
351 as the number of traits the SNP has an interaction effect with increases (or, similarly, the sparsity
352 of the interaction effect decreases). However, across all four sub-figures, SCAMPI consistently
353 shows improved power over Multivariate Levene's test. We note that such improved power of
354 SCAMPI over Multivariate Levene's test holds even when the SNP has an interaction effect on
355 only one of the traits under study (Figure 2a), which suggests that the inclusion of traits with no
356 interaction effects still contributes valuable information to the SCAMPI test via their correlation
357 with the trait that does have an interaction effect. We do see that, in Figure 2b and Figure 2c,
358 SCAMPI experiences a pattern at $\delta=0.4$ and $\delta=0.25$, respectively, where power dips slightly at the
359 parameter value; this pattern emerges under conditions where interaction effects are present in
360 multiple, but not all, traits. It results from randomly assigning interaction effects to a subset of
361 traits, provided that the pairwise correlation among the traits are distinct. While the Multivariate
362 Levene's test does not exhibit this behavior (since it only considers the variance of the traits under
363 study), we find that SCAMPI is still more powerful in these situations. We also overlay the power
364 curve of SCAMPI and Multivariate Levene's test with varying sparsity for better visualization in
365 the same plot in Supplemental Figure S5.

366 In addition to the power simulations inspired by the UKBB, Supplemental Figure S6
367 provides power results for SCAMPI and multivariate Levene's test under a broader range of
368 models that vary the number of traits considered, the sparsity of the interaction effect, the
369 correlation among traits, and the main effect of the variable interacting with genotype. Overall, we
370 find the power of SCAMPI increases with a decrease in the sparsity of the interaction effect, a
371 decrease in the trait correlation, and an increase in the effect size of the interaction variable.

372 Assuming these three inputs are fixed, we find that the power of SCAMPI increases as the number
373 of traits modeled increases. Regarding the power comparisons between SCAMPI and the
374 Multivariate Levene’s test under this broader range of models, Supplemental Figure S6 also
375 reaffirms the trends observed in our UKBB-inspired power simulations. Across the spectrum of
376 scenarios tested, SCAMPI consistently exhibited superior performance when compared to the
377 Multivariate Levene’s test, largely because the former method accounts for correlation among
378 traits that the latter method ignores.

379 Application to UKBB: Figure 3 provides the Manhattan plot of SCAMPI results for
380 detecting interaction effects on four lipid-related traits. SCAMPI identified 210 SNPs across 68
381 genes and intergenic regions at a study-wide significance level ($\alpha = 1.67 \times 10^{-7}$, i.e., multiple
382 comparison correction for 300,000 SNPs). Table 2 highlights the SNPs with the smallest SCAMPI
383 p-value on each chromosome from the 210 SNPs. A comprehensive list of the 210 SNPs is
384 available in the Supplemental Table S1. The Q-Q plot for SCAMPI (Supplemental Figure S7)
385 shows no evidence of inflation. SCAMPI is an omnibus test that, by aggregating p-values (outputs
386 of Step 3 in Figure 1) from association tests of trait correlation, pinpoints the specific traits that
387 influence the overall signal. Thus, for every lead SNP in Table 2, we examined the p-values linked
388 to each trait variance and cross-trait correlation at a genome-wide significance threshold of
389 1.67×10^{-7} . Significant variance and correlation terms among traits are noted in the “Significant
390 Variance/Correlation Components” column of the Table. For example, SNP rs7528419 on
391 CELSR2 is significantly associated with the correlation of triglycerides and LDL, as well as the
392 variance of LDL alone, suggesting the SNP may have an interaction effect with other genetic or
393 environmental factors on these two specific traits that merit further investigation.

394 We also cross-referenced our findings in Table 2 with PheWAS results based on the GWAS
395 Catalog or UK Biobank from the Open Targets Platform (v22.10) , which confirmed many of our
396 initial findings.⁶⁰ For instance, SNP rs738409 in PNPLA3 (which SCAMPI identified to be
397 associated with the correlation of triglycerides and BMI as well as triglyceride variance) is reported
398 by Open Targets Platform to be significantly linked with BMI. These results of the lead SNPs are
399 cross listed in the “PheWAS” column of Table 2. Beyond the lead SNPs, Supplemental Table S2
400 includes the p-values for all correlation components related to the 210 SNPs.

401 Overall, SCAMPI identified several established lipid- and BMI-related genes that also
402 demonstrate potential interaction effects. For example, *APOC1*, which contained the smallest
403 SCAMPI p-value ($p=8.1 \times 10^{-61}$), has pleiotropic effects on lipid metabolism, influencing
404 various processes through its actions on lipoprotein receptors and enzyme activity modulation. By
405 controlling the lipids plasma level, the influence of *APOC1* spans several disease areas, including
406 cardiovascular physiology, inflammation, immunity, sepsis, diabetes, cancer, viral infectivity, and
407 cognition.⁶¹ Furthermore, *CETP*, which contained a SNP demonstrating a possible interaction
408 effect with HDL ($p=7.61 \times 10^{-39}$), may prevent plaque buildup and protect from atherosclerotic
409 cardiovascular disease.⁶² There are also mixed results regarding the modifying effects of *CETP* on
410 cardiovascular events.⁶³⁻⁶⁵ Another top gene identified by SCAMPI was *LIPC*. Evidence suggests
411 the *LIPC* promoter polymorphism (T-514C) affects the activity of Hepatic lipase (HL) and, in
412 concert with other factors, modifies the therapeutic response in coronary artery disease (CAD)
413 patients, with those having the CC genotype benefiting the most from intensive lipid-lowering
414 treatments due to their predisposition to high HL activity and smaller, denser LDL particles.⁶⁶
415 SCAMPI also identified SNPs in *CELSR2* with interaction effects predominantly on lipids.
416 Research has shown *CELSR2* deficiency impacts intracellular Ca^{2+} levels, possibly due to

417 compromised endoplasmic reticulum (ER) function and unfolded protein response (UPR). The
418 depletion of *CELSR2* affects the expression of UPR sensors and the splicing of XBP-1, a critical
419 transcription factor for hepatic lipogenesis, as demonstrated by reactions to various cellular
420 stresses.⁶⁷

421 Interestingly, SCAMPI identified several SNPs (shown in Supplemental Table S3)
422 exclusively through the correlation among traits (such that they were not detected by the
423 multivariate Levene's test that only considered variance terms). Noteworthy among these are
424 rs2228603 (*NCAN*), rs58542926 (*TM6SF2*), and rs10415849 (*GATAD2A*). For each of these three
425 SNPs, SCAMPI detected a significant effect exclusively via the correlation of BMI and
426 triglycerides (each $p < 10^{-8}$); the SNP was not significantly associated with the variance of either
427 trait and, as such, was not picked up by Levene's test. Prior PheWAS studies show an association
428 between these SNPs and triglycerides.⁶⁸⁻⁷¹ A similar pattern is observed for three SNPs in
429 *NECTIN2*; each SNP is associated with the correlation of LDL and HDL (each $p < 10^{-8}$) but not
430 with the variance of either trait. PheWAS analysis previously demonstrated the association of these
431 SNPs with LDL. Beyond PheWAS, we also want to highlight that the SNPs identified by SCAMPI
432 have been implicated in other studies of lipid traits and BMI. For example, numerous studies
433 suggest that rs2228603 and rs58542926 are risk alleles associated with an increased likelihood of
434 liver inflammation and fibrosis that is closely associated with weight change, indeed impacting
435 BMI.⁷²⁻⁷⁴ rs10415849 is significantly associated with α -Tocopherol (one type of vitamin E), which
436 interacts with biological sex to modify BMI.⁷⁵ The two SNPs rs519113 and rs6859, which are
437 *BCL3-PVRL2-TOMM40* SNPs, imply gene-gene and gene-environment interactions on
438 dyslipidemia, which pathophysiology is characterized by reverse cholesterol transport in HDL
439 metabolism.^{76; 77} Even though there are not many direct studies showing the association between

440 rs3852860 and HDL, rs3852860 is a well-known predictor in Alzheimer's disease, and
441 Alzheimer's disease progressed with HDL change.⁷⁸⁻⁸⁰

442 SCAMPI Analysis in UKBB Adjusting for APOE: In our applied analyses of lipid traits and
443 BMI in the UKBB, the strongest signal detected by SCAMPI was located within *APOC1*, which
444 is in close physical proximity to *APOE*, a gene with established relevance to the lipid traits we
445 examined. Given APOE's prominence as a biomarker in lipid panels,⁸¹ we determined whether the
446 signals we observed at *APOC1* were independent of those at *APOE*. To assess this, we repeated
447 our SCAMPI analyses conditioning on the main and variance effects of *APOE* SNPs. Specifically,
448 we selected all SNPs on *APOE*, located within 45,409,113 and 45,412,532 on chromosome 19,
449 based on the Genome Reference Consortium Human Build 37 (GRCh37). Five SNPs (rs440446,
450 rs769449, rs769450, rs429358, and rs7412) within this region passed the SNP level QC. We
451 adjusted for the effects of the five *APOE* SNPs on the phenotypic outcomes' mean and variance
452 and then reapplied the SCAMPI methodology. We note that the sample size for our adjusted
453 SCAMPI analysis dropped from 288,709 samples to 241,167 samples due to missing genotypes at
454 the five *APOE* SNPs.

455 We provide the Manhattan and Q-Q plots for the APOE-adjusted SCAMPI analyses in
456 Supplemental Figure S8. Overall, SCAMPI identified 150 SNPs (see Supplemental Table S4) that
457 remained significant after adjusting for *APOE* genotypes. Our original top hits in *APOC1* remain
458 significant after adjusting for *APOE* genotypes (minimum $p = 3.35 \times 10^{-38}$), which suggests an
459 independent relationship between this gene and lipid traits. This underscores the potential for
460 *APOC1* to be a locus of interest in interaction analyses, with implications for lipid metabolism and
461 associated phenotypes. We note that the initial UKBB analysis identified *APOC1* as the top gene
462 and *LDLR* as the second top gene on Chromosome 19. Upon adjusting for *APOE*, we note that the

463 rankings of the two genes switch; the SNP with the lowest SCAMPI p-value is now rs55791371
464 ($p = 4.11 \times 10^{-48}$), located in an intergenic region near *LDLR*.

465 Computational Performance: We benchmarked the computational performance of
466 SCAMPI across varying sample sizes and numbers of traits for analyzing a single genotype using
467 the High-Performance Computing (HPC) cluster hosted by Emory University Rollins School of
468 Public Health (RSPH), whose infrastructure consists of 25 nodes: twenty-four equipped with 32
469 compute cores and 192GB of RAM, and one outlier with 1.5TB of RAM. We provide average
470 computational run times per genotype in Figure 4. For instance, in our applied analysis of UKB
471 data, SCAMPI processed a single genotype in an average of 20.17 seconds for four lipid-related
472 traits with 300,000 participants. In general, computational run time of SCAMPI increased linearly
473 with sample size and exhibited quadratic growth with the number of traits. While using SCAMPI
474 on the RSPH HPC, we distribute the computational workload into one job array with 1,000
475 simultaneous job instances (1,000 job instances are the maximum allowance per job array on
476 RSPH cluster), which effectively partitions the analysis of 300,000 SNPs into 1000 instances of
477 300 SNPs each. Figure 4 also depicts the number of hours required to complete analyses under
478 various sample sizes and trait quantities by assigning 1,000 job instances on RSPH HPC. Notably,
479 our computational configuration can complete the UKB analysis in approximately 1.68 hours. The
480 figure also shows that processing times grow only modestly with the expansion of the dataset; for
481 instance, a dataset featuring 8 traits and 300,000 samples is estimated to take about 3.98 hours,
482 underscoring SCAMPI's effectiveness for large-scale genetic analyses. Moreover, for the users
483 who are interested in applying SCAMPI to analyze the UKB imputed dataset of over 90 million
484 SNPs, which has approximately 6,000,000 SNPs after QC using the same QC procedure we have
485 discussed in the previous session,⁴⁹ supplemental Figure S9 depicts the number of hours required

486 to complete analyses of 6,000,000 SNPs under various sample sizes and trait quantities by
487 assigning 1,000 job instances on RSPH HPC. Notably, our computational workload configuration
488 can complete the UKBB analysis in approximately 33.62 hours for 6,000,000 SNPs.

489 It should be noted that optimizing the HPC system with a more powerful processing
490 configuration could significantly decrease computational time. Enhancements such as increasing
491 CPU count and expanding storage and memory would contribute to this efficiency. Our evaluation
492 of SCAMPI's computational performance on a single genotype, across various sample sizes and
493 trait numbers, also utilized a MacBook Pro with an Apple M1 chip. This analysis, detailed in
494 Supplemental Figure S10 (a)-(b), mirrors the one in Figure 4 and Figure S9, where SCAMPI
495 processed a single SNP for four lipid-related traits among 300,000 participants in an average of
496 8.65 seconds. An HPC system powered with the M1 chip could presumably and feasibly complete
497 our UKBB analysis, involving 300,000 samples and 4 traits, in just about 0.72 hours. Moreover, it
498 will take 14.41 hours to analyze 6,000,000 SNPs.

499 **Discussion**

500 The observation that narrow-sense heritability estimates of complex traits are often
501 considerably larger when estimated from close relatives than distant relatives points to a potential
502 role of variants with interactive effects on such traits. In this work, we develop our method
503 SCAMPI to help screen for such variants that can then be prioritized for subsequent interaction
504 analyses using standard tools. By studying correlation patterns among multiple traits, we showed
505 using simulated data that SCAMPI has improved power relative to univariate variance-based
506 screening procedures. Like variance-based procedures, SCAMPI does not require the specification
507 of the factor that interacts with the variant to influence the traits under study. This means that users
508 do not need prior knowledge of potential interacting factors, which can often be overlooked,

509 unavailable, or difficult to collect. Furthermore, while SCAMPI produces an omnibus test to assess
510 whether a SNP has an interactive effect on at least one of the traits under study, the method allows
511 a user to identify the specific traits that are driving the signal by inspection of the individual cross-
512 product p-values that are aggregated to form the omnibus test. The method, implemented in R
513 code, is scalable to biobank-scale data and can handle many phenotypes.

514 While we developed SCAMPI with the intent of identifying variants harboring interaction
515 effects with other genetic variants or environments, the method generally detects any variants with
516 non-additive effects, which can also include dominance effects or parent-of-origin effects. To help
517 delineate dominance effects from potential gene-gene or gene-environment effects, one can rerun
518 SCAMPI regressing out the dominance effect of the variant in the DGLM model prior to analysis
519 and observing whether the original interaction signal remains. For parent-of-origin testing, one can
520 recode the SCAMPI regression framework to assess whether the trait correlation among
521 heterozygotes is significantly different from the two homozygote categories.⁸² We note that the
522 appearance of a variant with a possible interaction effect can also arise if the variant is in linkage
523 disequilibrium (LD) with a nearby variant that has a marginal effect on the traits under study.³² In
524 this situation, we suggest identifying such variants with marginal effects in LD with the test variant
525 prior to analysis and regressing the effects of such variants out of the DGLM mean model prior to
526 analysis using SCAMPI.

527 SCAMPI makes a few modeling assumptions that warrant further discussion. By
528 implementing a DGLM model that assumes a Gaussian distribution to standardize traits, the
529 SCAMPI framework inherently assumes the trait values under study follow a multivariate normal
530 distribution. To meet this assumption in the main analysis, we transform the traits to normality
531 using a non-parametric rank-based method, the Inverse Normal Transformation (INT), prior to

532 SCAMPI analysis. We also explored whether transforming the traits before residualizing on the
533 main effects of genotype and confounders (which we refer to as Direct INT or D-INT) led to
534 different inference from transforming after residualizing (which we refer to as Indirect INT or I-
535 INT)⁵³ and found no marked difference in results (see Tables S5-S7). Rather than conducting a
536 rank-based inverse normal transformation, we could also explore trait standardization on the
537 original scale using a different form of a DGLM that assumes the trait outcome follows a gamma
538 distribution. An additional SCAMPI assumption is that the sample size is large enough and the
539 minor allele frequency of the tested variant common enough to enable p-value derivation of the
540 cross-product regression test using asymptotic theory. For SCAMPI analysis of less-common
541 variants in modest sample sizes, we recommend deriving the p-values of the cross-product
542 regression tests using resampling procedures (which randomly shuffle genotypes across subjects)
543 rather than relying on asymptotic theory to ensure valid inference.

544 Our SCAMPI framework complements a recent kernel-based method Latent Interaction
545 Testing (LIT) for interaction testing that used kernel distance covariance techniques to test whether
546 similarity of sample trait correlation patterns correlate with genotype similarity at a test SNP.³⁸
547 SCAMPI has practical features that LIT lacks, including the ability to directly assess which
548 phenotypes among those modeled demonstrate interaction effects (as illustrated in Table 2 and
549 Supplemental Table S3). Additionally, because SCAMPI is based on aggregating results across
550 multiple cross-trait regression tests, it can handle missing data more efficiently than LIT (which
551 requires complete information on all traits for inference). To illustrate, suppose we have a sample
552 where N subjects possess information on two phenotypes while only half of these subjects further
553 possess additional information on a third phenotype. For joint analysis of all 3 phenotypes, LIT
554 only considers the N/2 subjects with complete trait data for inference. SCAMPI, on the other hand,

555 can incorporate the remaining $N/2$ subjects that have only information on phenotypes 1 and 2
556 within its cross-trait statistic. The flexible regression framework that forms the backbone of
557 SCAMPI also enables extensions to perform interaction screening for a variety of other study
558 designs used in genetic projects, including longitudinal and family-based designs. Moreover,
559 SCAMPI can be extended to meta-analysis settings where individual-level data cannot be shared
560 across studies. We will explore these SCAMPI extensions in future work.

561

562 SCAMPI R Package is available for installation on GitHub: [https://github.com/epstein-
564 software/SCAMPI](https://github.com/epstein-
563 software/SCAMPI)

564

565 **Funding Support**

566 This work was supported by NIH grants R01 AG071170 (AJB, SB, DJC, MPE) and R01
567 AG075827 (TSW and APW).

568

569

570

571

572

573

574

575

576

577

578 Reference

- 579 1. Yang, J., Benyamin, B., McEvoy, B.P., Gordon, S., Henders, A.K., Nyholt, D.R., Madden,
580 P.A., Heath, A.C., Martin, N.G., Montgomery, G.W., et al. (2010). Common SNPs
581 explain a large proportion of the heritability for human height. *Nat Genet* 42, 565-569.
- 582 2. Ni, G., Moser, G., Wray, N.R., and Lee, S.H. (2018). Estimation of Genetic Correlation via
583 Linkage Disequilibrium Score Regression and Genomic Restricted Maximum
584 Likelihood. *Am J Hum Genet* 102, 1185-1194.
- 585 3. Wainschtein, P., Jain, D., Zheng, Z., Cupples, L.A., Shadyab, A.H., McKnight, B.,
586 Shoemaker, B.M., Mitchell, B.D., Psaty, B.M., Kooperberg, C., et al. (2022). Assessing
587 the contribution of rare variants to complex trait heritability from whole-genome
588 sequence data. *Nat Genet* 54, 263-273.
- 589 4. Elks, C.E., den Hoed, M., Zhao, J.H., Sharp, S.J., Wareham, N.J., Loos, R.J., and Ong, K.K.
590 (2012). Variability in the heritability of body mass index: a systematic review and meta-
591 regression. *Front Endocrinol (Lausanne)* 3, 29.
- 592 5. Ridge, P.G., Hoyt, K.B., Boehme, K., Mukherjee, S., Crane, P.K., Haines, J.L., Mayeux, R.,
593 Farrer, L.A., Pericak-Vance, M.A., Schellenberg, G.D., et al. (2016). Assessment of the
594 genetic variance of late-onset Alzheimer's disease. *Neurobiol Aging* 41, 200.e213-
595 200.e220.
- 596 6. Seshadri, S., Fitzpatrick, A.L., Ikram, M.A., DeStefano, A.L., Gudnason, V., Boada, M., Bis,
597 J.C., Smith, A.V., Carassquillo, M.M., Lambert, J.C., et al. (2010). Genome-wide
598 analysis of genetic loci associated with Alzheimer disease. *Jama* 303, 1832-1840.
- 599 7. Naj, A.C., and Schellenberg, G.D. (2017). Genomic variants, genes, and pathways of
600 Alzheimer's disease: An overview. *Am J Med Genet B Neuropsychiatr Genet* 174, 5-26.
- 601 8. Naj, A.C., Jun, G., Beecham, G.W., Wang, L.S., Vardarajan, B.N., Buross, J., Gallins, P.J.,
602 Buxbaum, J.D., Jarvik, G.P., Crane, P.K., et al. (2011). Common variants at
603 MS4A4/MS4A6E, CD2AP, CD33 and EPHA1 are associated with late-onset Alzheimer's
604 disease. *Nat Genet* 43, 436-441.
- 605 9. Lambert, J.C., Ibrahim-Verbaas, C.A., Harold, D., Naj, A.C., Sims, R., Bellenguez, C.,
606 DeStafano, A.L., Bis, J.C., Beecham, G.W., Grenier-Boley, B., et al. (2013). Meta-
607 analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer's
608 disease. *Nat Genet* 45, 1452-1458.
- 609 10. Kunkle, B.W., Grenier-Boley, B., Sims, R., Bis, J.C., Damotte, V., Naj, A.C., Boland, A.,
610 Vronskaya, M., van der Lee, S.J., Amlie-Wolf, A., et al. (2019). Genetic meta-analysis of
611 diagnosed Alzheimer's disease identifies new risk loci and implicates A β , tau, immunity
612 and lipid processing. *Nat Genet* 51, 414-430.
- 613 11. Jones, L., Harold, D., and Williams, J. (2010). Genetic evidence for the involvement of lipid
614 metabolism in Alzheimer's disease. *Biochim Biophys Acta* 1801, 754-761.
- 615 12. Hollingworth, P., Harold, D., Sims, R., Gerrish, A., Lambert, J.C., Carrasquillo, M.M.,
616 Abraham, R., Hamshere, M.L., Pahwa, J.S., Moskvina, V., et al. (2011). Common
617 variants at ABCA7, MS4A6A/MS4A4E, EPHA1, CD33 and CD2AP are associated with
618 Alzheimer's disease. *Nat Genet* 43, 429-435.
- 619 13. de la Fuente, J., Grotzinger, A.D., Marioni, R.E., Nivard, M.G., and Tucker-Drob, E.M.
620 (2021). Multivariate Modeling of Direct and Proxy GWAS Indicates Substantial
621 Common Variant Heritability of Alzheimer's Disease. *medRxiv*,
622 2021.2005.2006.21256747.

- 623 14. Keller, M.F., Ferrucci, L., Singleton, A.B., Tienari, P.J., Laaksovirta, H., Restagno, G., Chiò,
624 A., Traynor, B.J., and Nalls, M.A. (2014). Genome-wide analysis of the heritability of
625 amyotrophic lateral sclerosis. *JAMA neurology* 71, 1123-1134.
- 626 15. Visscher, P.M., Brown, M.A., McCarthy, M.I., and Yang, J. (2012). Five years of GWAS
627 discovery. *Am J Hum Genet* 90, 7-24.
- 628 16. Røysamb, E., Moffitt, T.E., Caspi, A., Ystrøm, E., and Nes, R.B. (2023). Worldwide Well-
629 Being: Simulated Twins Reveal Genetic and (Hidden) Environmental Influences.
630 *Perspect Psychol Sci* 18, 1562-1574.
- 631 17. Kempthorne, O. (1955). The Theoretical Values of Correlations between Relatives in
632 Random Mating Populations. *Genetics* 40, 153-167.
- 633 18. Robinson, M.R., English, G., Moser, G., Lloyd-Jones, L.R., Triplett, M.A., Zhu, Z., Nolte,
634 I.M., van Vliet-Ostaptchouk, J.V., Snieder, H., Esko, T., et al. (2017). Genotype-
635 covariate interaction effects and the heritability of adult body mass index. *Nature*
636 *Genetics* 49, 1174-1181.
- 637 19. Binder, E.B., Bradley, R.G., Liu, W., Epstein, M.P., Deveau, T.C., Mercer, K.B., Tang, Y.,
638 Gillespie, C.F., Heim, C.M., Nemeroff, C.B., et al. (2008). Association of FKBP5
639 polymorphisms and childhood abuse with risk of posttraumatic stress disorder symptoms
640 in adults. *Jama* 299, 1291-1305.
- 641 20. Bradley, R.G., Binder, E.B., Epstein, M.P., Tang, Y., Nair, H.P., Liu, W., Gillespie, C.F.,
642 Berg, T., Evces, M., Newport, D.J., et al. (2008). Influence of child abuse on adult
643 depression: moderation by the corticotropin-releasing hormone receptor gene. *Arch Gen*
644 *Psychiatry* 65, 190-200.
- 645 21. Bailey, J.M., Colón-Rodríguez, A., and Atchison, W.D. (2017). Evaluating a Gene-
646 Environment Interaction in Amyotrophic Lateral Sclerosis: Methylmercury Exposure and
647 Mutated SOD1. *Current Environmental Health Reports* 4, 200-207.
- 648 22. Morahan, J.M., Yu, B., Trent, R.J., and Pamphlett, R. (2007). A gene-environment study of
649 the paraoxonase 1 gene and pesticides in amyotrophic lateral sclerosis. *NeuroToxicology*
650 28, 532-540.
- 651 23. Dunn, A.R., O'Connell, K.M.S., and Kaczorowski, C.C. (2019). Gene-by-environment
652 interactions in Alzheimer's disease and Parkinson's disease. *Neuroscience &*
653 *Biobehavioral Reviews* 103, 73-80.
- 654 24. McAllister, K., Mechanic, L.E., Amos, C., Aschard, H., Blair, I.A., Chatterjee, N., Conti, D.,
655 Gauderman, W.J., Hsu, L., Hutter, C.M., et al. (2017). Current Challenges and New
656 Opportunities for Gene-Environment Interaction Studies of Complex Diseases. *American*
657 *Journal of Epidemiology* 186, 753-761.
- 658 25. Hutter, C.M., Mechanic, L.E., Chatterjee, N., Kraft, P., and Gillanders, E.M. (2013). Gene-
659 environment interactions in cancer epidemiology: a National Cancer Institute Think Tank
660 report. *Genet Epidemiol* 37, 643-657.
- 661 26. Spiegelman, D. (2010). Approaches to uncertainty in exposure assessment in environmental
662 epidemiology. *Annual review of public health* 31, 149-163.
- 663 27. Aschard, H., Lutz, S., Maus, B., Duell, E.J., Fingerlin, T.E., Chatterjee, N., Kraft, P., and
664 Van Steen, K. (2012). Challenges and opportunities in genome-wide environmental
665 interaction (GWEL) studies. *Human genetics* 131, 1591-1613.
- 666 28. Lindström, S., Yen, Y.-C., Spiegelman, D., and Kraft, P. (2009). The impact of gene-
667 environment dependence and misclassification in genetic association studies
668 incorporating gene-environment interactions. *Human heredity* 68, 171-181.

- 669 29. Kraft, P., and Aschard, H. (2015). Finding the missing gene–environment interactions.
670 *European journal of epidemiology* 30, 353-355.
- 671 30. Paré, G., Cook, N.R., Ridker, P.M., and Chasman, D.I. (2010). On the use of variance per
672 genotype as a tool to identify quantitative trait interaction effects: a report from the
673 Women's Genome Health Study. *PLoS genetics* 6, e1000981.
- 674 31. Brown, M.B., and Forsythe, A.B. (1974). Robust Tests for the Equality of Variances. *Journal*
675 *of the American Statistical Association* 69, 364-367.
- 676 32. Wang, H., Zhang, F., Zeng, J., Wu, Y., Kemper, K.E., Xue, A., Zhang, M., Powell, J.E.,
677 Goddard, M.E., Wray, N.R., et al. (2019). Genotype-by-environment interactions inferred
678 from genetic effects on phenotypic variability in the UK Biobank. *Science Advances* 5,
679 eaaw3538.
- 680 33. Westerman, K.E., Majarian, T.D., Giulianini, F., Jang, D.-K., Miao, J., Florez, J.C., Chen, H.,
681 Chasman, D.I., Udler, M.S., Manning, A.K., et al. (2022). Variance-quantitative trait loci
682 enable systematic discovery of gene-environment interactions for cardiometabolic serum
683 biomarkers. *Nature Communications* 13, 3993.
- 684 34. Zhu, X., Feng, T., Tayo, B.O., Liang, J., Young, J.H., Franceschini, N., Smith, J.A., Yanek,
685 L.R., Sun, Y.V., Edwards, T.L., et al. (2015). Meta-analysis of correlated traits via
686 summary statistics from GWASs with an application in hypertension. *Am J Hum Genet*
687 96, 21-36.
- 688 35. Turley, P., Walters, R.K., Maghzian, O., Okbay, A., Lee, J.J., Fontana, M.A., Nguyen-Viet,
689 T.A., Wedow, R., Zacher, M., Furlotte, N.A., et al. (2018). Multi-trait analysis of
690 genome-wide association summary statistics using MTAG. *Nat Genet* 50, 229-237.
- 691 36. O'Reilly, P.F., Hoggart, C.J., Pomyen, Y., Calboli, F.C., Elliott, P., Jarvelin, M.R., and Coin,
692 L.J. (2012). MultiPhen: joint model of multiple phenotypes can increase discovery in
693 GWAS. *PLoS One* 7, e34861.
- 694 37. Broadaway, K.A., Cutler, D.J., Duncan, R., Moore, J.L., Ware, E.B., Jhun, M.A., Bielak,
695 L.F., Zhao, W., Smith, J.A., Peyser, P.A., et al. (2016). A Statistical Approach for Testing
696 Cross-Phenotype Effects of Rare Variants. *Am J Hum Genet* 98, 525-540.
- 697 38. Bass, A.J., Bian, S., Wingo, A.P., Wingo, T.S., Cutler, D.J., and Epstein, M.P. (2024).
698 Identifying latent genetic interactions in genome-wide association studies using multiple
699 traits. *Genome Medicine* 16, 62.
- 700 39. Liu, Y., Chen, S., Li, Z., Morrison, A.C., Boerwinkle, E., and Lin, X. (2019). ACAT: a fast
701 and powerful p value combination method for rare-variant analysis in sequencing studies.
702 *The American Journal of Human Genetics* 104, 410-421.
- 703 40. Liu, Y., and Xie, J. (2020). Cauchy Combination Test: A Powerful Test With Analytic p-
704 Value Calculation Under Arbitrary Dependency Structures. *Journal of the American*
705 *Statistical Association* 115, 393-402.
- 706 41. Moore, Camille M., Jacobson, Sean A., and Fingerlin, Tasha E. (2020). Power and Sample
707 Size Calculations for Genetic Association Studies in the Presence of Genetic Model
708 Misspecification. *Human Heredity* 84, 256-271.
- 709 42. Joo, J., Kwak, M., Chen, Z., and Zheng, G. (2010). Efficiency robust statistics for genetic
710 linkage and association studies under genetic model uncertainty. *Statistics in medicine*
711 29, 158-180.
- 712 43. Zheng, G., Freidlin, B., and Gastwirth, J.L. (2006). Comparison of robust tests for genetic
713 association using case-control studies. *Lecture Notes-Monograph Series*, 253-265.

- 714 44. Joo, J., Kwak, M., Ahn, K., and Zheng, G. (2009). A Robust Genome-Wide Scan Statistic of
715 the Wellcome Trust Case–Control Consortium. *Biometrics* 65, 1115-1122.
- 716 45. Lea, A., Subramaniam, M., Ko, A., Lehtimäki, T., Raitoharju, E., Kähönen, M., Seppälä, I.,
717 Mononen, N., Raitakari, O.T., Ala-Korpela, M., et al. (2019). Genetic and environmental
718 perturbations lead to regulatory decoherence. *eLife* 8, e40538.
- 719 46. Musharoff, S., Park, D., Dahl, A., Galanter, J., Liu, X., Huntsman, S., Eng, C., Burchard,
720 E.G., Ayroles, J.F., and Zaitlen, N. (2018). Existence and implications of population
721 variance structure. *bioRxiv*, 439661.
- 722 47. Smyth, G.K. (1989). Generalized linear models with varying dispersion. *Journal of the royal*
723 *statistical society series b-methodological* 51, 47-60.
- 724 48. Murphy, M.D., Fernandes, S.B., Morota, G., and Lipka, A.E. (2022). Assessment of two
725 statistical approaches for variance genome-wide association studies in plants. *Heredity*, 1-
726 10.
- 727 49. Marderstein, A.R., Davenport, E.R., Kulm, S., Van Hout, C.V., Elemento, O., and Clark,
728 A.G. (2021). Leveraging phenotypic variability to identify genetic interactions in human
729 phenotypes. *Am J Hum Genet* 108, 49-67.
- 730 50. Collister, J.A., Liu, X., and Clifton, L. (2022). Calculating Polygenic Risk Scores (PRS) in
731 UK Biobank: A Practical Guide for Epidemiologists. *Front Genet* 13, 818574.
- 732 51. Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L.T., Sharp, K., Motyer, A.,
733 Vukcevic, D., Delaneau, O., and O’Connell, J. (2018). The UK Biobank resource with
734 deep phenotyping and genomic data. *Nature* 562, 203-209.
- 735 52. Wigmore, E.M., Clarke, T.-K., Howard, D., Adams, M., Hall, L., Zeng, Y., Gibson, J.,
736 Davies, G., Fernandez-Pujals, A., and Thomson, P.A. (2017). Do regional brain volumes
737 and major depressive disorder share genetic architecture? A study of Generation Scotland
738 (n= 19 762), UK Biobank (n= 24 048) and the English Longitudinal Study of Ageing (n=
739 5766). *Translational psychiatry* 7, e1205-e1205.
- 740 53. McCaw, Z.R., Lane, J.M., Saxena, R., Redline, S., and Lin, X. (2020). Operating
741 characteristics of the rank-based inverse normal transformation for quantitative trait
742 analysis in genome-wide association studies. *Biometrics* 76, 1262-1272.
- 743 54. Scuteri, A., Sanna, S., Chen, W.-M., Uda, M., Albai, G., Strait, J., Najjar, S., Nagaraja, R.,
744 Orrú, M., and Usala, G. (2007). Genome-wide association scan shows genetic variants in
745 the FTO gene are associated with obesity-related traits. *PLoS genetics* 3, e115.
- 746 55. Barber, M.J., Mangravite, L.M., Hyde, C.L., Chasman, D.I., Smith, J.D., McCarty, C.A., Li,
747 X., Wilke, R.A., Rieder, M.J., and Williams, P.T. (2010). Genome-wide association of
748 lipid-lowering response to statins in combined study populations. *PloS one* 5, e9763.
- 749 56. Cade, B.E., Chen, H., Stilp, A.M., Gleason, K.J., Sofer, T., Ancoli-Israel, S., Arens, R., Bell,
750 G.I., Below, J.E., and Bjornnes, A.C. (2016). Genetic associations with obstructive sleep
751 apnea traits in Hispanic/Latino Americans. *American journal of respiratory and critical*
752 *care medicine* 194, 886-897.
- 753 57. Li, X., Li, Z., Zhou, H., Gaynor, S.M., Liu, Y., Chen, H., Sun, R., Dey, R., Arnett, D.K.,
754 Aslibekyan, S., et al. (2020). Dynamic incorporation of multiple in silico functional
755 annotations empowers rare variant association analysis of large whole-genome
756 sequencing studies at scale. *Nat Genet* 52, 969-983.
- 757 58. Li, X., Quick, C., Zhou, H., Gaynor, S.M., Liu, Y., Chen, H., Selvaraj, M.S., Sun, R., Dey,
758 R., Arnett, D.K., et al. (2023). Powerful, scalable and resource-efficient meta-analysis of

- 759 rare variant associations in large whole genome sequencing studies. *Nat Genet* 55, 154-
760 164.
- 761 59. Li, Z., Li, X., Zhou, H., Gaynor, S.M., Selvaraj, M.S., Arapoglou, T., Quick, C., Liu, Y.,
762 Chen, H., Sun, R., et al. (2022). A framework for detecting noncoding rare-variant
763 associations of large-scale whole-genome sequencing studies. *Nat Methods* 19, 1599-
764 1611.
- 765 60. Ochoa, D., Hercules, A., Carmona, M., Suveges, D., Baker, J., Malangone, C., Lopez, I.,
766 Miranda, A., Cruz-Castillo, C., Fumis, L., et al. (2022). The next-generation Open
767 Targets Platform: reimaged, redesigned, rebuilt. *Nucleic Acids Research* 51, D1353-
768 D1359.
- 769 61. Fuior, E.V., and Gafencu, A.V. (2019). Apolipoprotein C1: Its Pleiotropic Effects in Lipid
770 Metabolism and Beyond. *Int J Mol Sci* 20.
- 771 62. Casula, M., Colpani, O., Xie, S., Catapano, A.L., and Baragetti, A. (2021). HDL in
772 Atherosclerotic Cardiovascular Disease: In Search of a Role. *Cells* 10.
- 773 63. Dullaart, R.P., Perton, F., van der Klauw, M.M., Hillege, H.L., Sluiter, W.J., and Group, P.S.
774 (2010). High plasma lecithin: cholesterol acyltransferase activity does not predict low
775 incidence of cardiovascular events: possible attenuation of cardioprotection associated
776 with high HDL cholesterol. *Atherosclerosis* 208, 537-542.
- 777 64. Mabuchi, H., Nohara, A., and Inazu, A. (2014). Cholesteryl ester transfer protein (CETP)
778 deficiency and CETP inhibitors. *Molecules and cells* 37, 777.
- 779 65. Rousset, X., Vaisman, B., Amar, M., Sethi, A.A., and Remaley, A.T. (2009). Lecithin:
780 cholesterol acyltransferase: from biochemistry to role in cardiovascular disease. *Current*
781 *opinion in endocrinology, diabetes, and obesity* 16, 163.
- 782 66. Deeb, S.S., Zambon, A., Carr, M.C., Ayyobi, A.F., and Brunzell, J.D. (2003). Hepatic lipase
783 and dyslipidemia: interactions among genetic variants, obesity, gender, and diet. *Journal*
784 *of Lipid Research* 44, 1279-1286.
- 785 67. Tan, J., Che, Y., Liu, Y., Hu, J., Wang, W., Hu, L., Zhou, Q., Wang, H., and Li, J. (2021).
786 CELSR2 deficiency suppresses lipid accumulation in hepatocyte by impairing the UPR
787 and elevating ROS level. *The FASEB Journal* 35, e21908.
- 788 68. Teslovich, T.M., Musunuru, K., Smith, A.V., Edmondson, A.C., Stylianou, I.M., Koseki, M.,
789 Pirruccello, J.P., Ripatti, S., Chasman, D.I., Willer, C.J., et al. (2010). Biological, clinical
790 and population relevance of 95 loci for blood lipids. *Nature* 466, 707-713.
- 791 69. Willer, C.J., Schmidt, E.M., Sengupta, S., Peloso, G.M., Gustafsson, S., Kanoni, S., Ganna,
792 A., Chen, J., Buchkovich, M.L., Mora, S., et al. (2013). Discovery and refinement of loci
793 associated with lipid levels. *Nat Genet* 45, 1274-1283.
- 794 70. Barton, A.R., Sherman, M.A., Mukamel, R.E., and Loh, P.-R. (2021). Whole-exome
795 imputation within UK Biobank powers rare coding variant association and fine-mapping
796 analyses. *Nat Genet* 53, 1260-1269.
- 797 71. Prins, B.P., Kuchenbaecker, K.B., Bao, Y., Smart, M., Zabaneh, D., Fatemifar, G., Luan, J.a.,
798 Wareham, N.J., Scott, R.A., Perry, J.R.B., et al. (2017). Genome-wide analysis of health-
799 related biomarkers in the UK Household Longitudinal Study reveals novel associations.
800 In *Scientific reports*. p 11008.
- 801 72. Gorden, A., Yang, R., Yerges-Armstrong, L.M., Ryan, K.A., Speliotes, E., Borecki, I.B.,
802 Harris, T.B., Chu, X., Wood, G.C., Still, C.D., et al. (2013). Genetic Variation at NCAN
803 Locus Is Associated with Inflammation and Fibrosis in Non-Alcoholic Fatty Liver
804 Disease in Morbid Obesity. *Human Heredity* 75, 34-43.

- 805 73. Li, X.Y., Liu, Z., Li, L., Wang, H.J., and Wang, H. (2022). TM6SF2 rs58542926 is related to
806 hepatic steatosis, fibrosis and serum lipids both in adults and children: A meta-analysis.
807 *Front Endocrinol (Lausanne)* 13, 1026901.
- 808 74. Ke, P., Xu, M., Feng, J., Tian, Q., He, Y., Lu, K., and Lu, Z. (2023). Association between
809 weight change and risk of liver fibrosis in adults with type 2 diabetes. *J Glob Health* 13,
810 04138.
- 811 75. Hamułka, J., Górnicka, M., Sulich, A., and Frackiewicz, J. (2019). Weight loss program is
812 associated with decrease α -tocopherol status in obese adults. *Clinical Nutrition* 38, 1861-
813 1870.
- 814 76. Miao, L., Yin, R.X., Pan, S.L., Yang, S., Yang, D.Z., and Lin, W.X. (2018). BCL3-PVRL2-
815 TOMM40 SNPs, gene-gene and gene-environment interactions on dyslipidemia. *Sci Rep*
816 8, 6189.
- 817 77. Urbina, E.M., and Daniels, S.R. (2008). Chapter 14 - Hyperlipidemia. In *Adolescent*
818 *Medicine*, G.B. Slap, ed. (Philadelphia, Mosby), pp 90-96.
- 819 78. Zhou, X., Chen, Y., Mok, K.Y., Kwok, T.C.Y., Mok, V.C.T., Guo, Q., Ip, F.C., Chen, Y.,
820 Mullapudi, N., Giusti-Rodríguez, P., et al. (2019). Non-coding variability at the APOE
821 locus contributes to the Alzheimer's risk. *Nat Commun* 10, 3310.
- 822 79. Jia, L., Li, F., Wei, C., Zhu, M., Qu, Q., Qin, W., Tang, Y., Shen, L., Wang, Y., Shen, L., et
823 al. (2020). Prediction of Alzheimer's disease using multi-variants from a Chinese
824 genome-wide association study. *Brain* 144, 924-937.
- 825 80. Button, E.B., Robert, J., Caffrey, T.M., Fan, J., Zhao, W., and Wellington, C.L. (2019). HDL
826 from an Alzheimer's disease perspective. *Curr Opin Lipidol* 30, 224-234.
- 827 81. Chasman, D.I., Kozlowski, P., Zee, R.Y., Kwiatkowski, D.J., and Ridker, P.M. (2006).
828 Qualitative and quantitative effects of APOE genetic variation on plasma C-reactive
829 protein, LDL-cholesterol, and apoE protein. *Genes & Immunity* 7, 211-219.
- 830 82. Head, S.T., Leslie, E.J., Cutler, D.J., and Epstein, M.P. (2023). POIROT: a powerful test for
831 parent-of-origin effects in unrelated samples leveraging multiple phenotypes.
832 *Bioinformatics* 39.
833

834

835

836

837

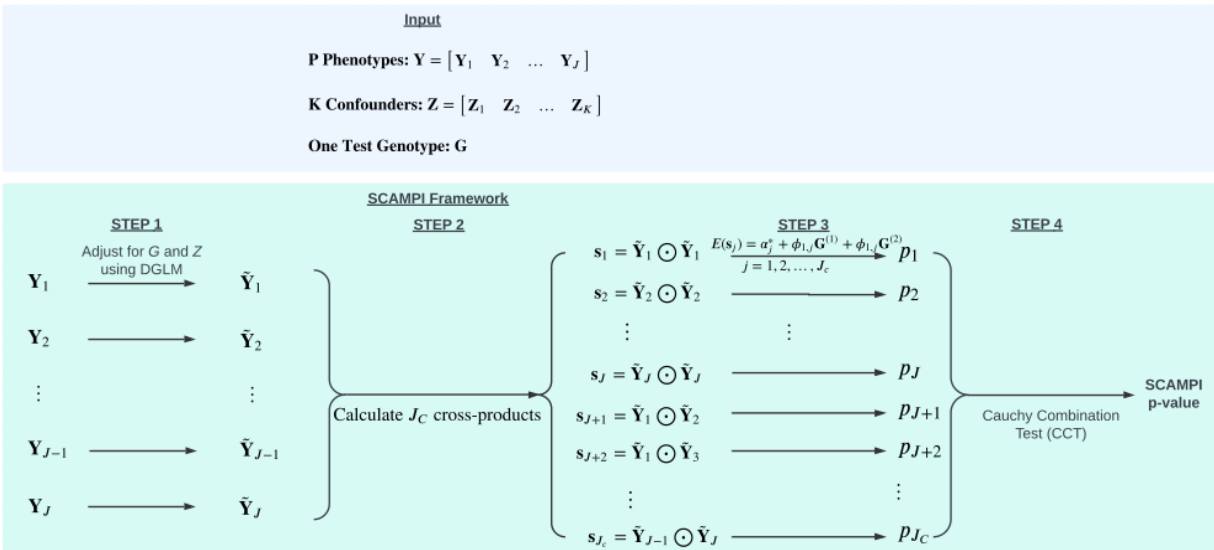
838

839

840

841

842 **Figure 1**



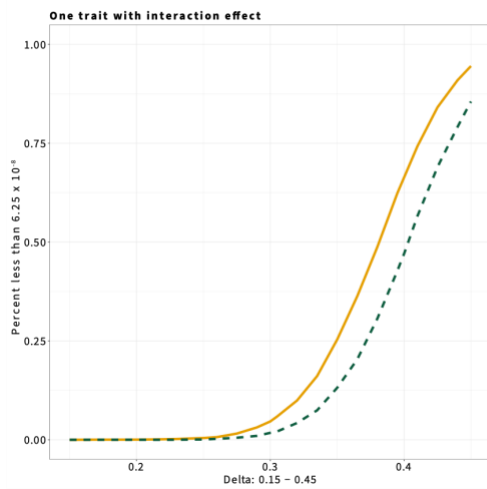
843

844 **Figure 1. Illustration of the SCAMPI framework.**

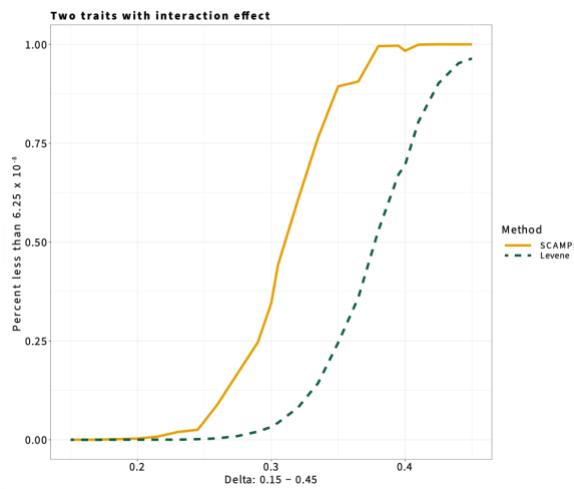
845 The SCAMPI framework involves a four-step process that consists of 1) adjustment of phenotypes
 846 for genotype and confounders, 2) calculation of cross products from adjusted phenotypes, 3)
 847 derivation of p-values from regression tests of cross products on test genotype, and 4) aggregating
 848 all p-values using the Cauchy Combination Test (CCT) to derive the final SCAMPI p-value to
 849 determine overall significance.

850
 851
 852
 853
 854
 855
 856
 857
 858
 859
 860
 861
 862
 863
 864
 865

866 **Figure 2**
867 (a)

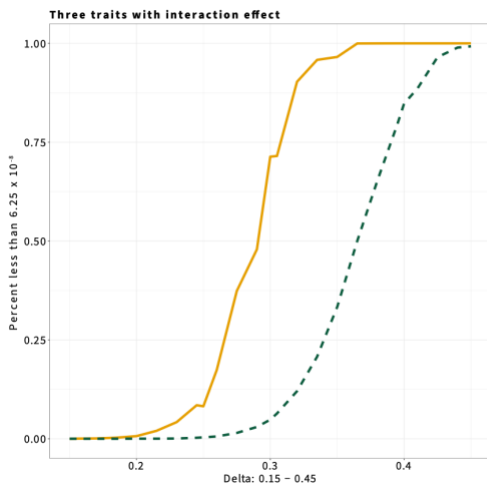


(b)

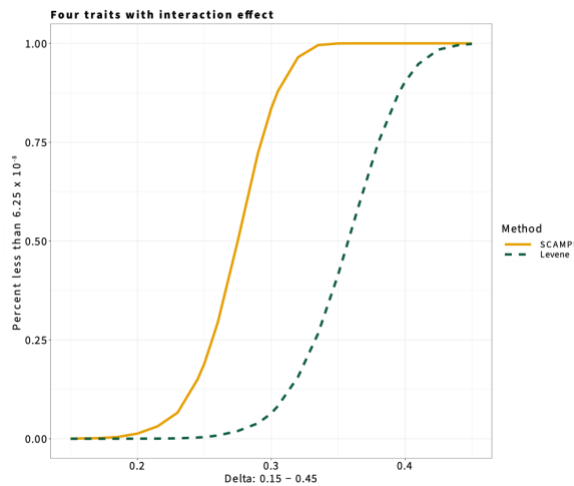


868

869 (c)



(d)

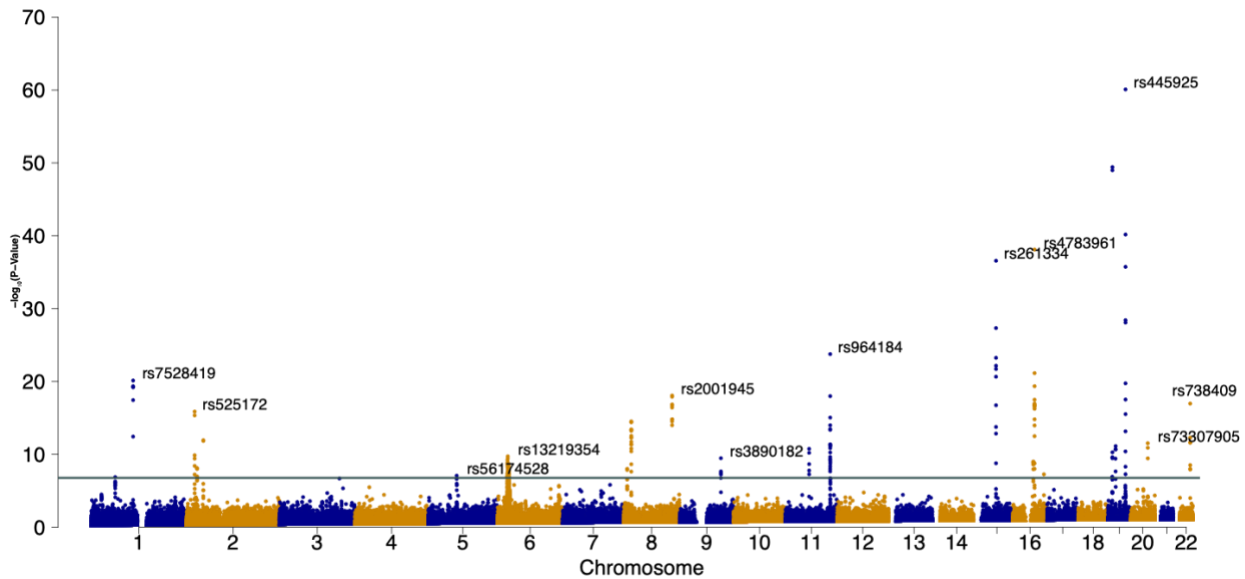


870

871 **Figure 2. UKBB-inspired power simulation of SCAMPI for four traits**

872 Power of SCAMPI at $\alpha = 6.25 \times 10^{-8}$ for four traits at a sample size of 20,000 with MAF=0.05
873 and $\gamma = 0.05$. The correlation among the four traits is inspired by correlation among lipid traits
874 considered in our applied UKBB analysis. Yellow solid line represents power of SCAMPI, while
875 dashed green line denotes power of the benchmark Multivariate Levene's test. Sub-figures (a) –
876 (d) examines power when interaction effects exist for one trait, two traits, three traits and four
877 traits, respectively. We analyzed 10,000 replicates under each model.

878 **Figure 3**



879

880 **Figure 3. Genome-wide results on lipid traits in UKBB using SCAMPI**

881 SCAMPI results for detecting latent interaction effects on high-density lipoprotein cholesterol
882 (HDL-C), low-density lipoproteins cholesterol (LDL-C), triglycerides (TGs), and Body Mass
883 Index (BMI). After SNP QC, 288,910 SNPs are included in the analysis with their MAF ≥ 0.05 .
884 SCAMPI successfully identified 210 SNPs from 68 genes and intergenic regions at the pre-
885 specified study-wide significance ($\alpha = 1.67 \times 10^{-7}$) represented by the green solid horizontal
886 line. The SNP with the smallest SCAMPI p-value is rs445925 located on *ApoC1*.

887

888

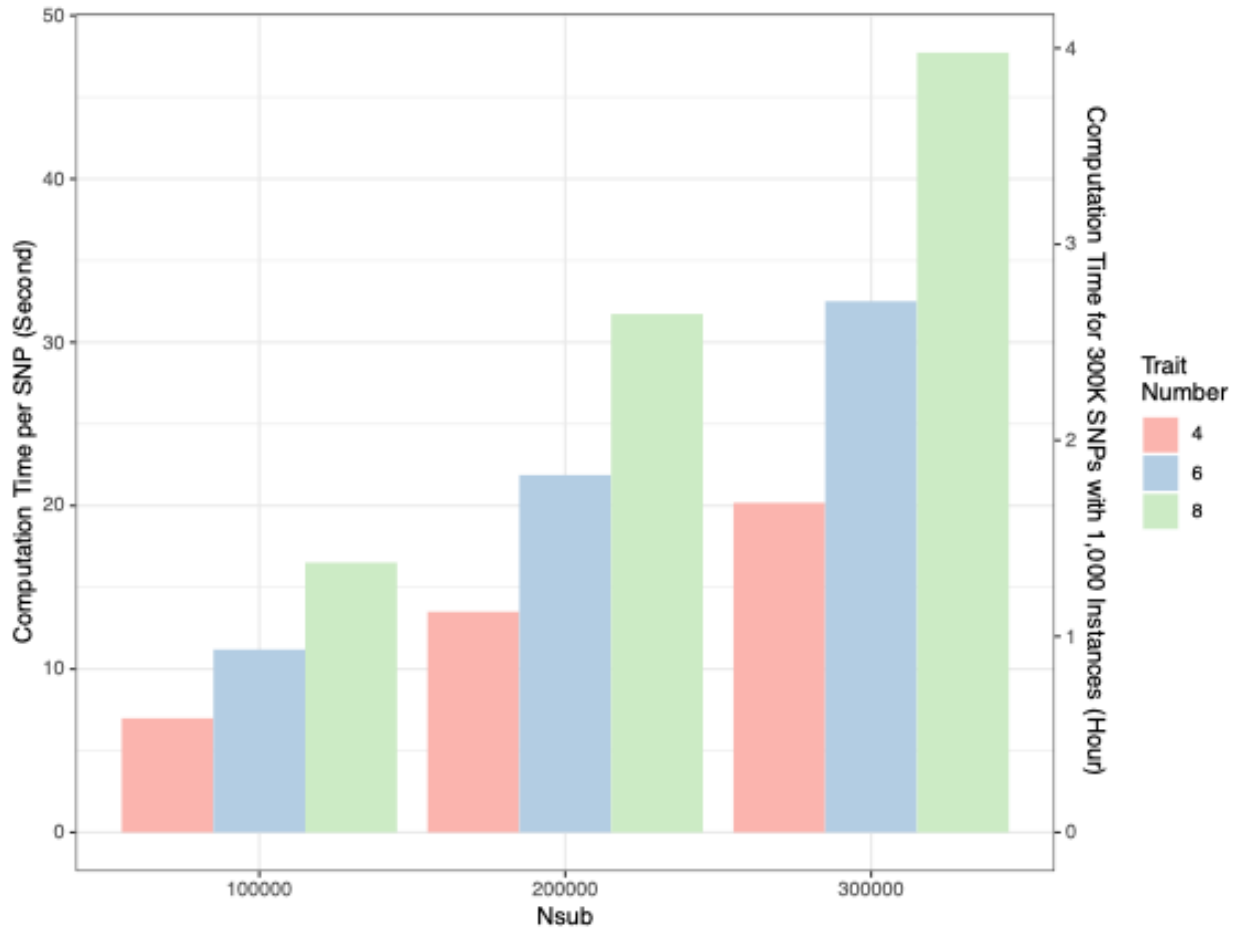
889

890

891

892

893 **Figure 4**



894

895 **Figure 4. Computational performance of SCAMPI**

896 Computational run time of SCAMPI for different sample sizes and number of traits using High-
897 Performance Computing (HPC) cluster hosted by Emory University Rollins School of Public
898 Health (RSPH). Computational run time is based on average of 1,000 simulations for different
899 scenarios with varying trait number and sample size. The first y-axis in Figure 4 displays the time
900 in seconds to complete SCAMPI for one SNP at different configurations. The second y-axis shows
901 the hours required to complete analyses assuming 1,000 job instances on a high-performance
902 cluster.

903

904 **Table 1**

N	MAF	Σ	<i>J</i>	LEQ 0.01	LEQ 0.001
3.00E+05	0.05	0.01	2	9.56E-03	1.01E-03
			4	1.00E-02	1.18E-03
			8	1.05E-02	1.13E-03
3.00E+05	0.05	0.25	2	9.96E-03	1.09E-03
			4	1.09E-02	9.30E-04
			8	1.12E-02	1.01E-03
3.00E+05	0.05	0.5	2	1.05E-02	1.16E-03
			4	1.16E-02	1.02E-03
			8	1.26E-02	1.23E-03
3.00E+05	0.25	0.01	2	9.51E-03	1.17E-03
			4	1.03E-02	1.13E-03
			8	1.00E-02	1.14E-03
3.00E+05	0.25	0.25	2	1.01E-02	1.12E-03
			4	1.09E-02	1.06E-03
			8	1.09E-02	9.40E-04
3.00E+05	0.25	0.5	2	1.03E-02	9.80E-04
			4	1.09E-02	1.03E-03
			8	1.28E-02	1.09E-03

905

906 **Table 1. Nominal rate of empirical type 1 error rates for SCAMPI.**

907 The empirical type I error rates for SCAMPI at nominal rates α of 10^{-2} and 10^{-3} in 100,000
 908 simulations with 300,000 observations. The result is presented across a range of conditions
 909 including varying Minor Allele Frequencies (MAF), numbers of phenotypes (*J*), and covariance
 910 of the phenotypes Σ when γ is simulated from a uniform distribution between 0 and 0.3. The value
 911 presented in the 'LEQ 0.01' and 'LEQ 0.001' columns reflect the pre-specified nominal error rates.

912

913 **Table 2**

Chr	Pos	Alt	Ref	RS #	Gene	SCAMPI P-value	Significant Variance/Correlation Components	PheWAS
1	109817192	G	A	rs7528419	<i>CELSR2</i>	7.33E-21	Corr(TRIG, LDL), Var(LDL)	TRIG, LDL
2	21382976	G	T	rs525172	Intergenic	1.33E-16	Corr(TRIG, HDL), Var(LDL)	TRIG, LDL
5	74400516	C	G	rs56174528	<i>ANKRD31</i>	8.26E-08	Var(LDL)	LDL
6	27185664	C	T	rs13219354	<i>PRSS16</i>	1.83E-10	Var(BMI)	BMI
8	126477978	C	G	rs2001945	(<i>TRIB1</i>)	8.34E-19	Var(LDL)	LDL
9	107647655	A	G	rs3890182	<i>ABCA1</i>	3.34E-10	Var(HDL)	HDL
11	116648917	C	G	rs964184	<i>ZPR1</i>	1.73E-24	Corr(TRIG, LDL), Corr(TRIG, HDL), Corr(LDL, HDL), Corr(LDL, BMI), Var(Trig), Var(LDL)	TRIG, LDL, HDL
15	58726744	C	G	rs261334	<i>LIPC</i> ; <i>LIPC-AS1</i>	2.63E-37	Corr(HDL, BMI), Var(TRIG)	TRIG, HDL
16	56994894	A	G	rs4783961	<i>CETP</i>	7.61E-39	Var(HDL)	HDL
19	45415640	A	G	rs445925	<i>APOC1</i>	8.10E-61	Corr(TRIG, LDL), Corr(TRIG, HDL), Corr(LDL, HDL), Corr(LDL, BMI), Var(Trig), Var(LDL), Var(HDL)	TRIG, LDL, HDL
20	44545773	C	A	rs73307905	(<i>PLTP</i>)	2.90E-12	Var(HDL)	HDL
22	44324727	G	C	rs738409	<i>PNPLA3</i>	1.07E-17	Corr(TRIG, BMI), Var(TRIG)	BMI

914

915 **Table 2. The lead SNPs, identified by SCAMPI within each chromosome, implies interaction effects for the four lipid traits in**916 **UKBB**

917 SCAMPI identified 12 lead SNPs from 12 chromosomes. The position of the SNPs is based on the Genome Reference Consortium

918 Human Build 37 (GRCh37). Column “Gene” indicates the gene where the SNP locates. Column “SCAMPI P-value” shows the SCAMPI

919 p-value. Column “Significant Variance/Correlation Components” indicates the variance or the correlation components of the four lipids

920 that are significantly associated with the corresponding SNP at the pre-specified study-wide significance ($\alpha = 1.67 \times 10^{-7}$). Column
921 “PheWAS” lists the traits involved in the significant variance and correlation components as noted in column “Significant
922 Variance/Correlation Components”, and these traits are also identified to be significant in PheWAS results, which is cross-referenced
923 based on the GWAS Catalog or UK Biobank from the Open Targets Platform.

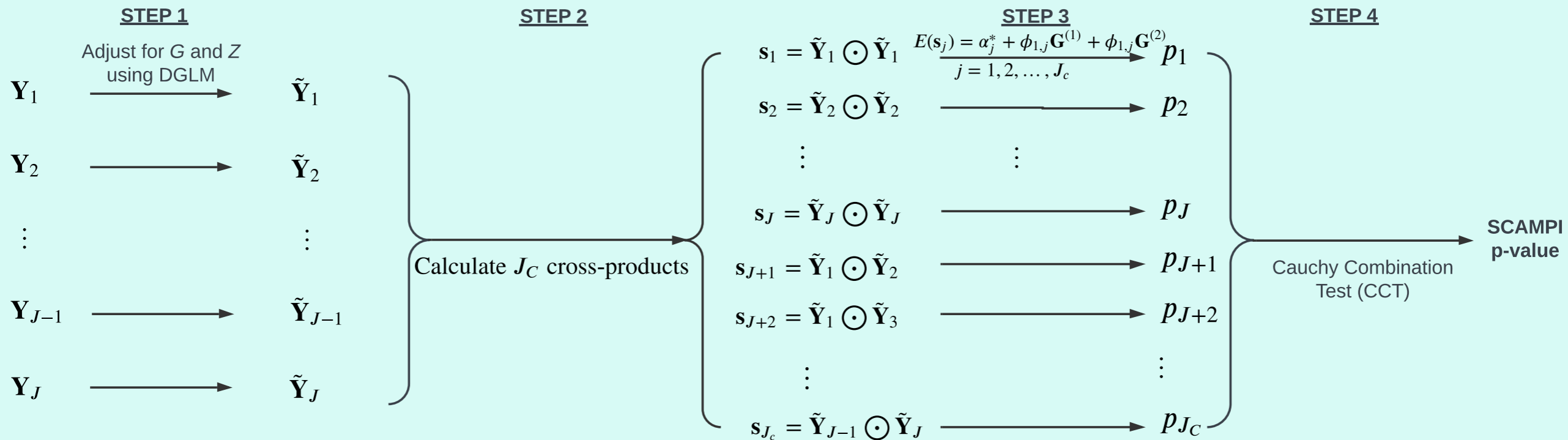
Input

P Phenotypes: $\mathbf{Y} = [\mathbf{Y}_1 \quad \mathbf{Y}_2 \quad \dots \quad \mathbf{Y}_J]$

K Confounders: $\mathbf{Z} = [\mathbf{Z}_1 \quad \mathbf{Z}_2 \quad \dots \quad \mathbf{Z}_K]$

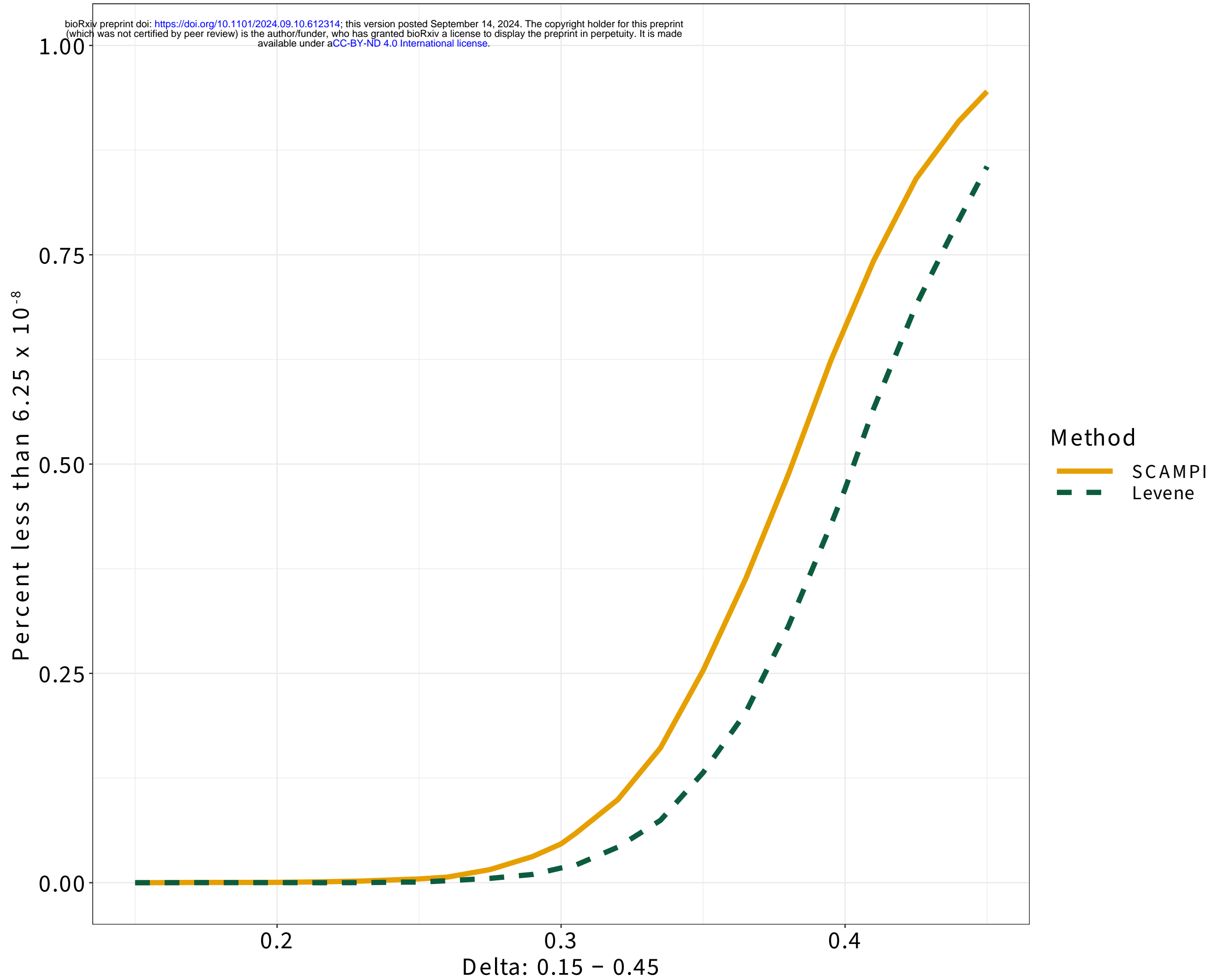
One Test Genotype: \mathbf{G}

SCAMPI Framework



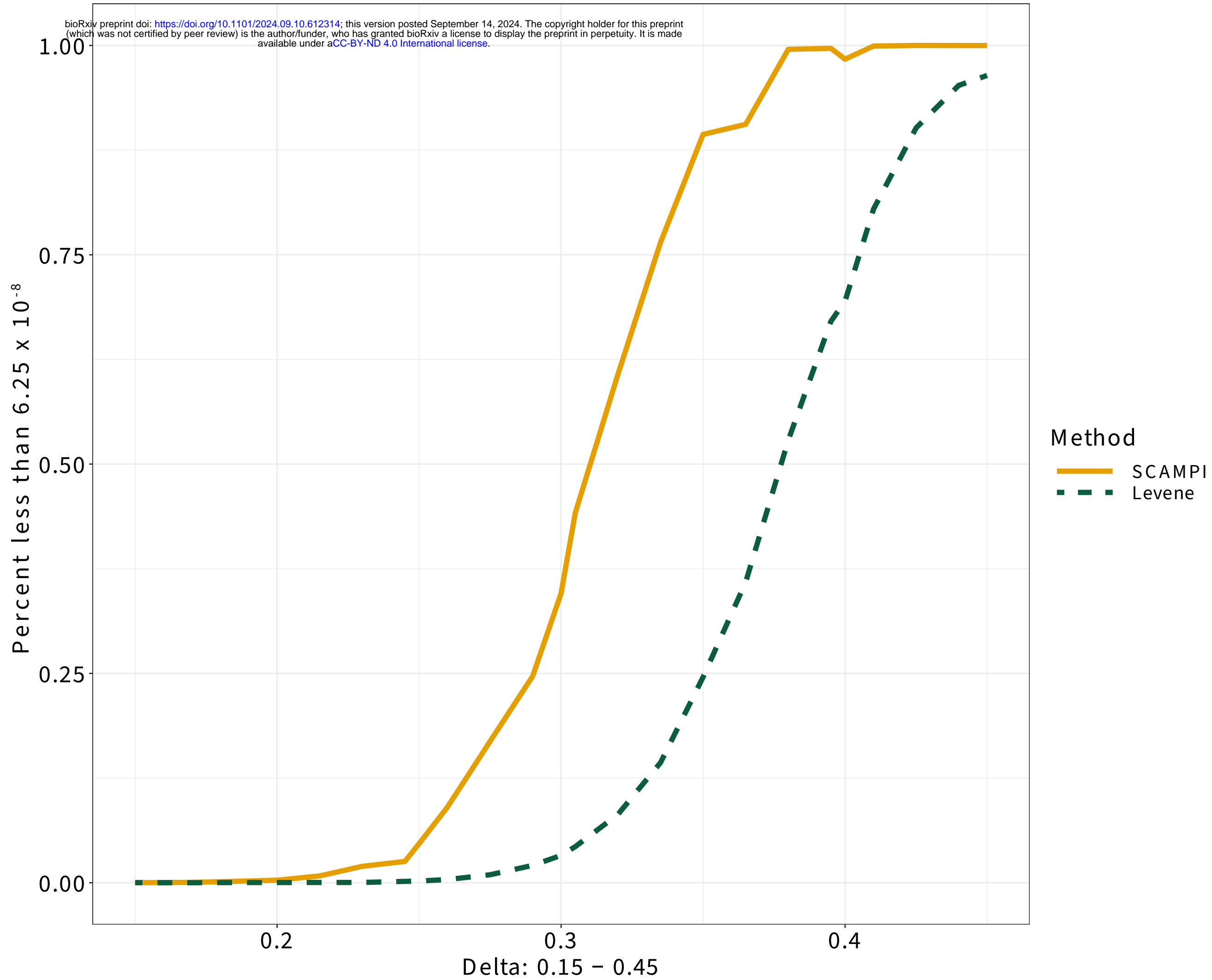
One trait with interaction effect

bioRxiv preprint doi: <https://doi.org/10.1101/2024.09.10.612314>; this version posted September 14, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).



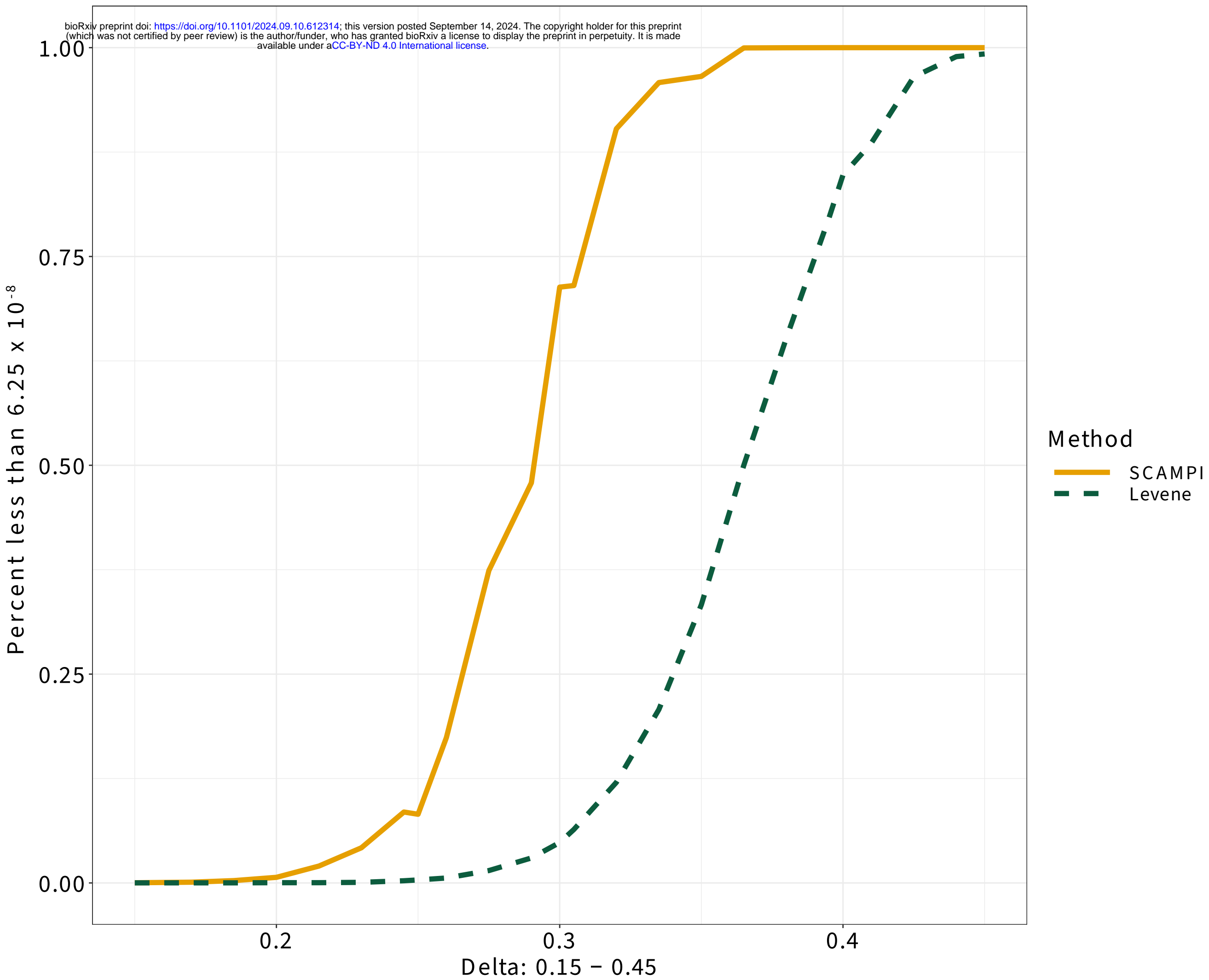
Two traits with interaction effect

bioRxiv preprint doi: <https://doi.org/10.1101/2024.09.10.612314>; this version posted September 14, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).



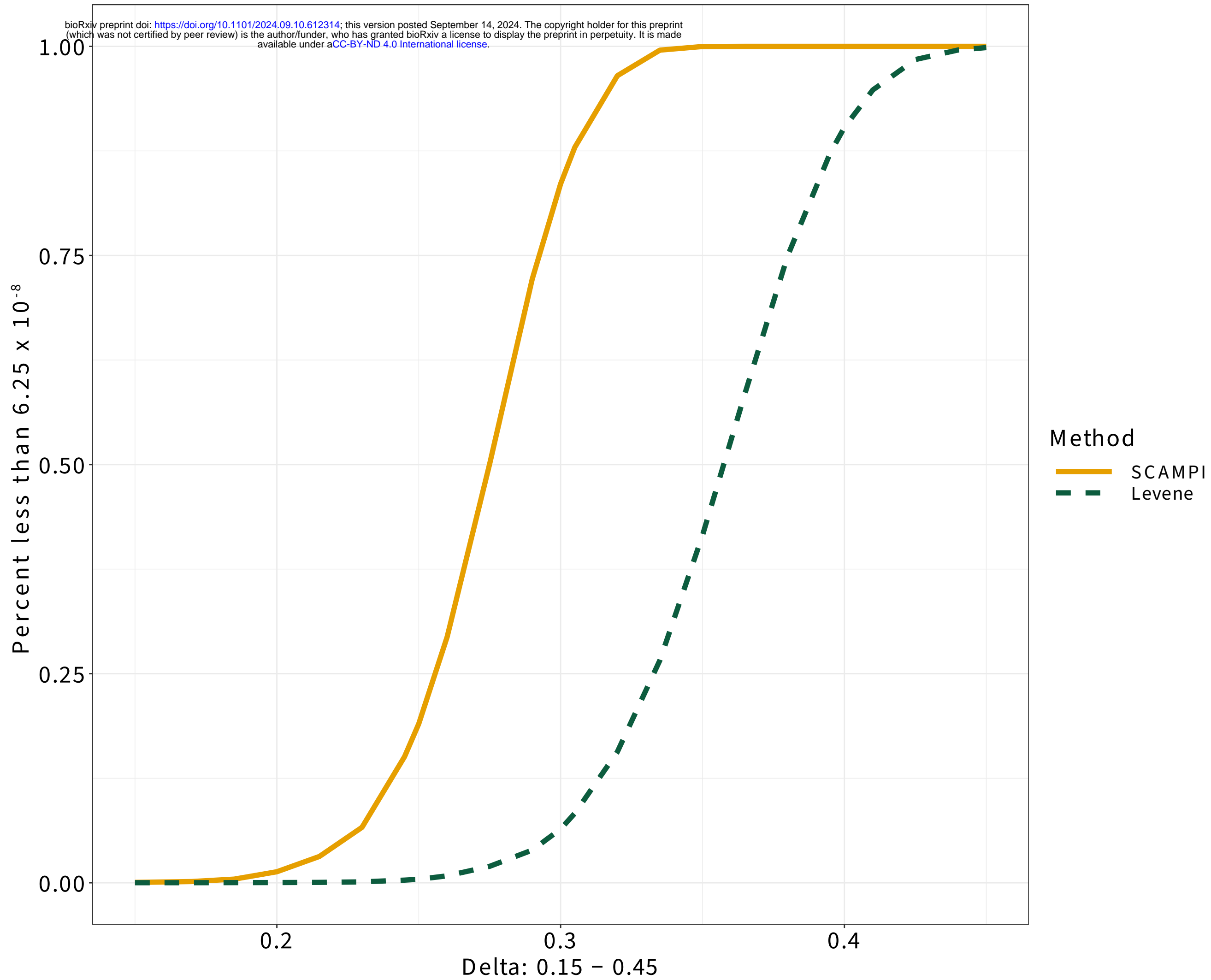
Three traits with interaction effect

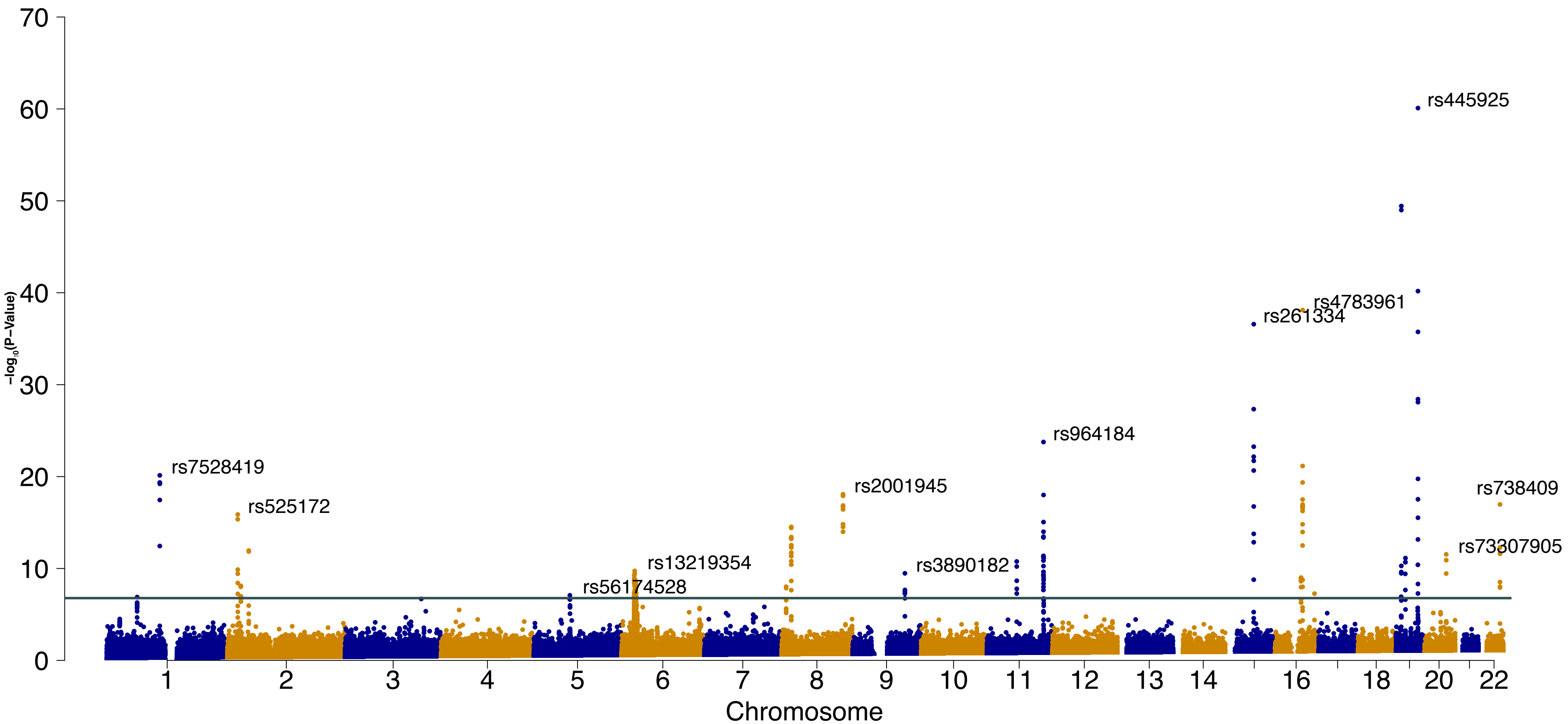
bioRxiv preprint doi: <https://doi.org/10.1101/2024.09.10.612314>; this version posted September 14, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

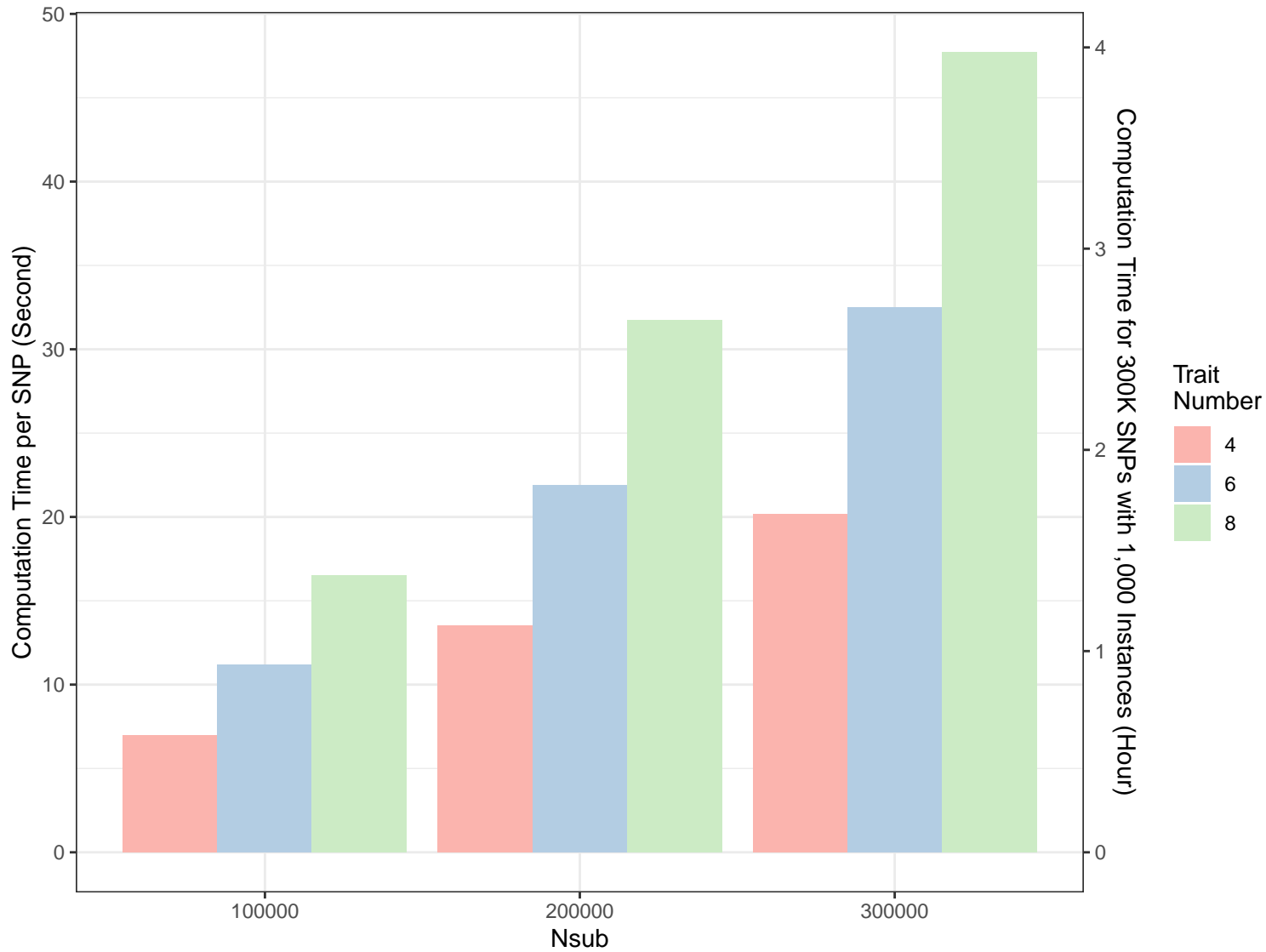


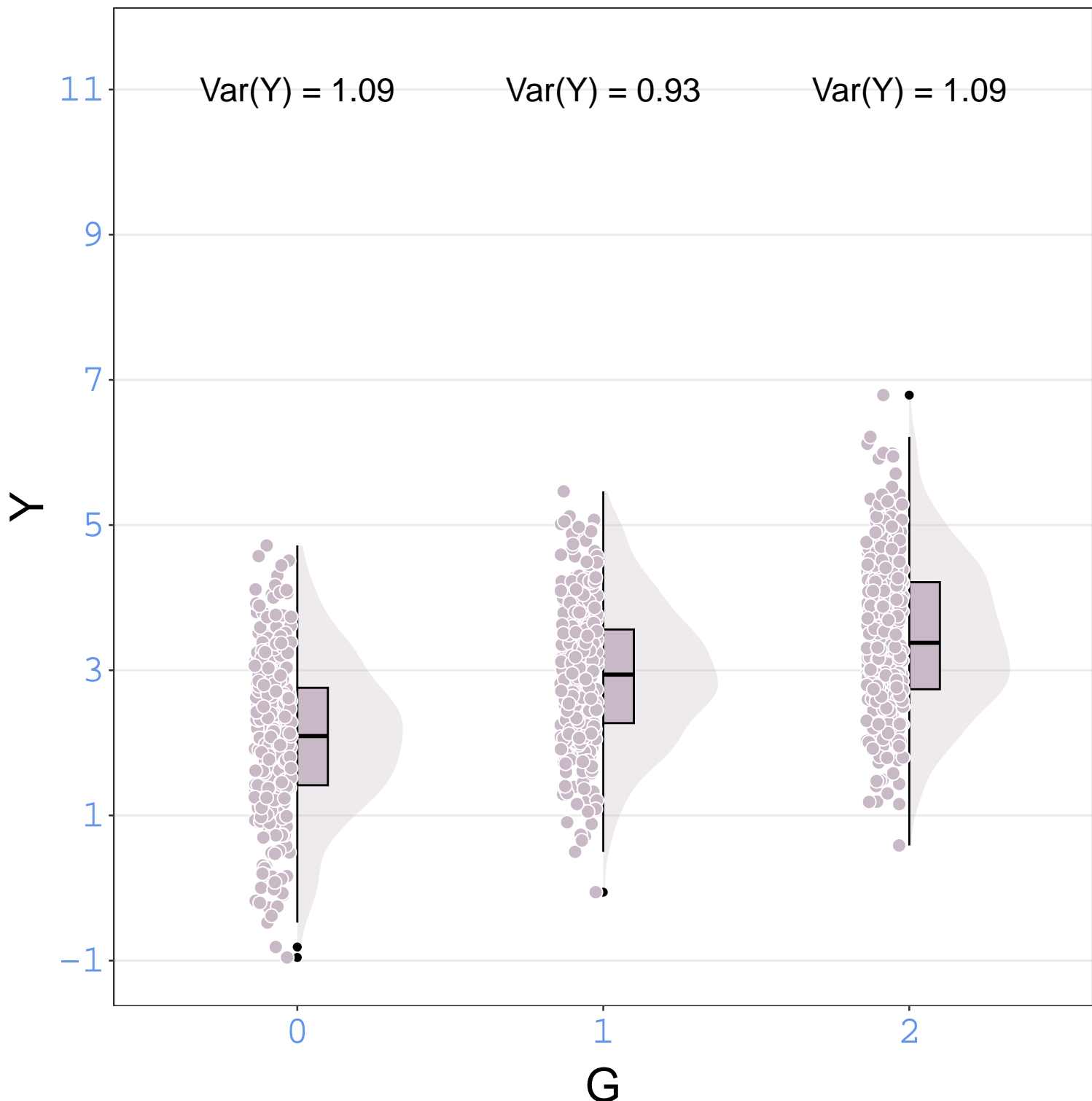
Four traits with interaction effect

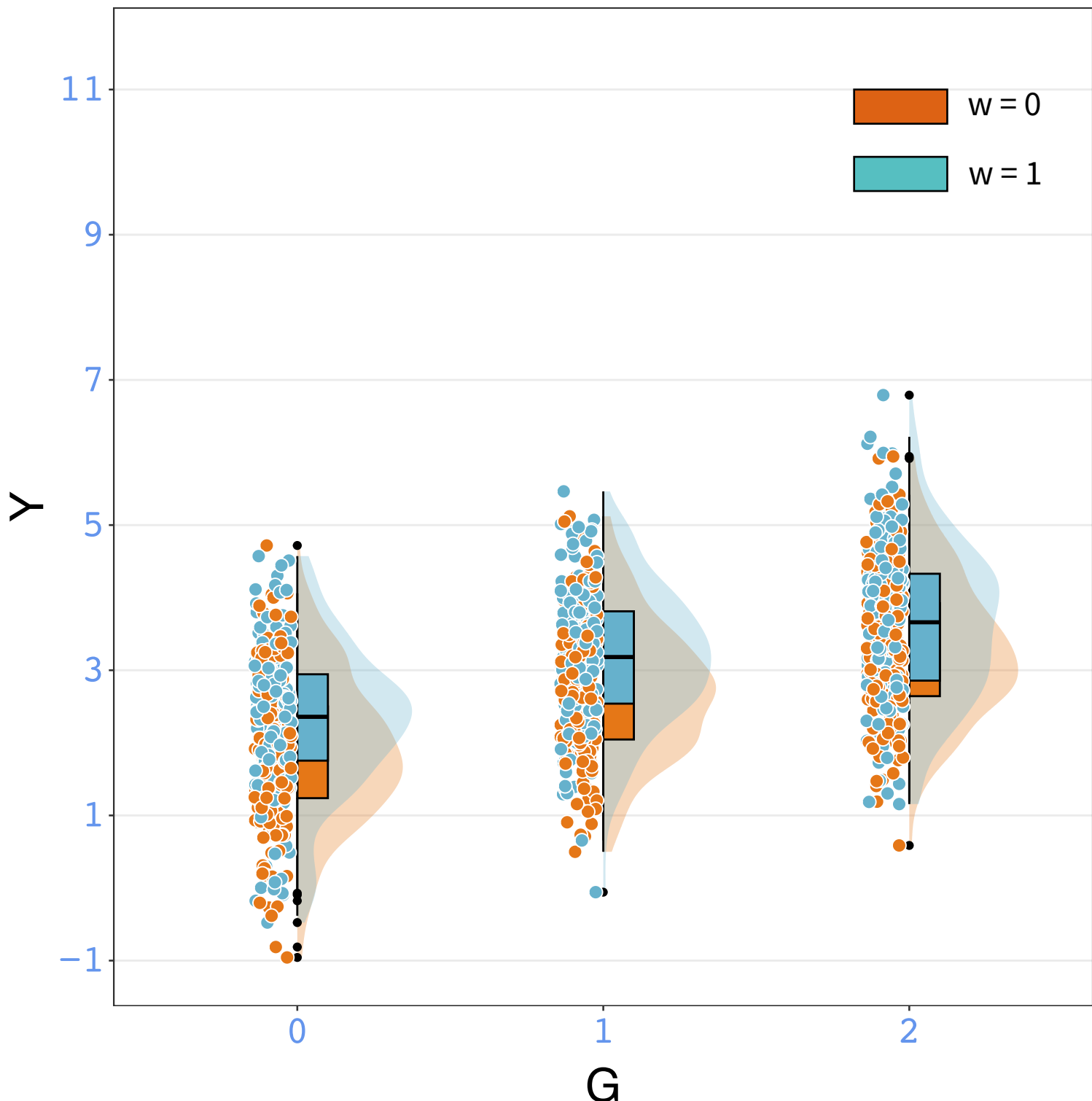
bioRxiv preprint doi: <https://doi.org/10.1101/2024.09.10.612314>; this version posted September 14, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

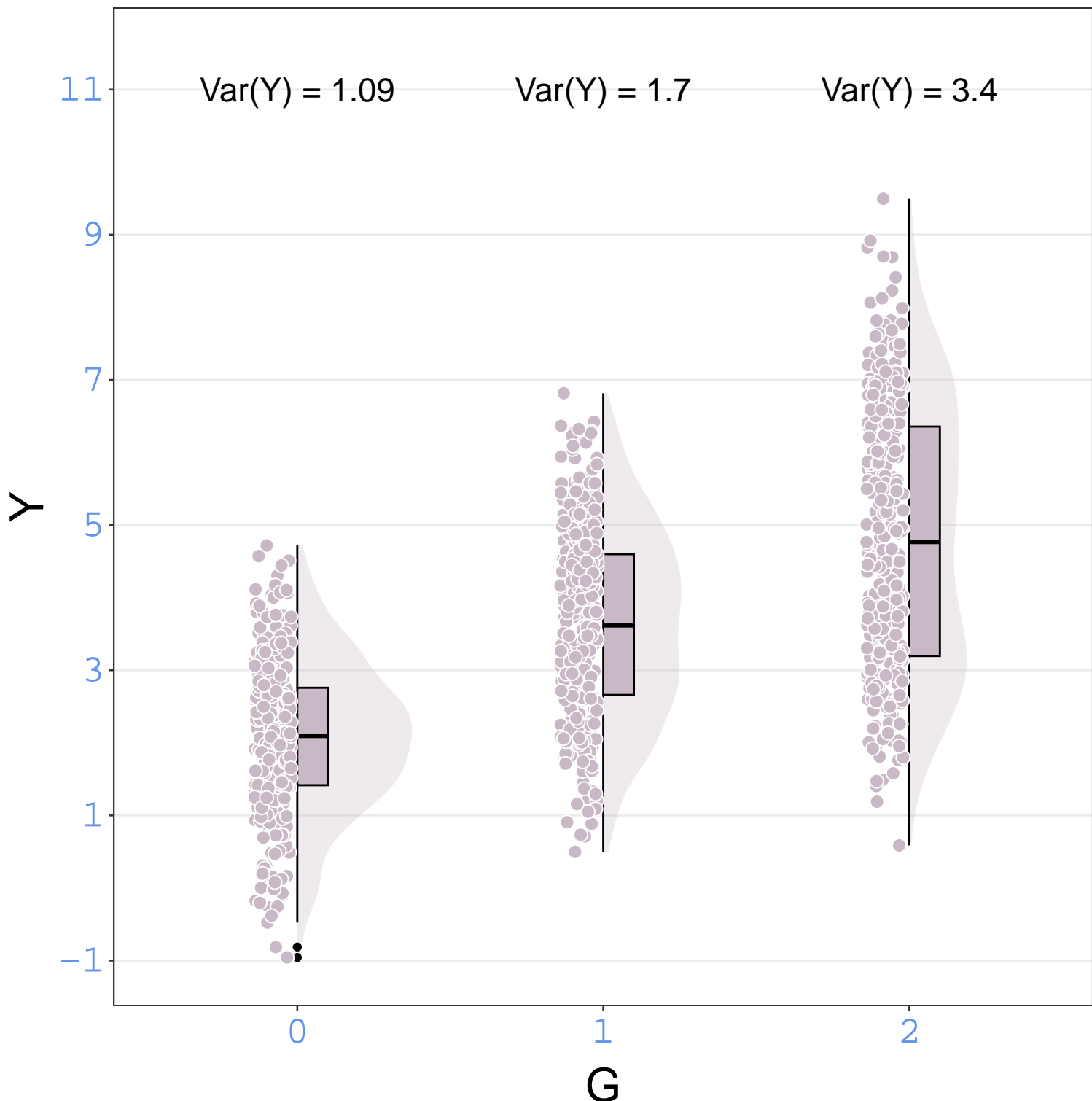


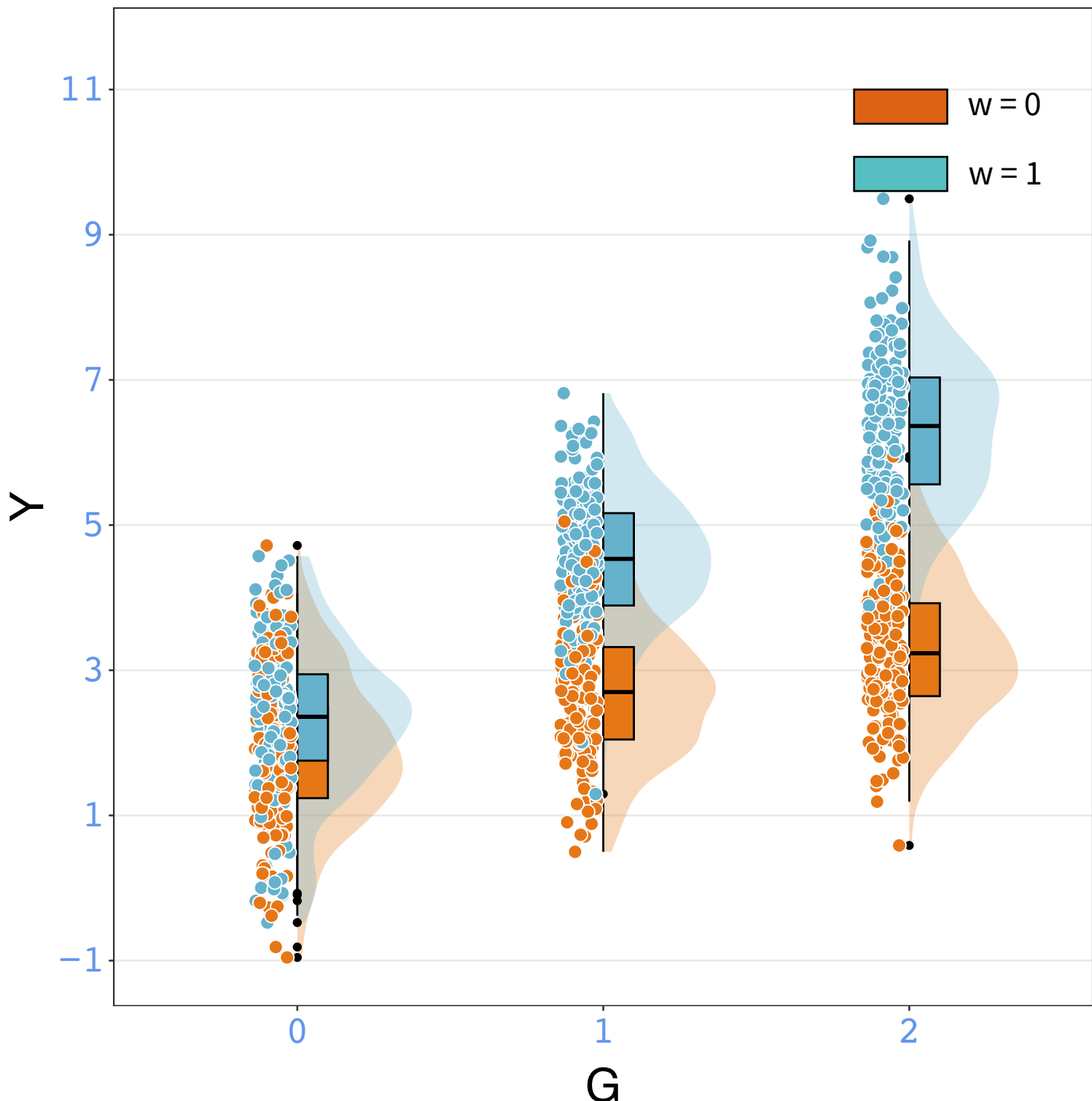






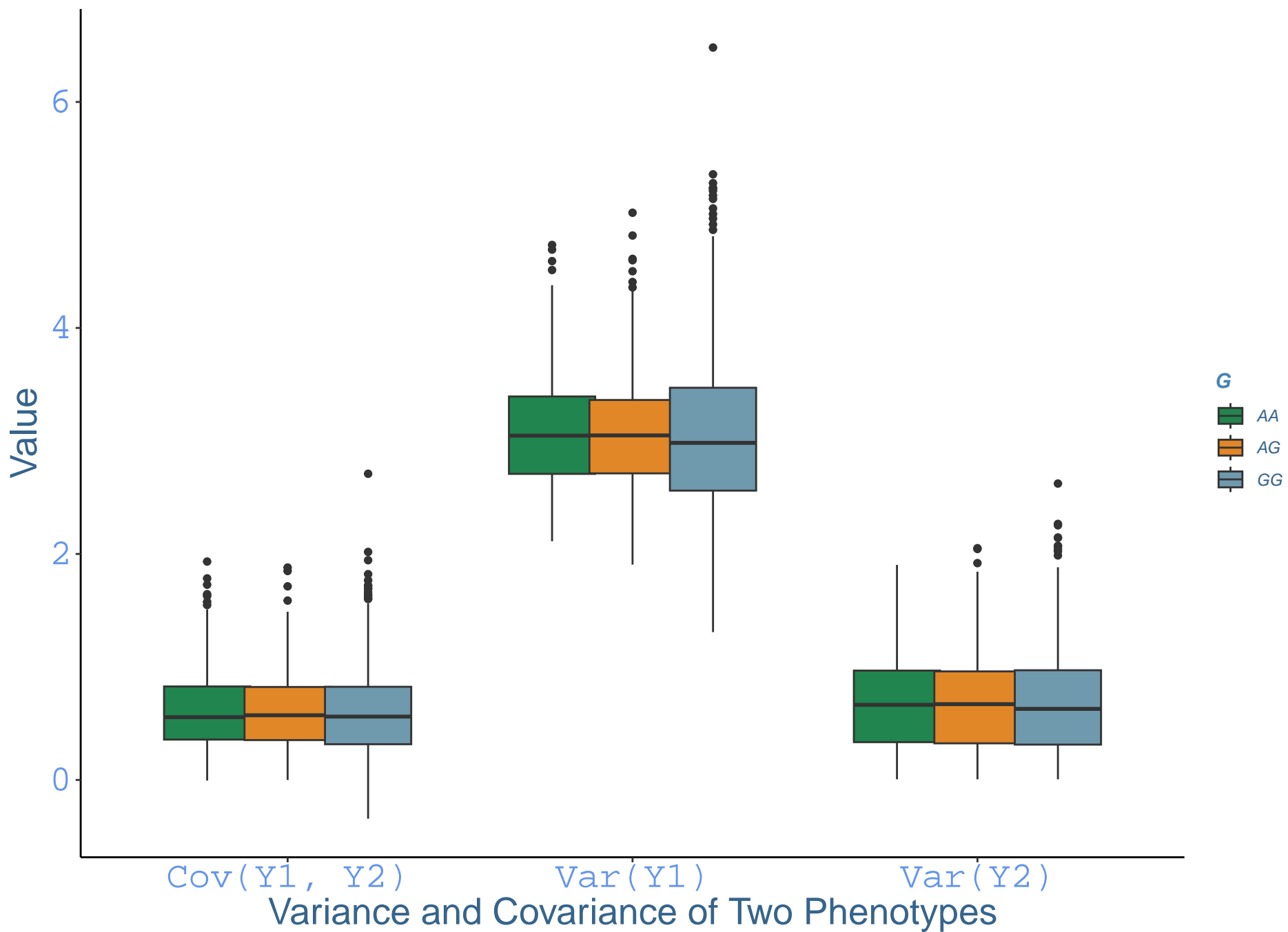






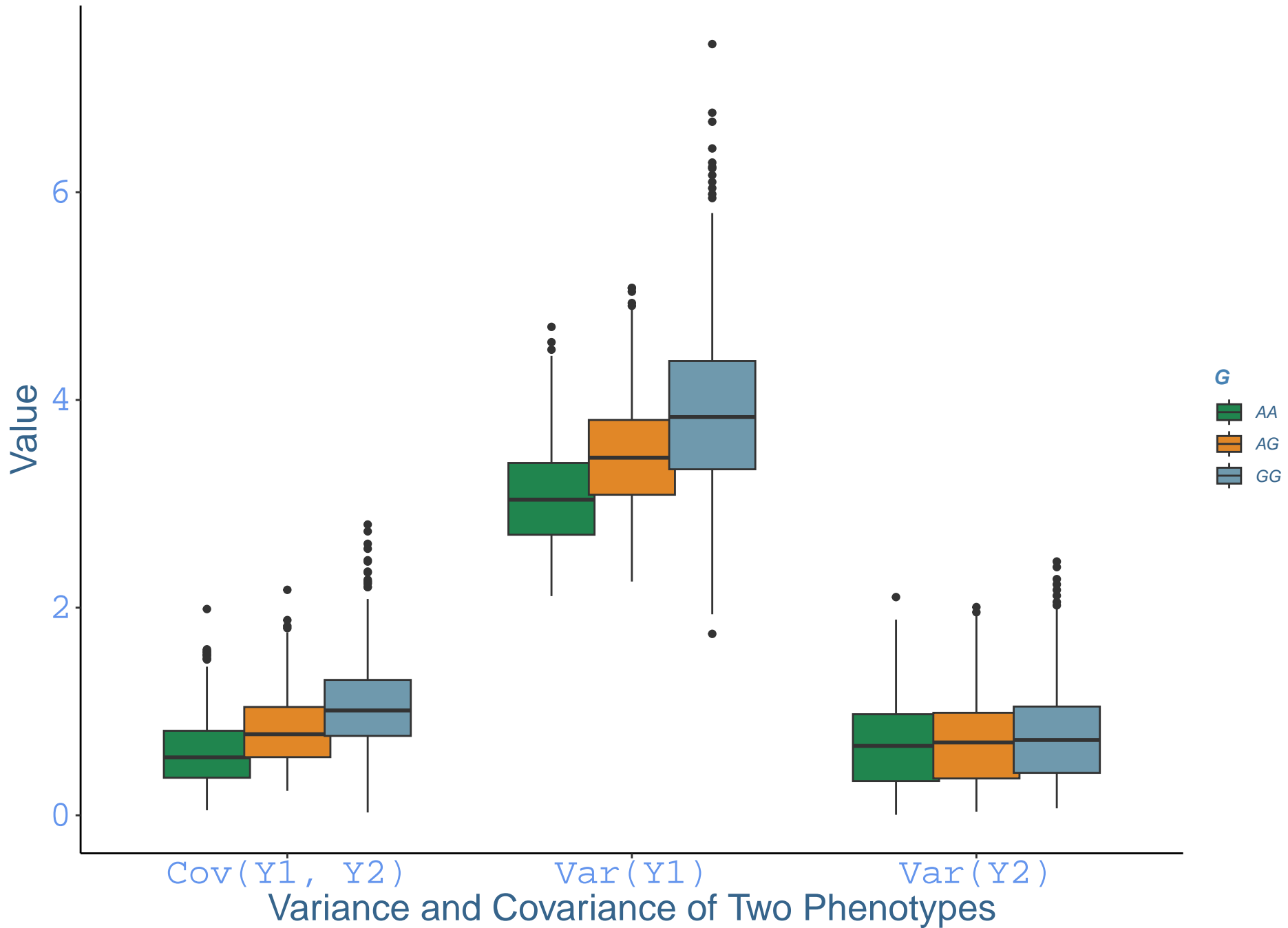
Variance and Covariance Distribution of Two Phenotypes Across Genotype Categories in 1,000 Simulated Datasets Without Interaction Effects

The differential covariance of the two phenotypes adds power while studying the interaction effects

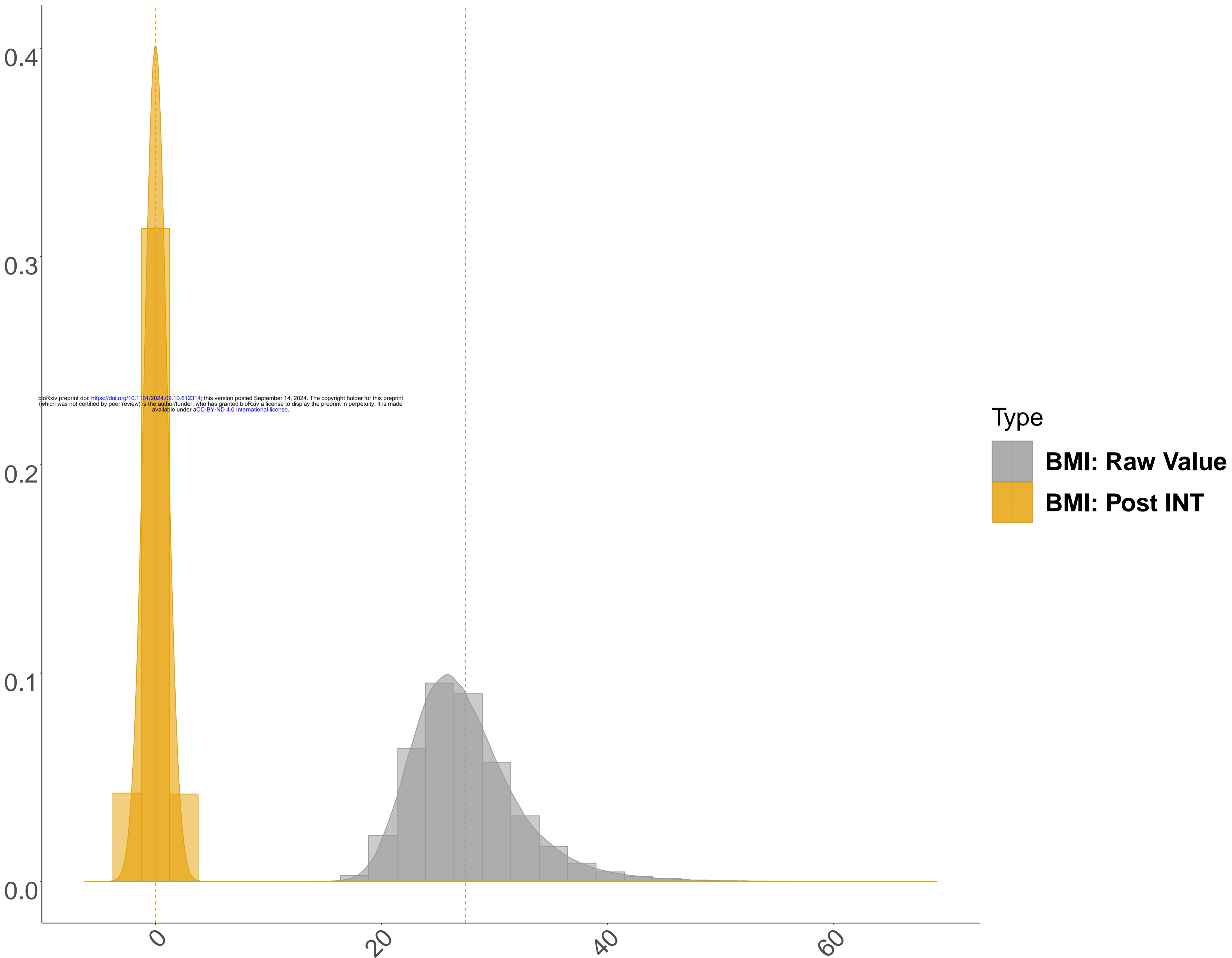


Variance and Covariance Distribution of Two Phenotypes Across Genotype Categories in 1,000 Simulated Datasets With Interaction Effects

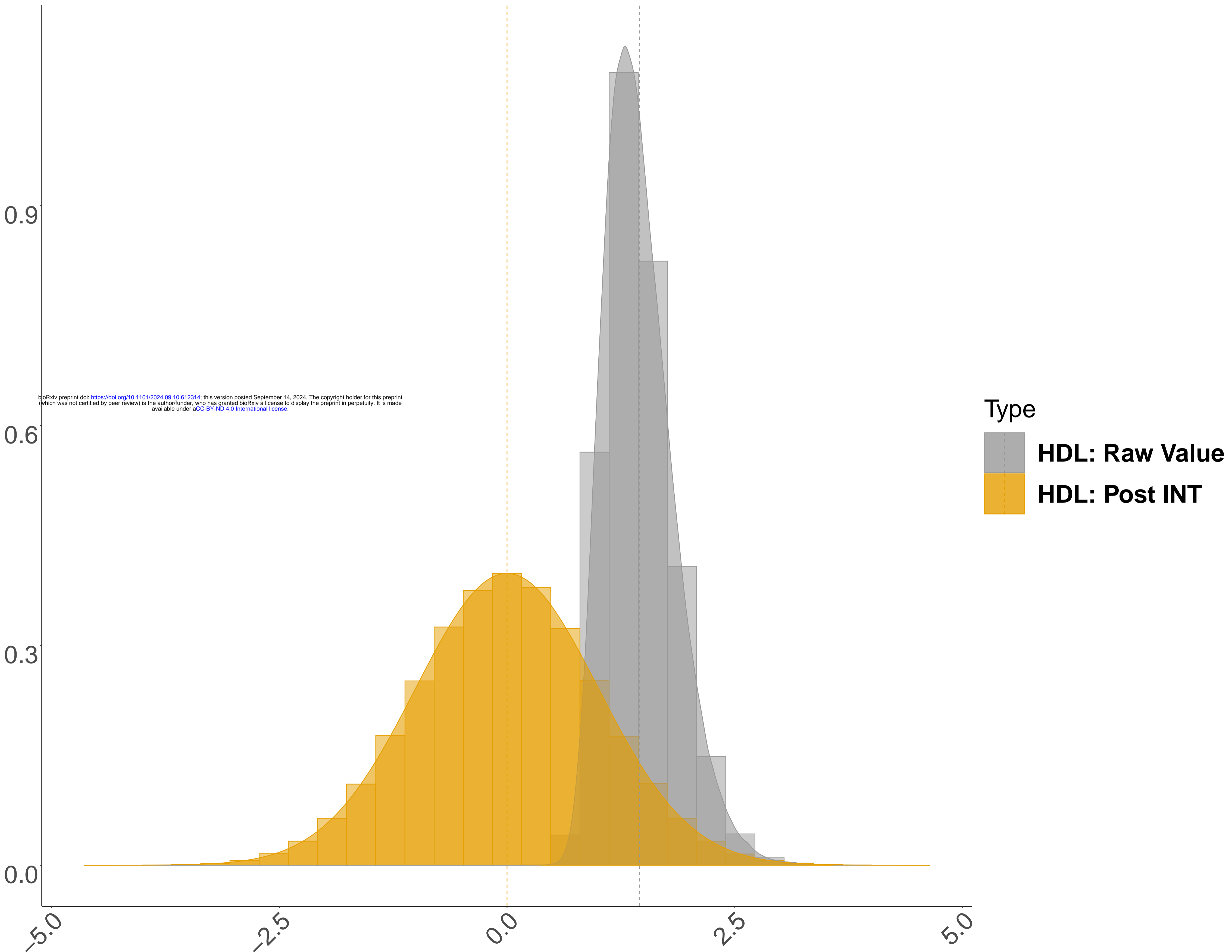
The differential covariance of the two phenotypes adds power while studying the interaction effects



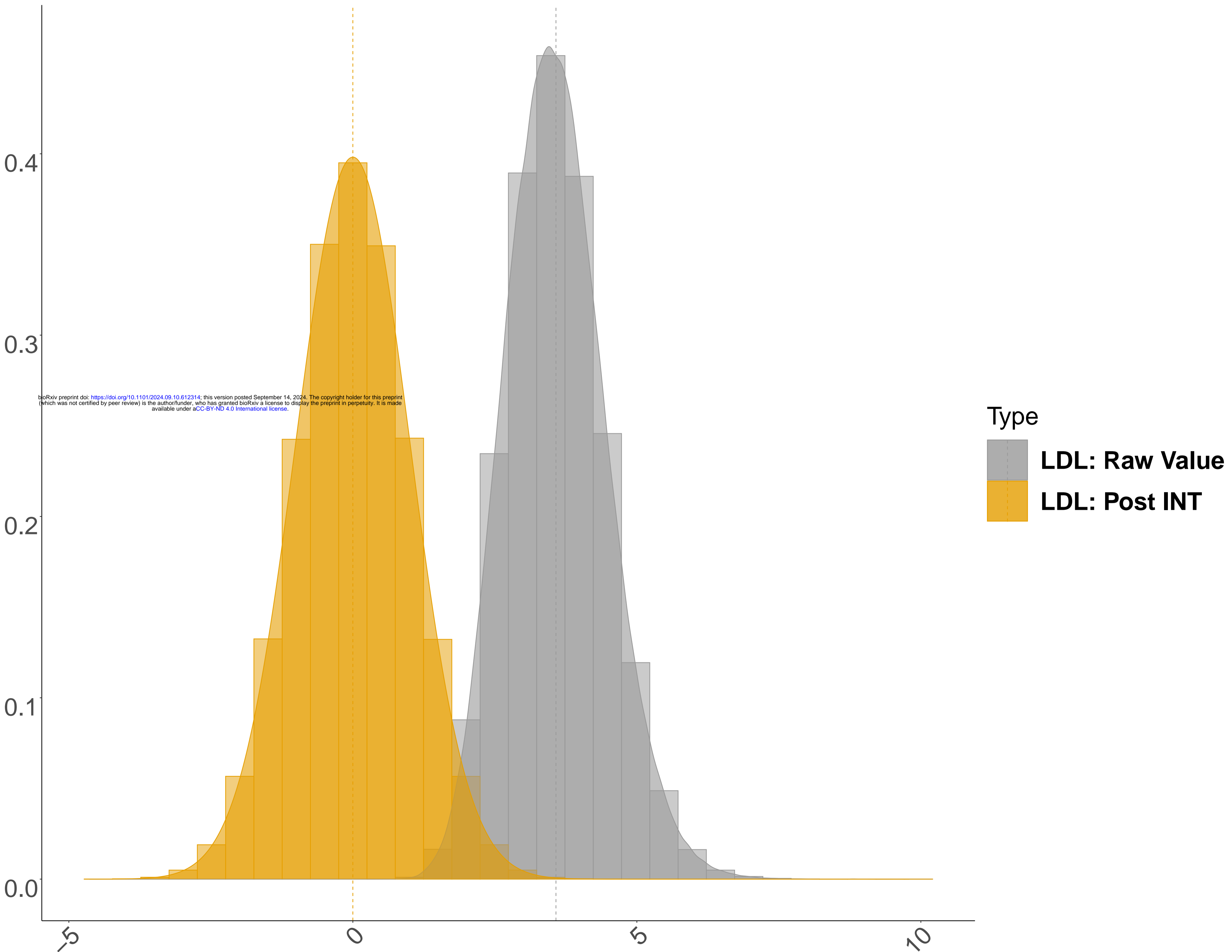
BMI: before and after inverse normal transformation (INT), N = 288,709



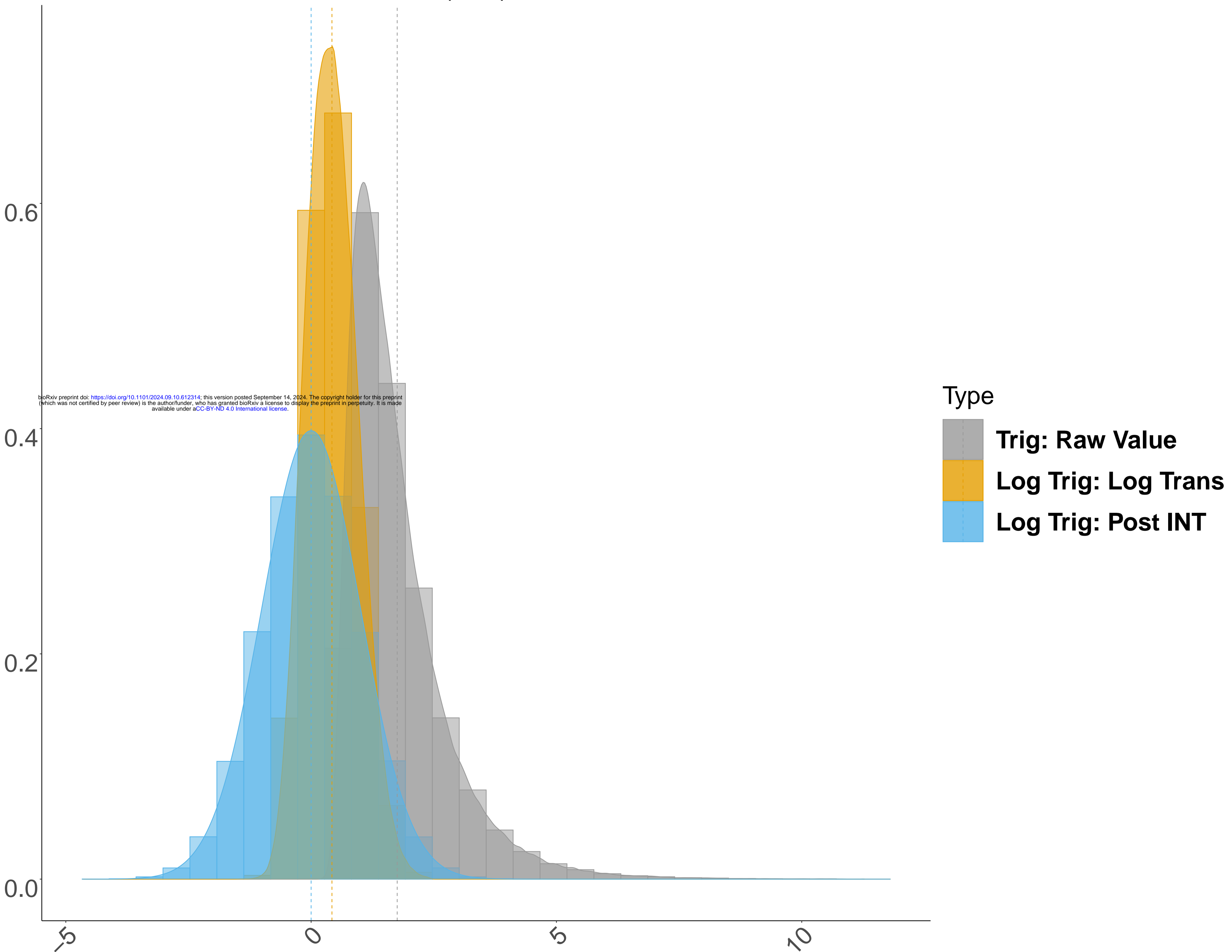
HDL: before and after inverse normal transformation (INT), N = 288,709



LDL: before and after inverse normal transformation (INT), N = 288,709



Natural Log Transformed Triglycerides: before and after inverse normal transformation (INT), N = 288,709



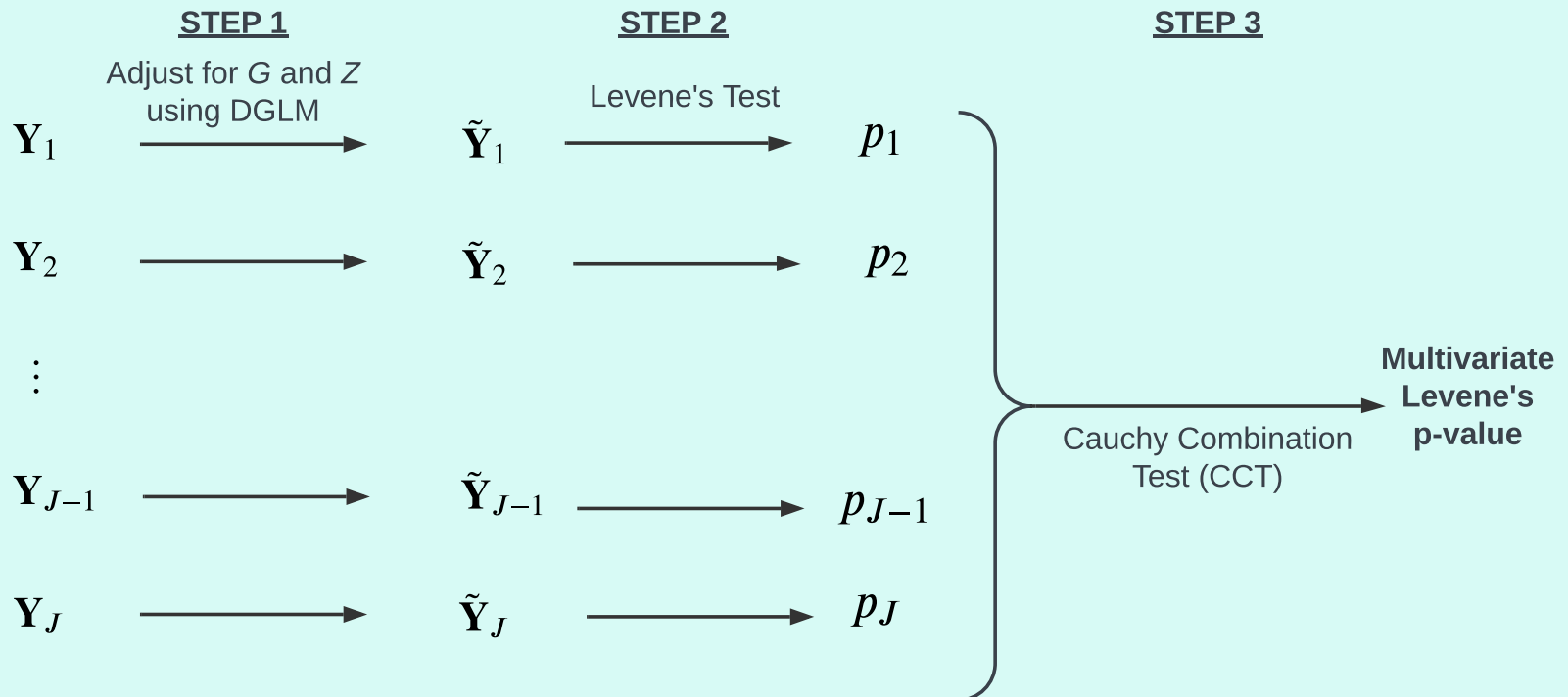
Input

P Phenotypes: $\mathbf{Y} = [\mathbf{Y}_1 \quad \mathbf{Y}_2 \quad \dots \quad \mathbf{Y}_J]$

K Confounders: $\mathbf{Z} = [\mathbf{Z}_1 \quad \mathbf{Z}_2 \quad \dots \quad \mathbf{Z}_K]$

One Test Genotype: G

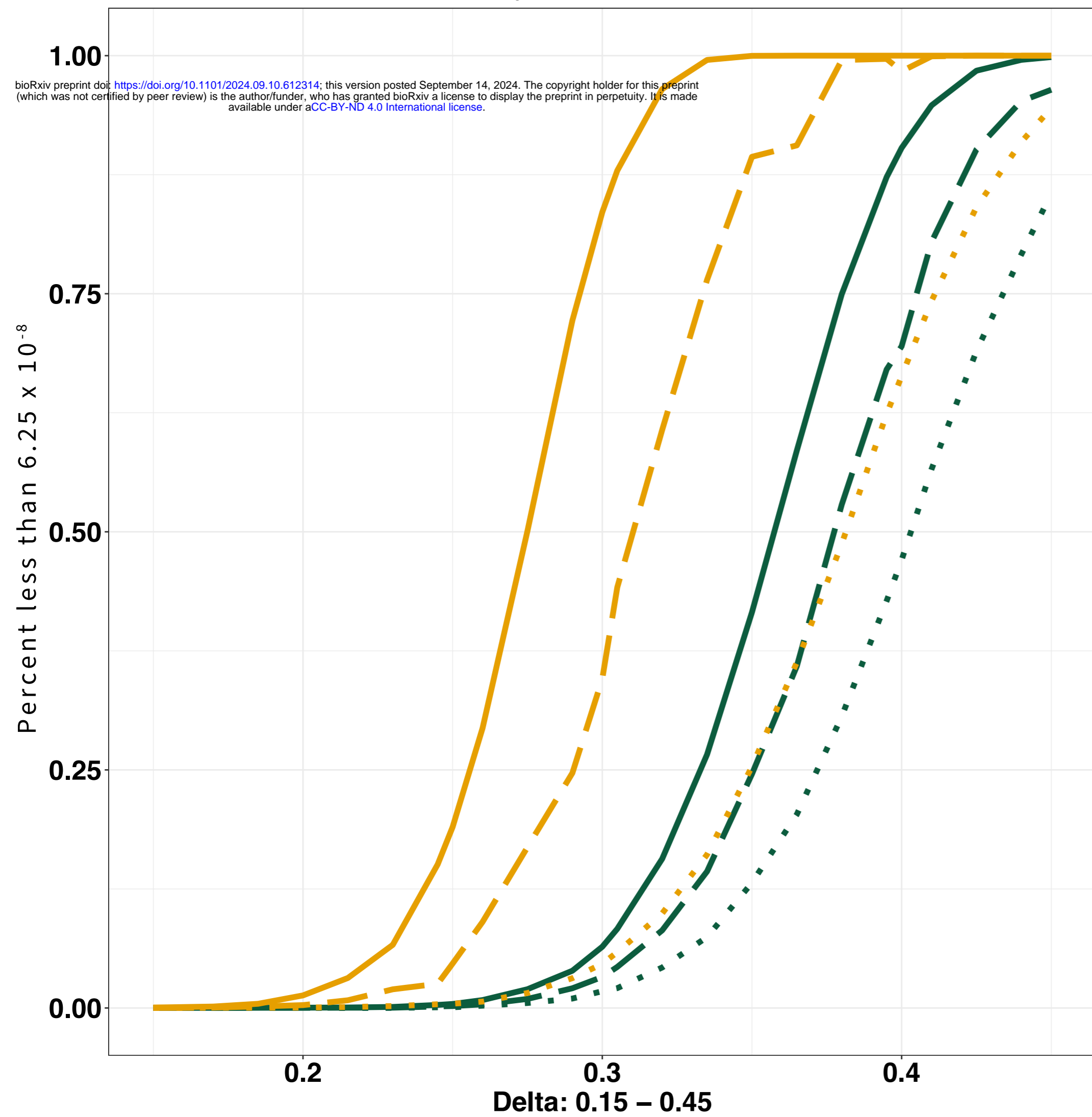
Multivariate Levene's Test Framework



J = 4

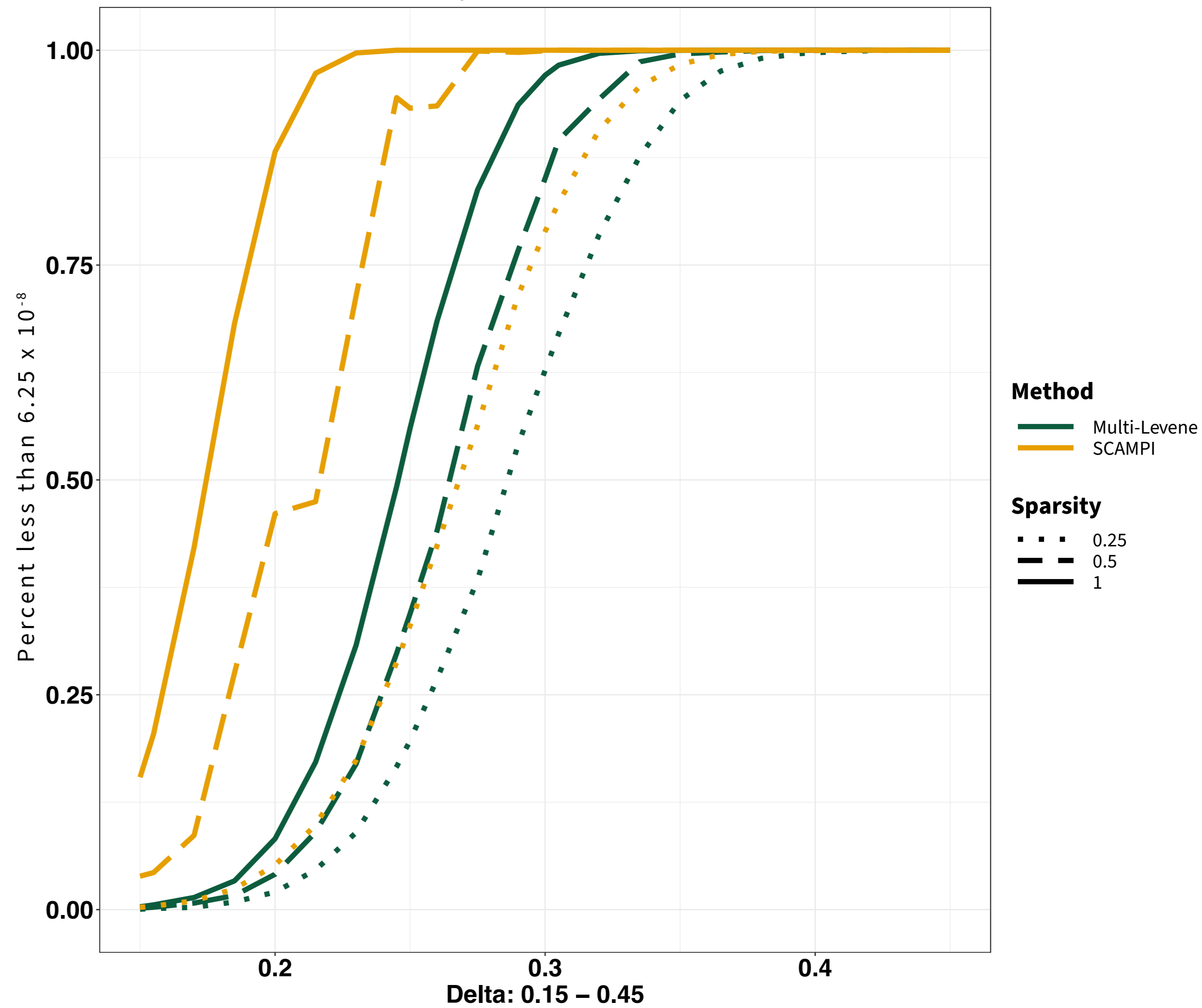
$\gamma = 0.05$

Sample size is 20000, MAF = 0.05, $\gamma = 0.05$



$\gamma = 0.25$

Sample size is 20000, MAF = 0.05, $\gamma = 0.25$

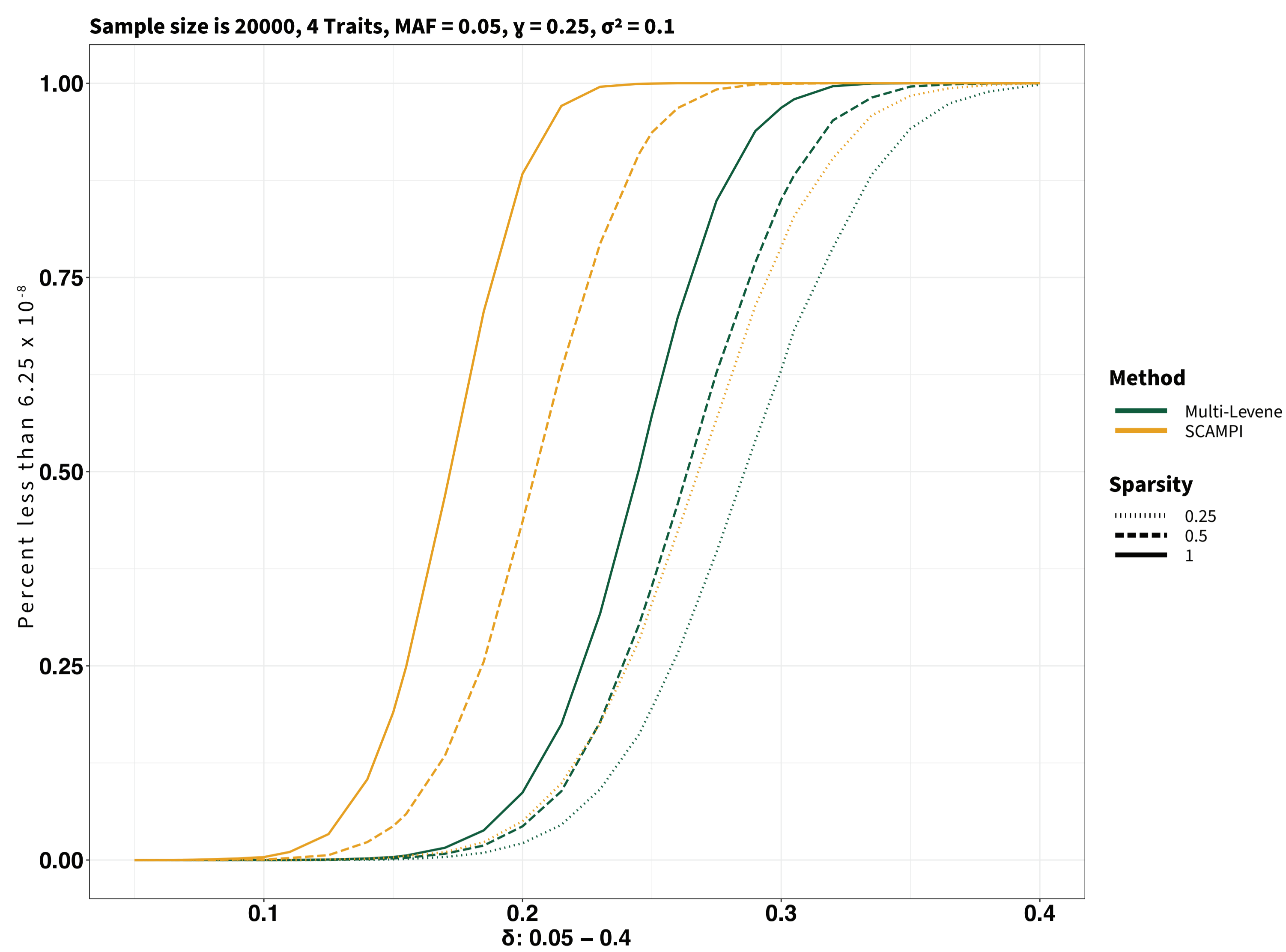
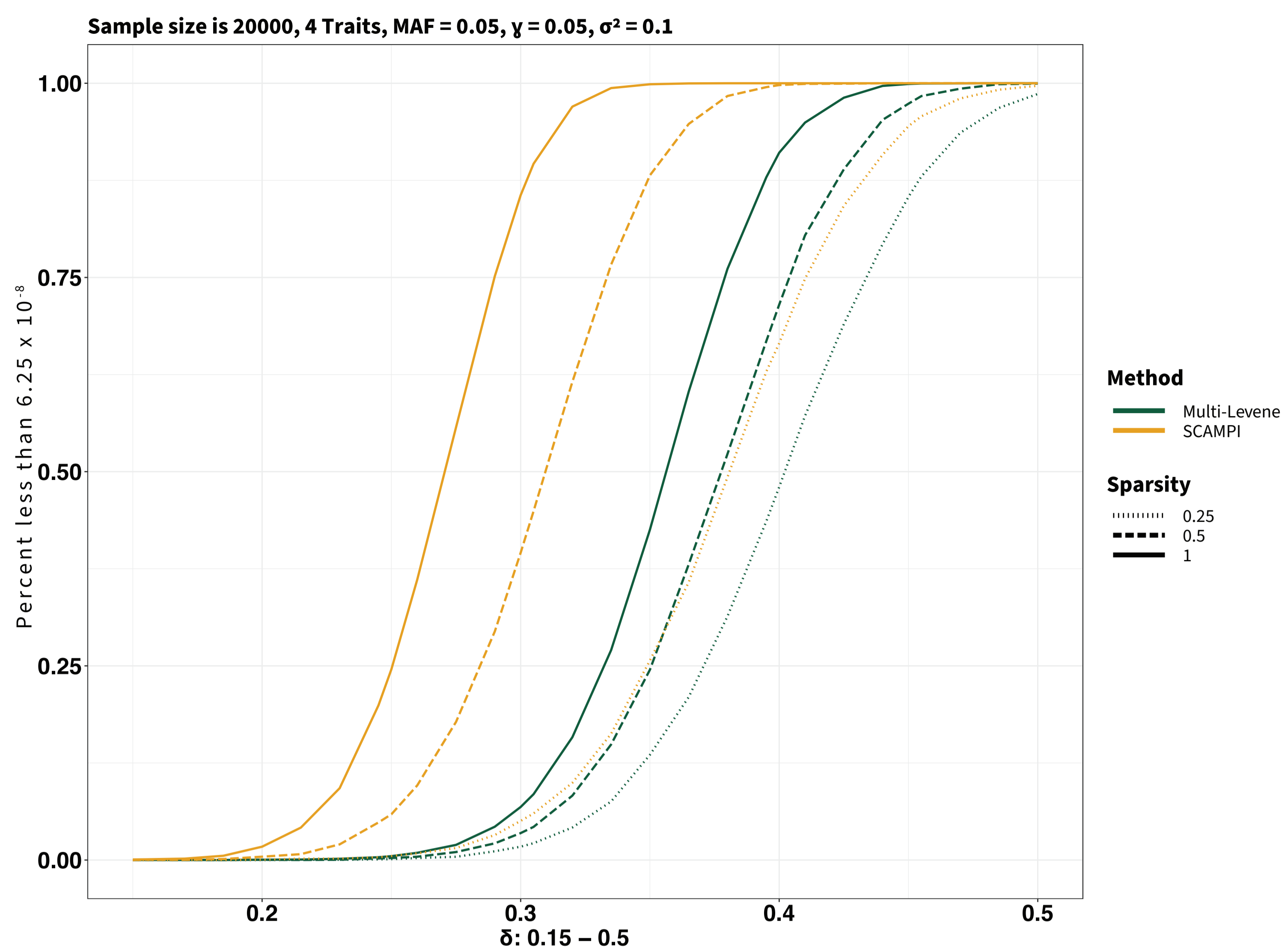


J = 4

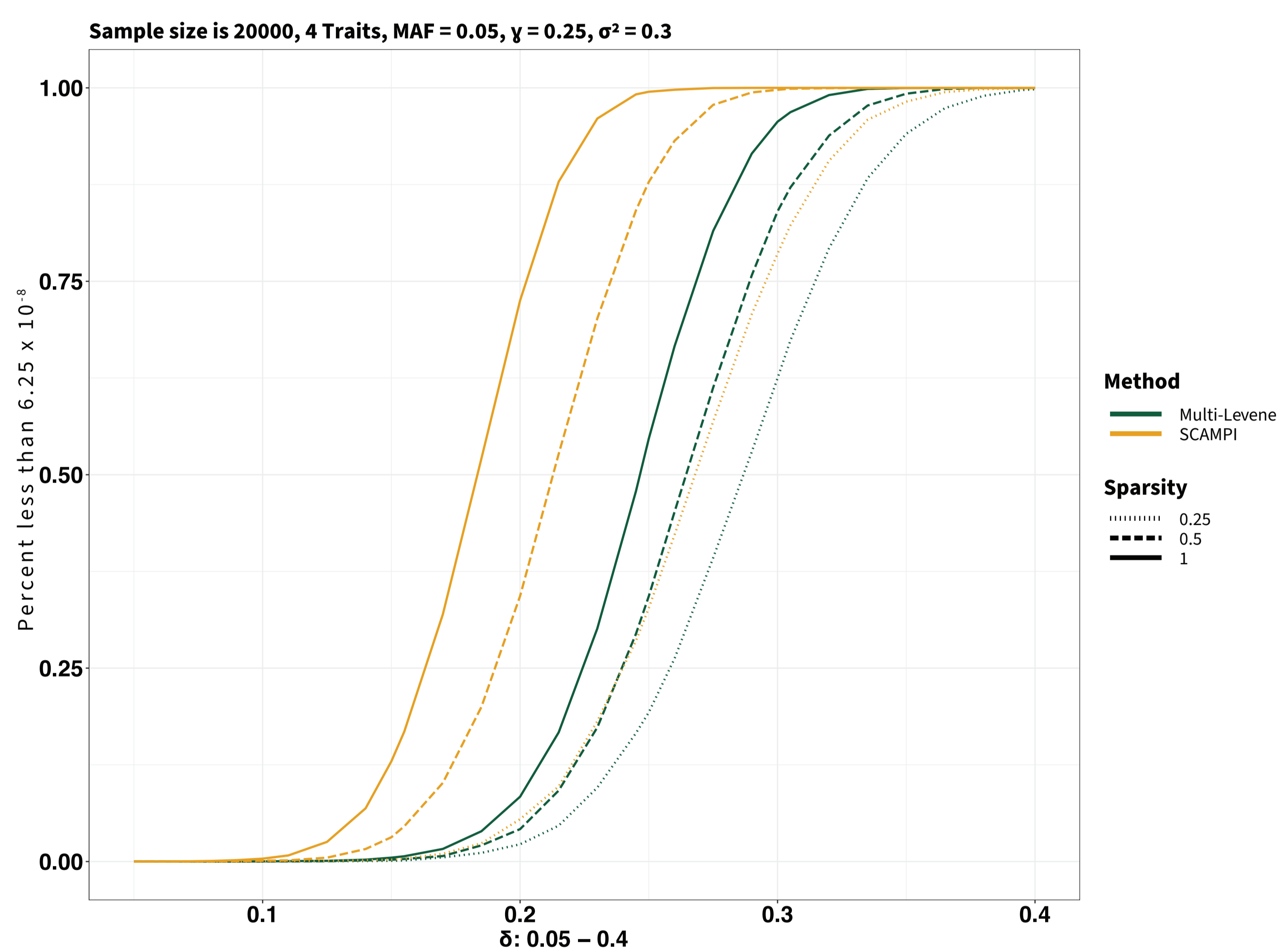
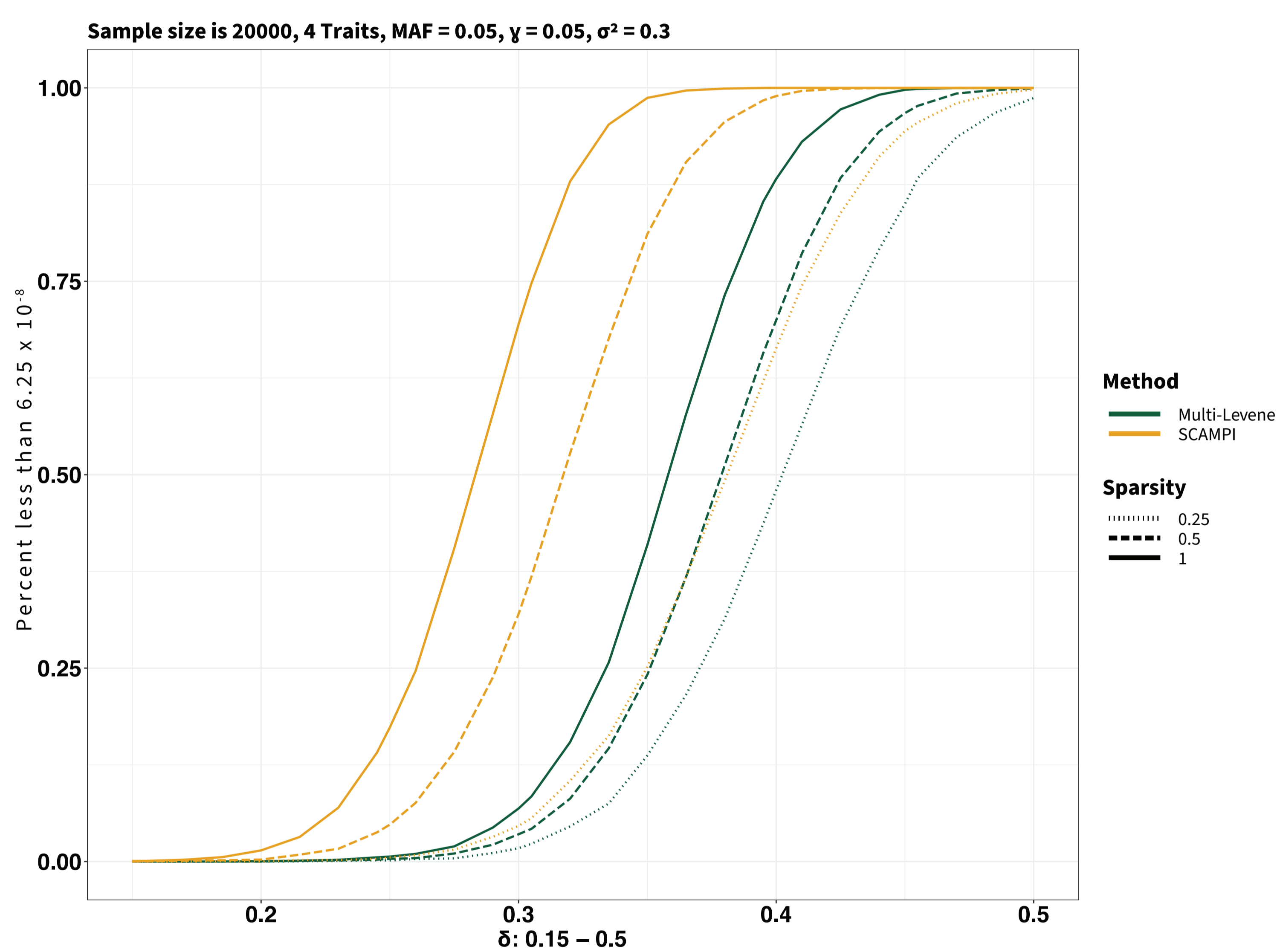
$\gamma = 0.05$

$\gamma = 0.25$

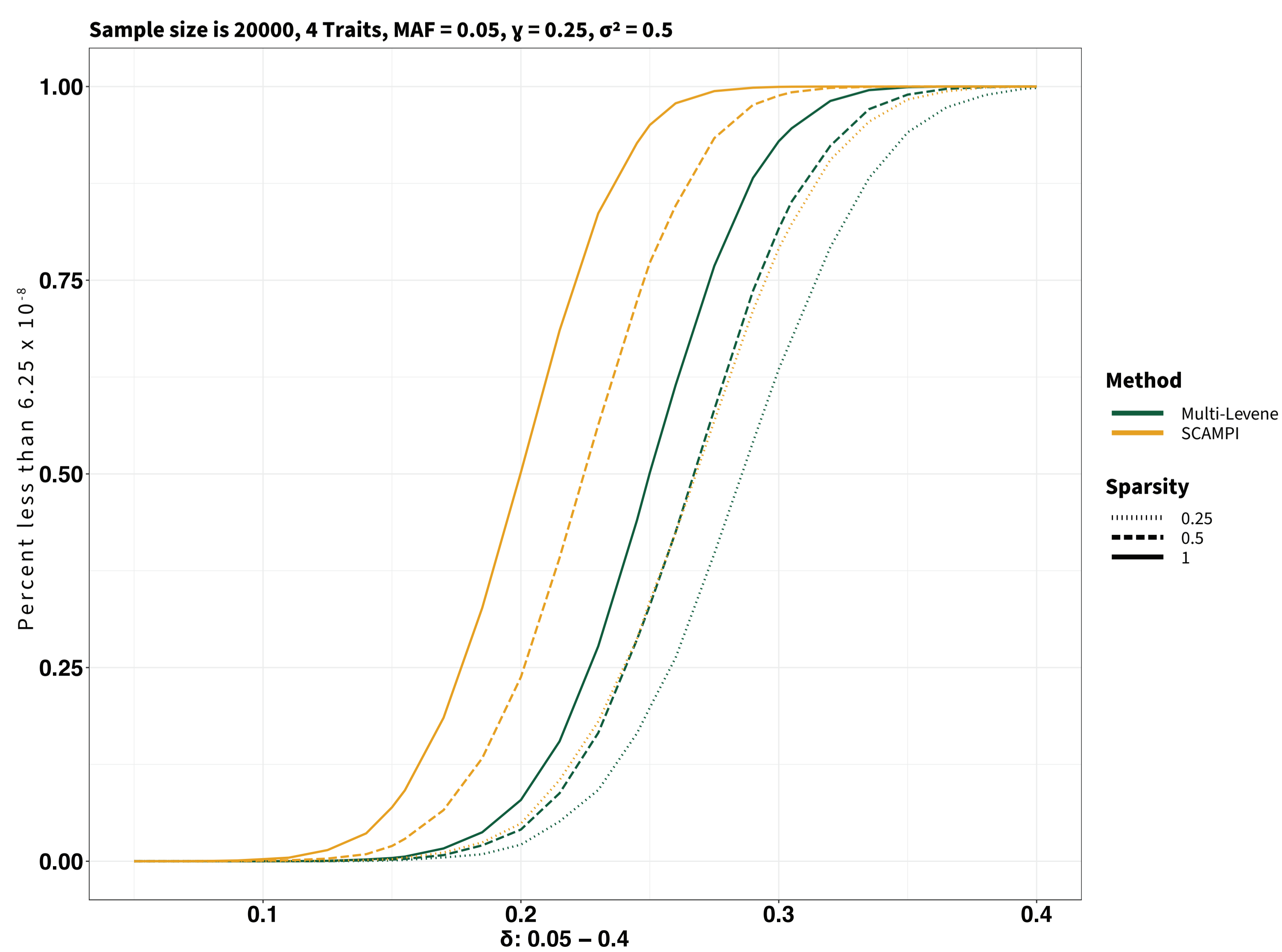
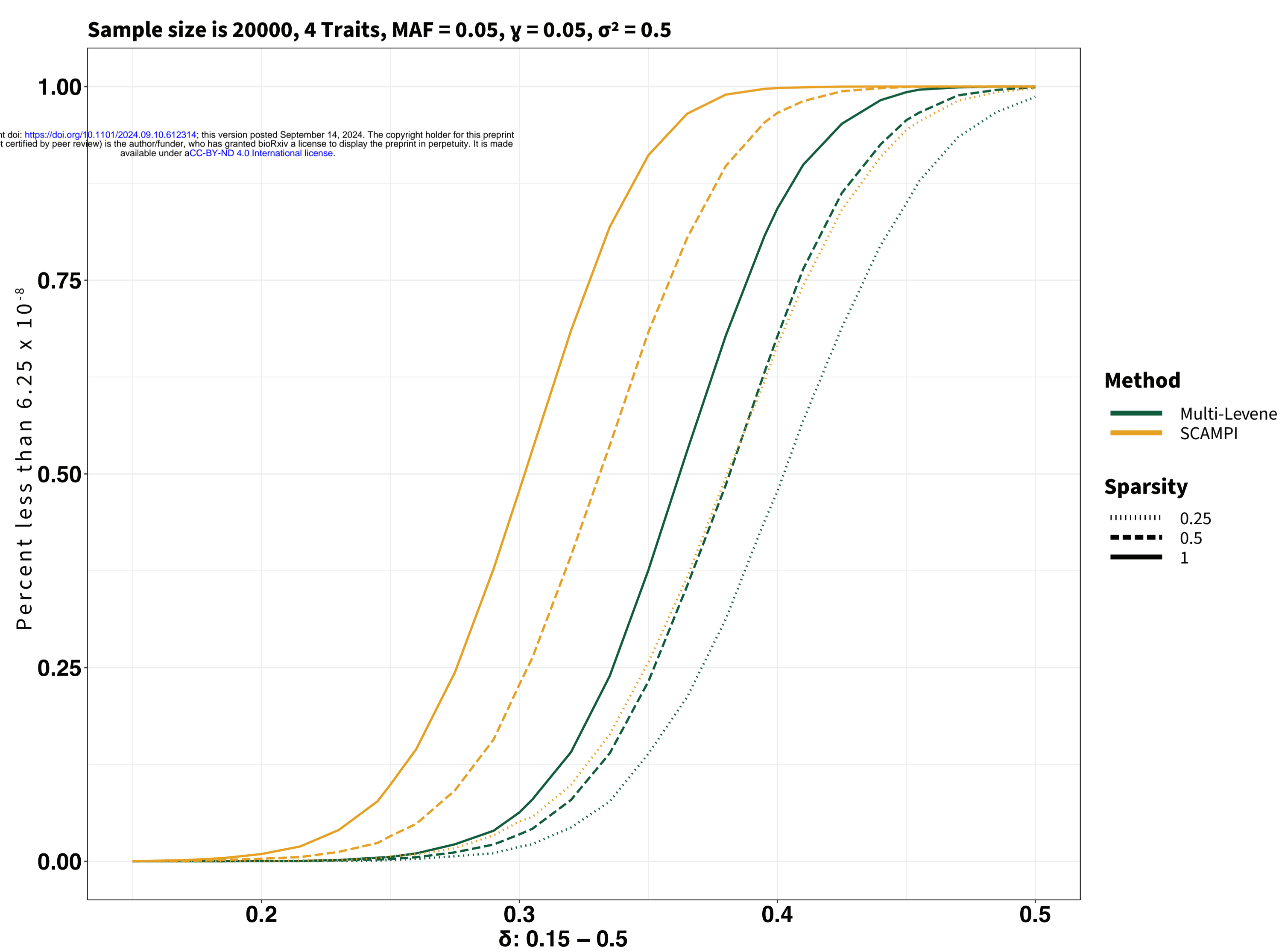
$\sigma^2 = 0.1$



$\sigma^2 = 0.3$



$\sigma^2 = 0.5$



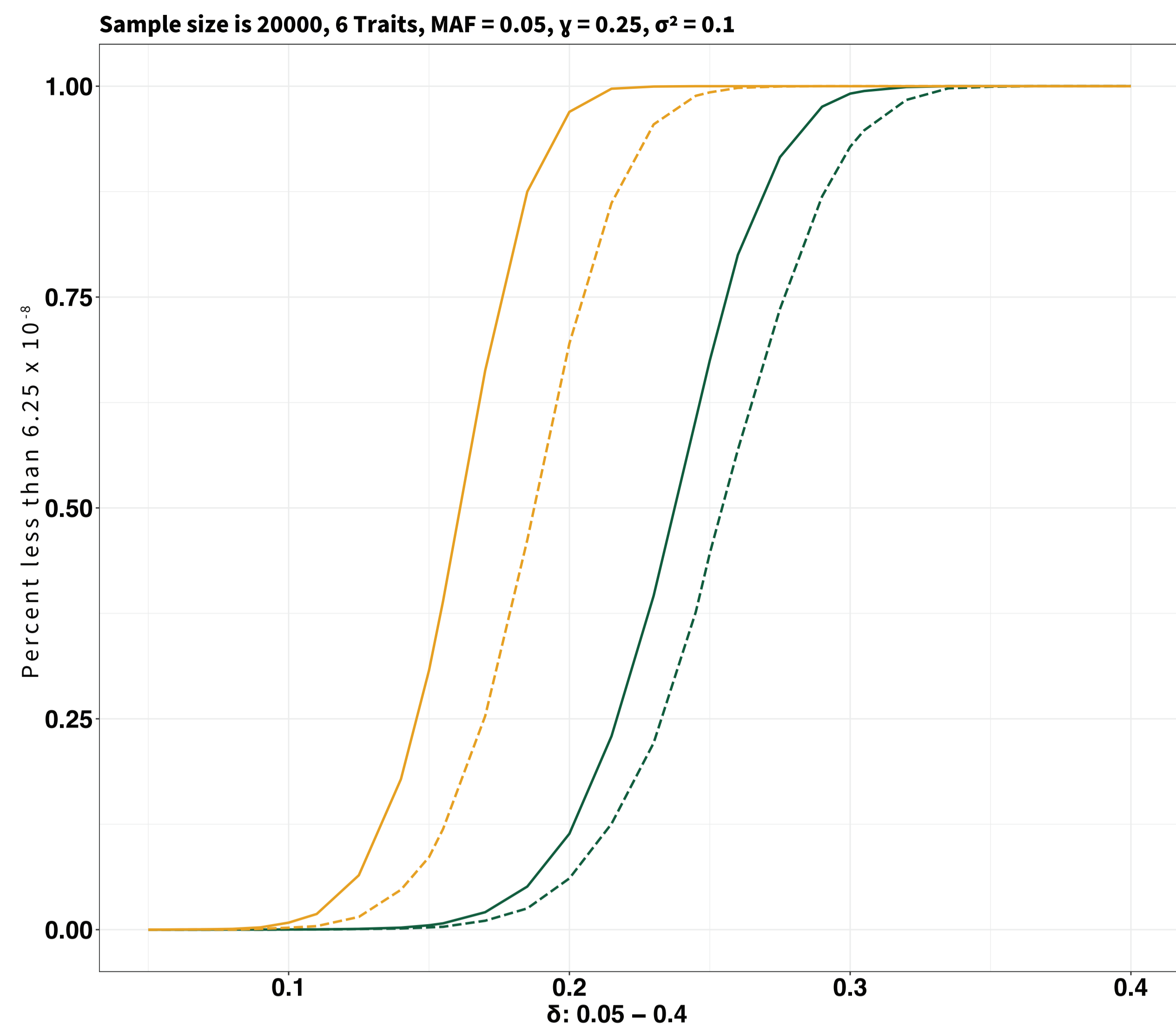
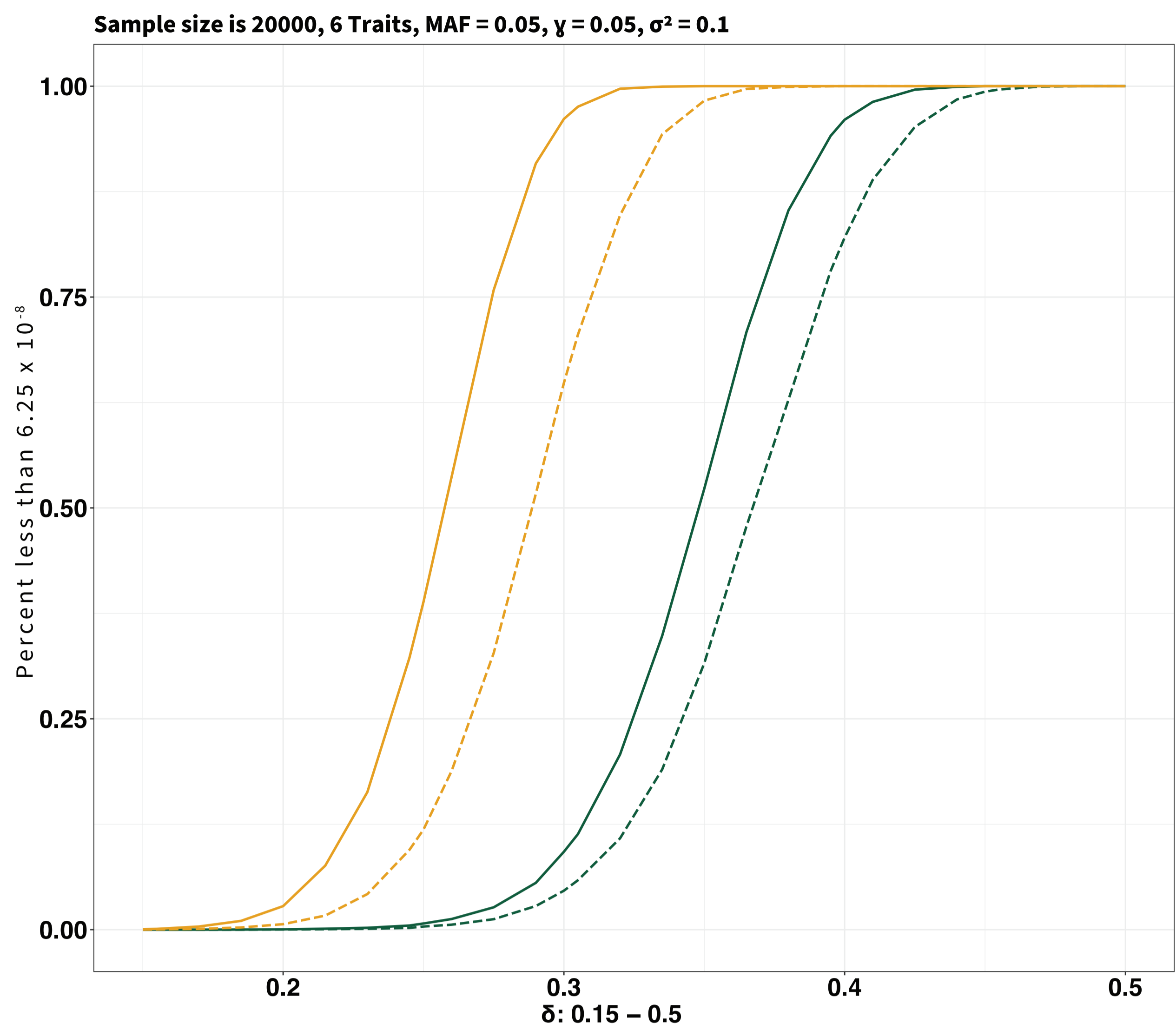
bioRxiv preprint doi: <https://doi.org/10.1101/2024.09.10.612314>; this version posted September 14, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

J = 6

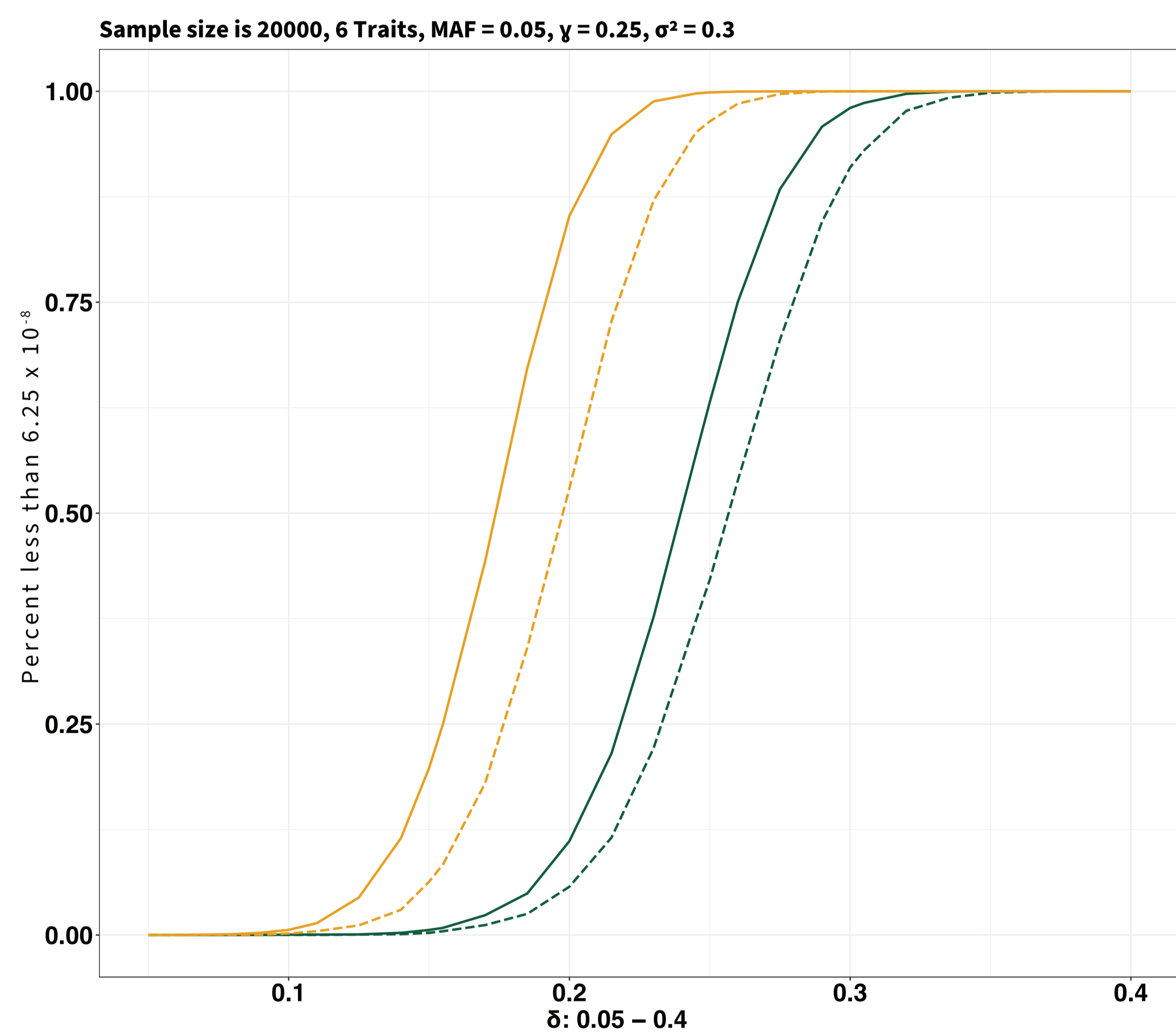
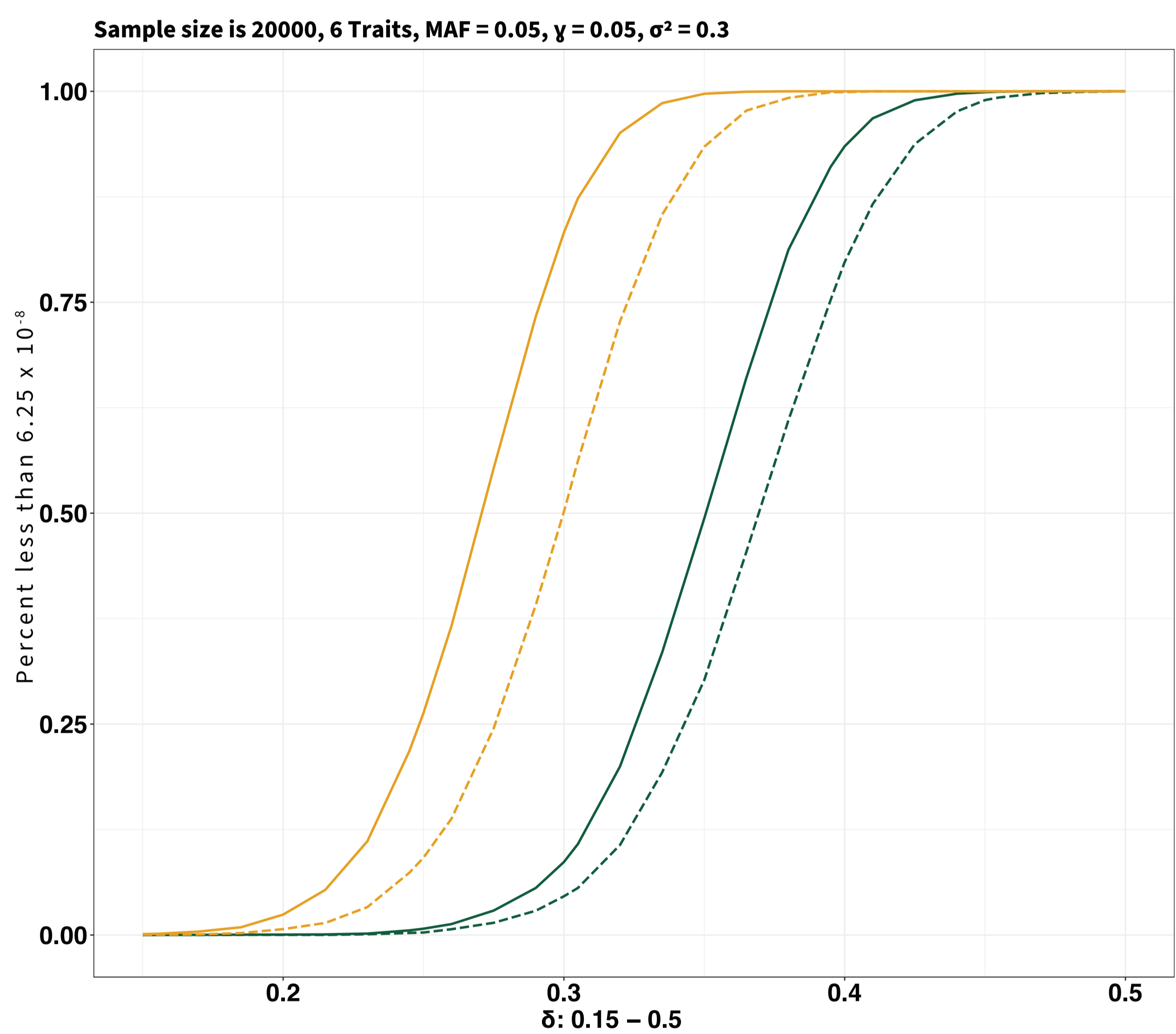
$\gamma = 0.05$

$\gamma = 0.25$

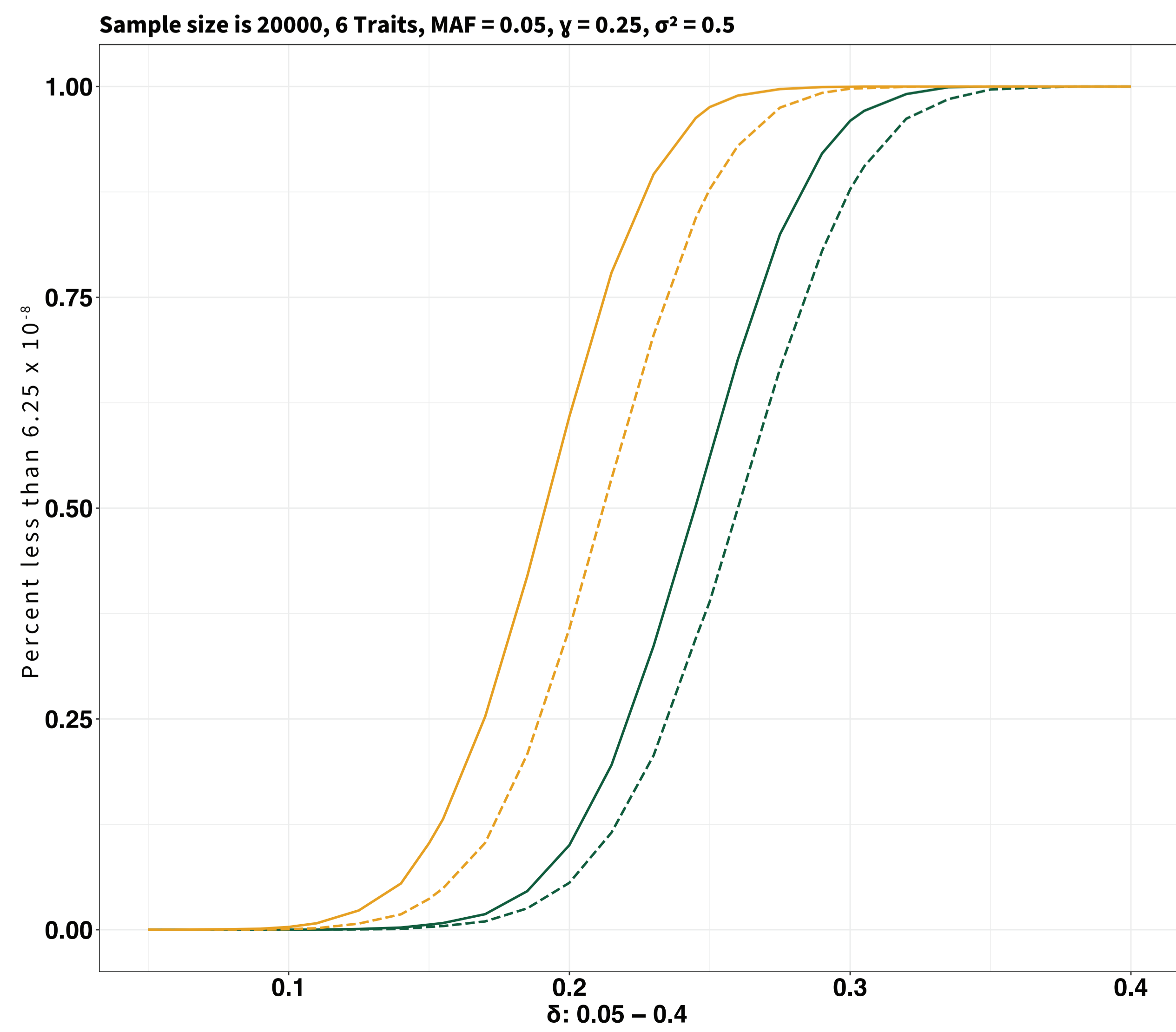
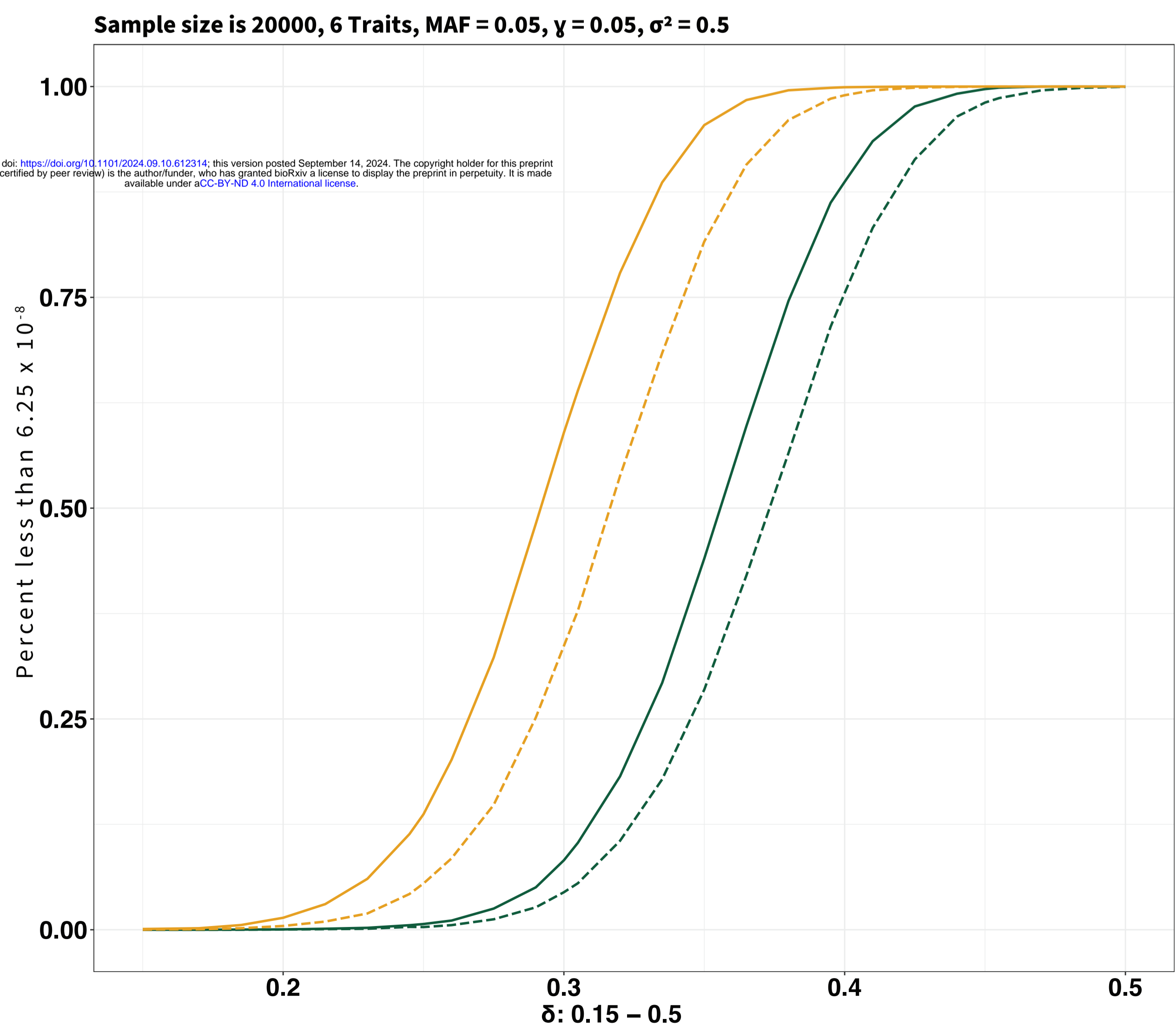
$\sigma^2 = 0.1$



$\sigma^2 = 0.3$



$\sigma^2 = 0.5$



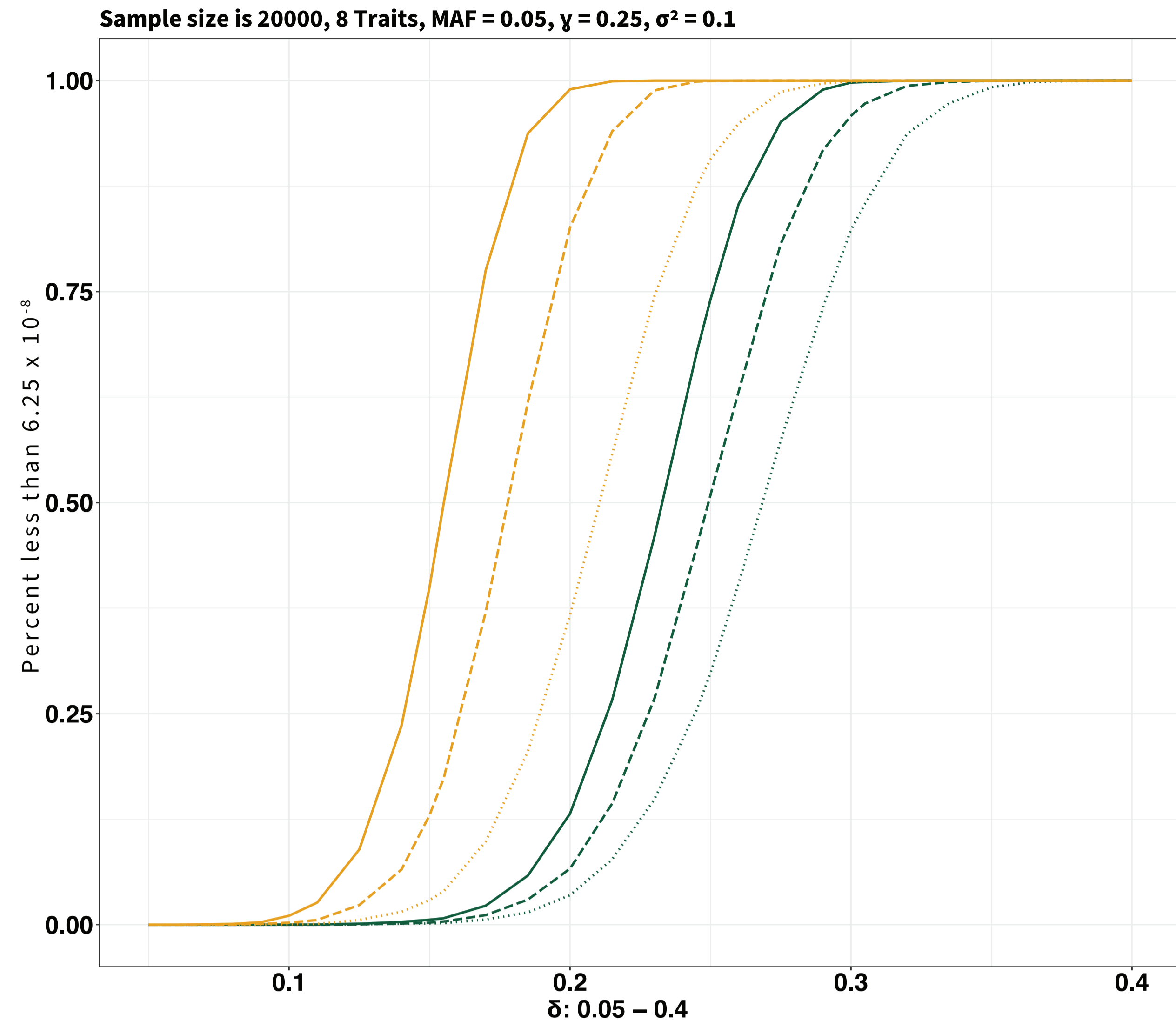
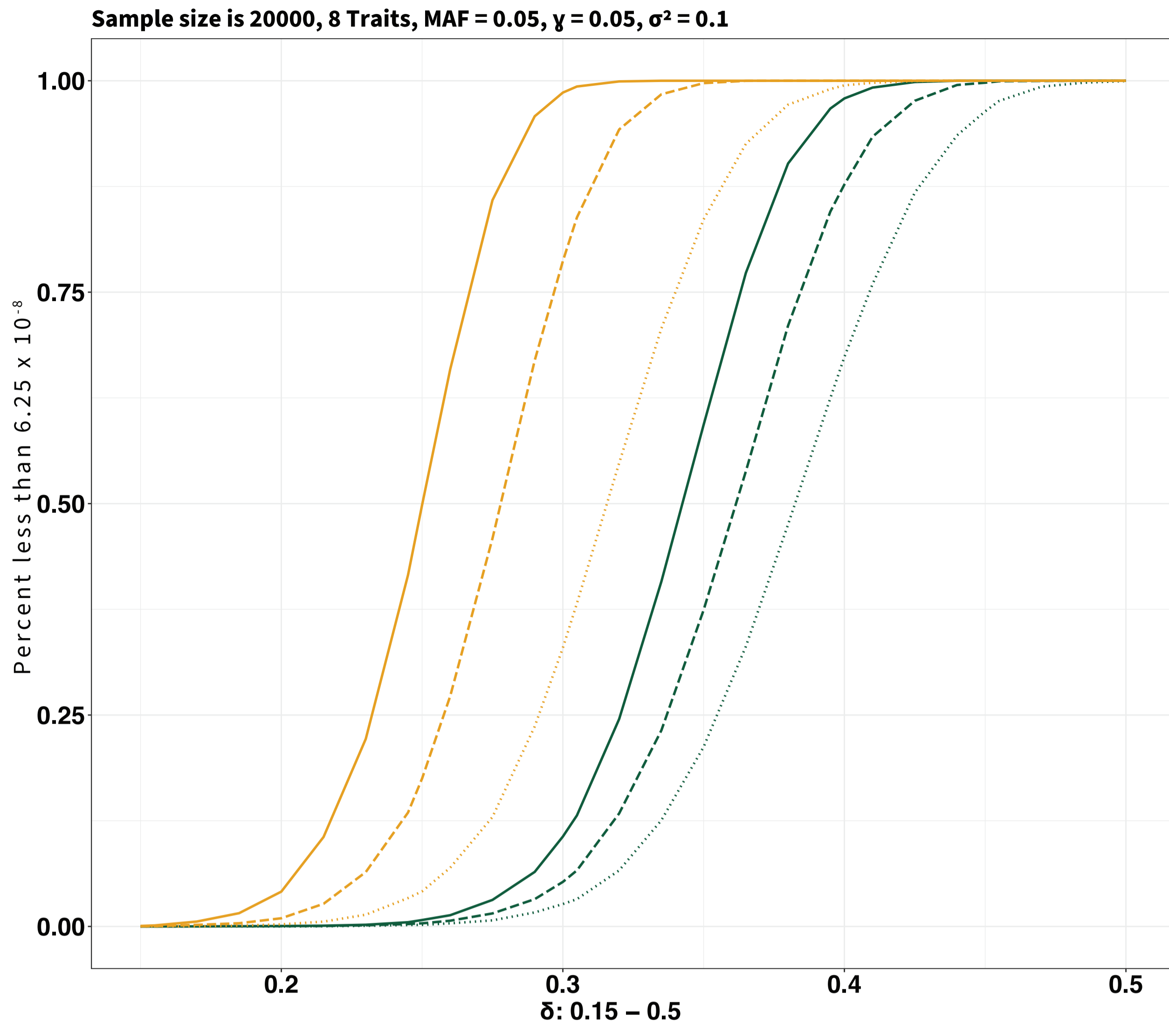
bioRxiv preprint doi: <https://doi.org/10.1101/2024.09.10.612314>; this version posted September 14, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

J = 8

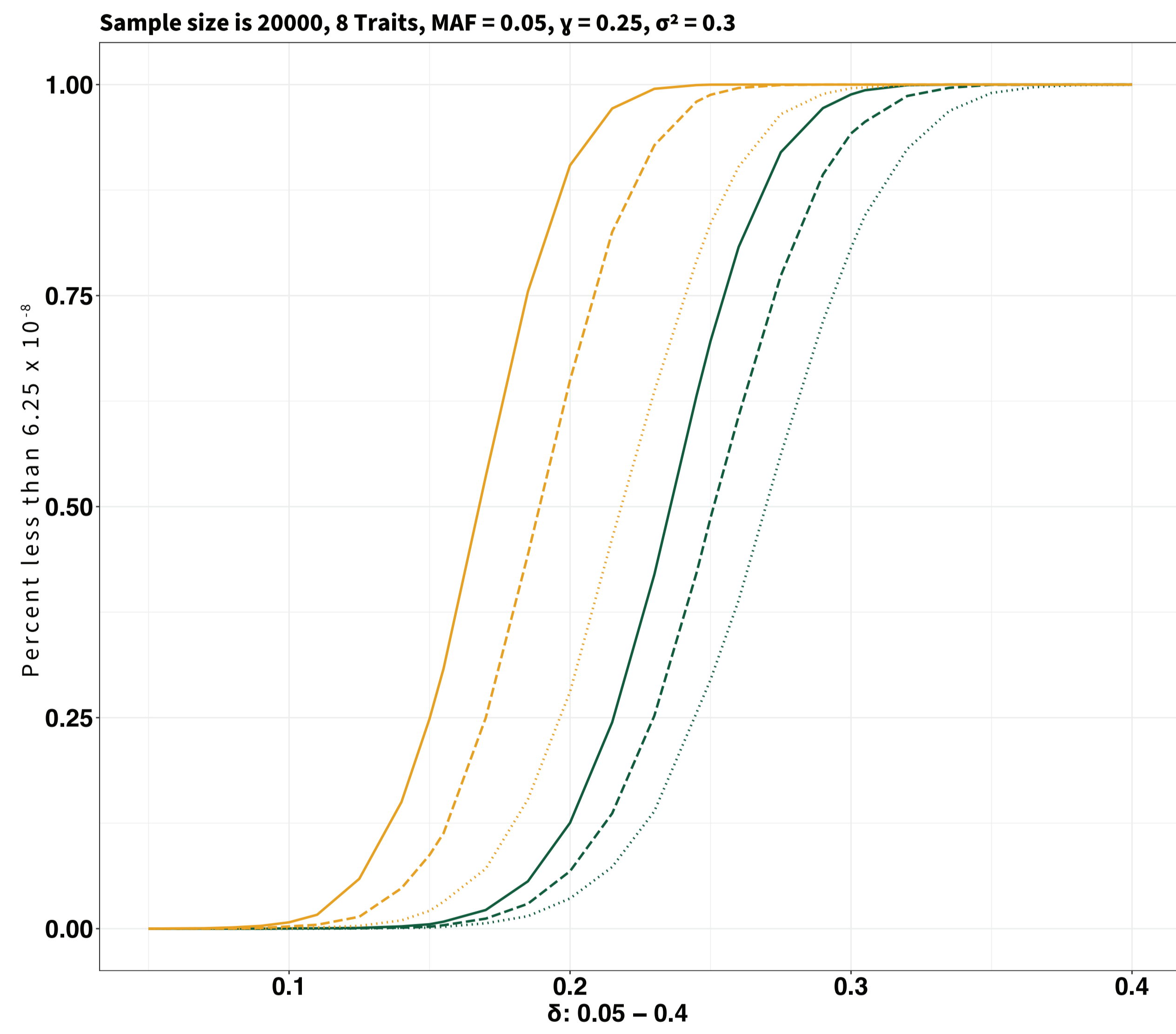
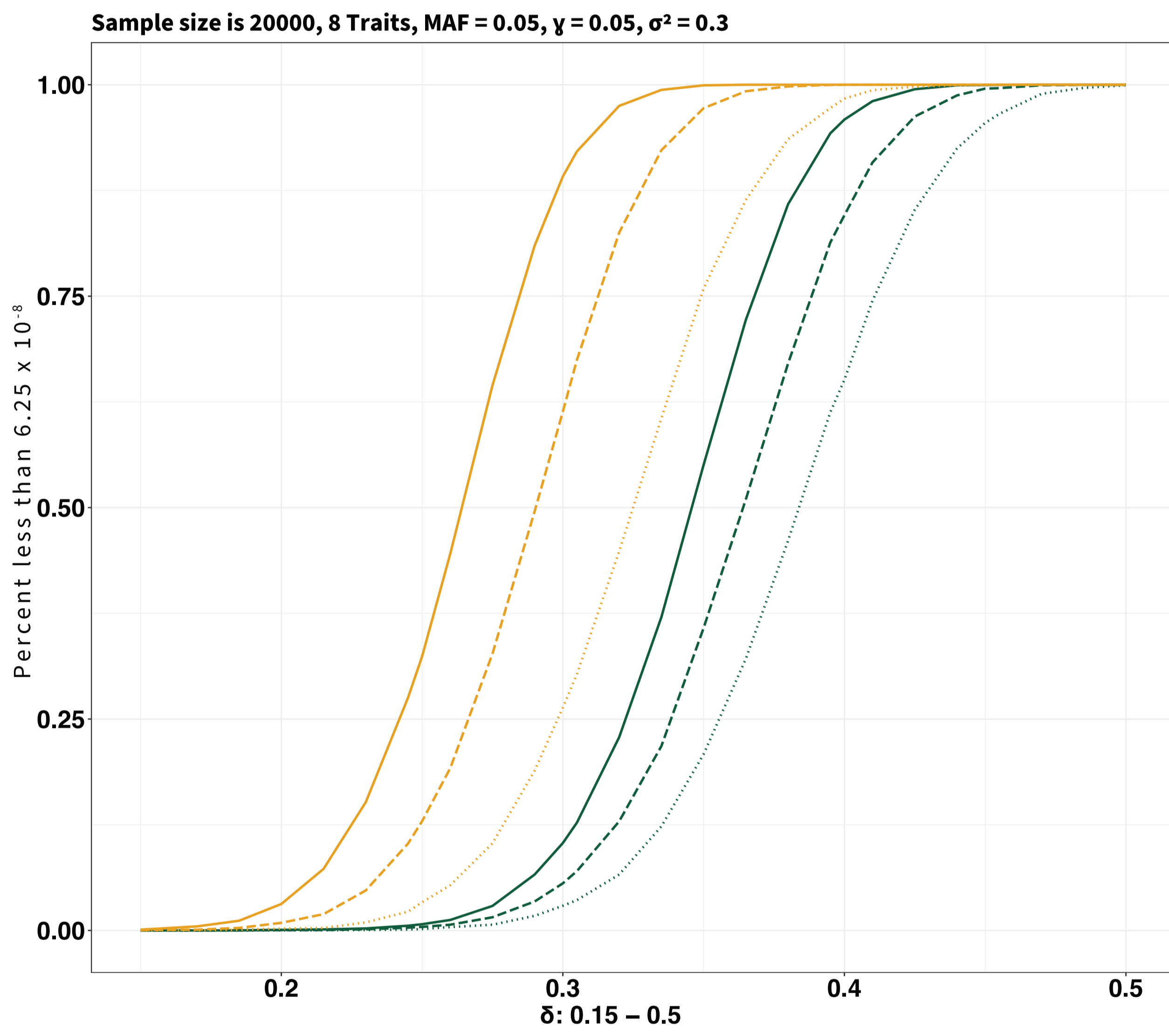
$\gamma = 0.05$

$\gamma = 0.25$

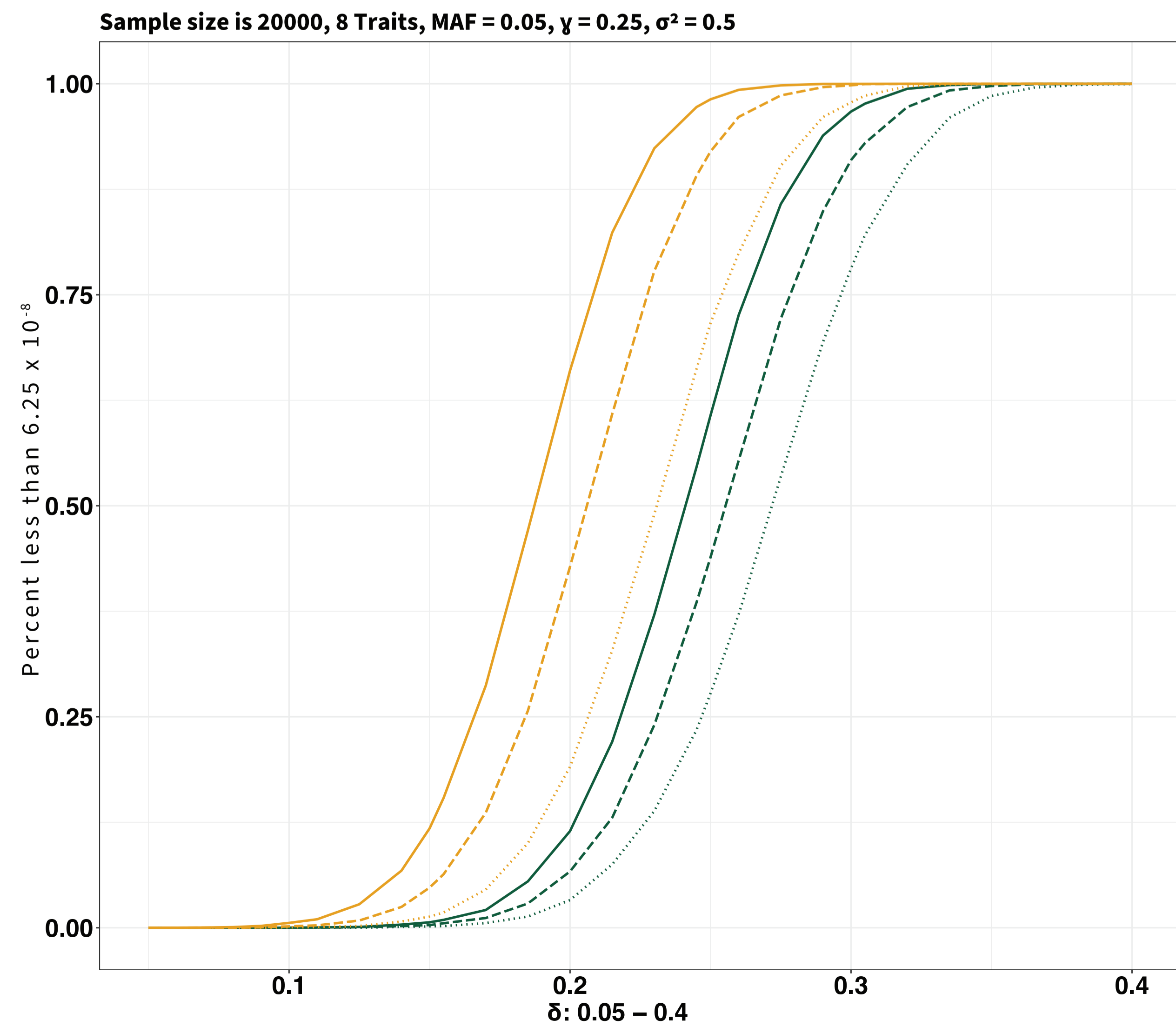
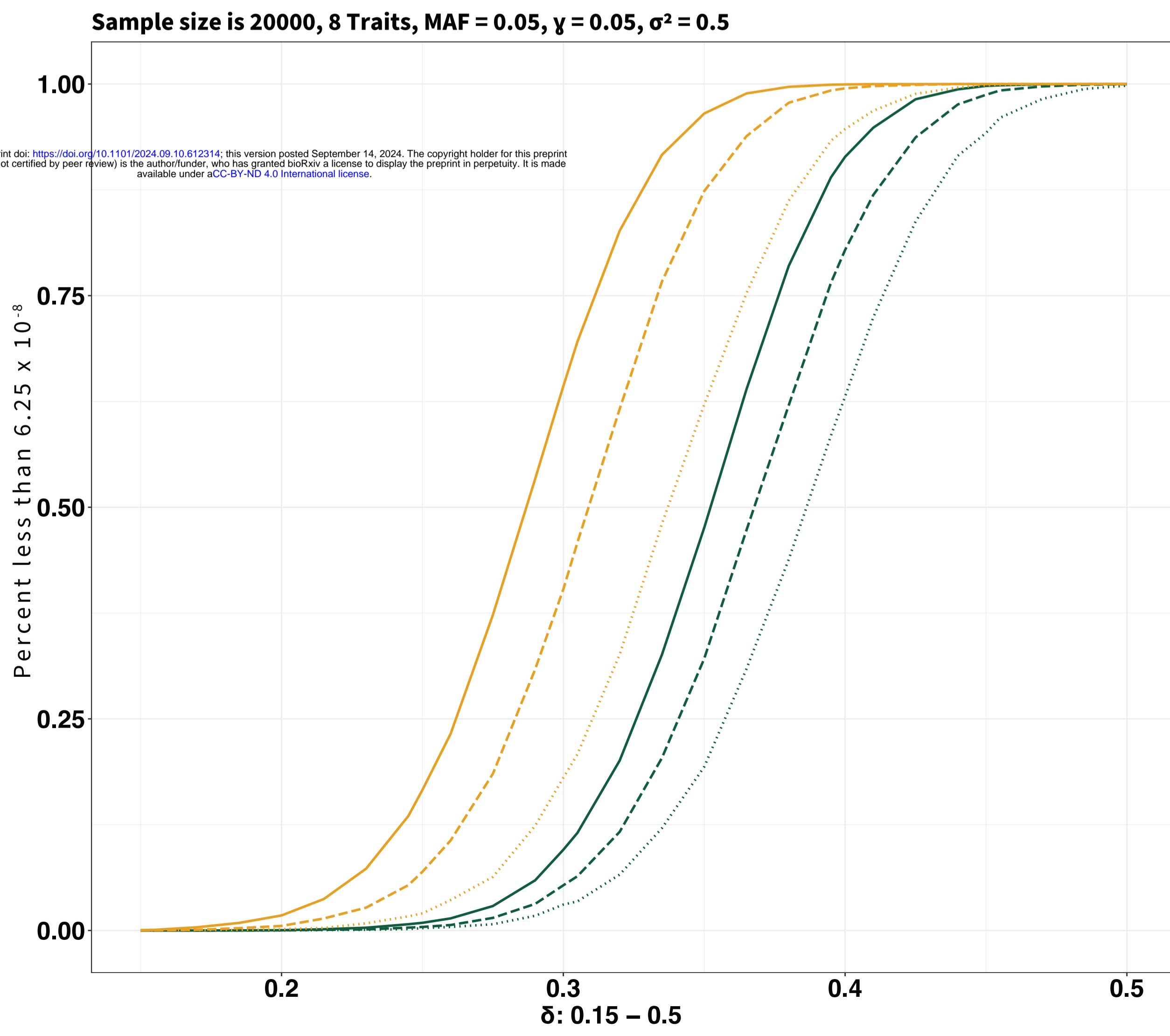
$\sigma^2 = 0.1$



$\sigma^2 = 0.3$

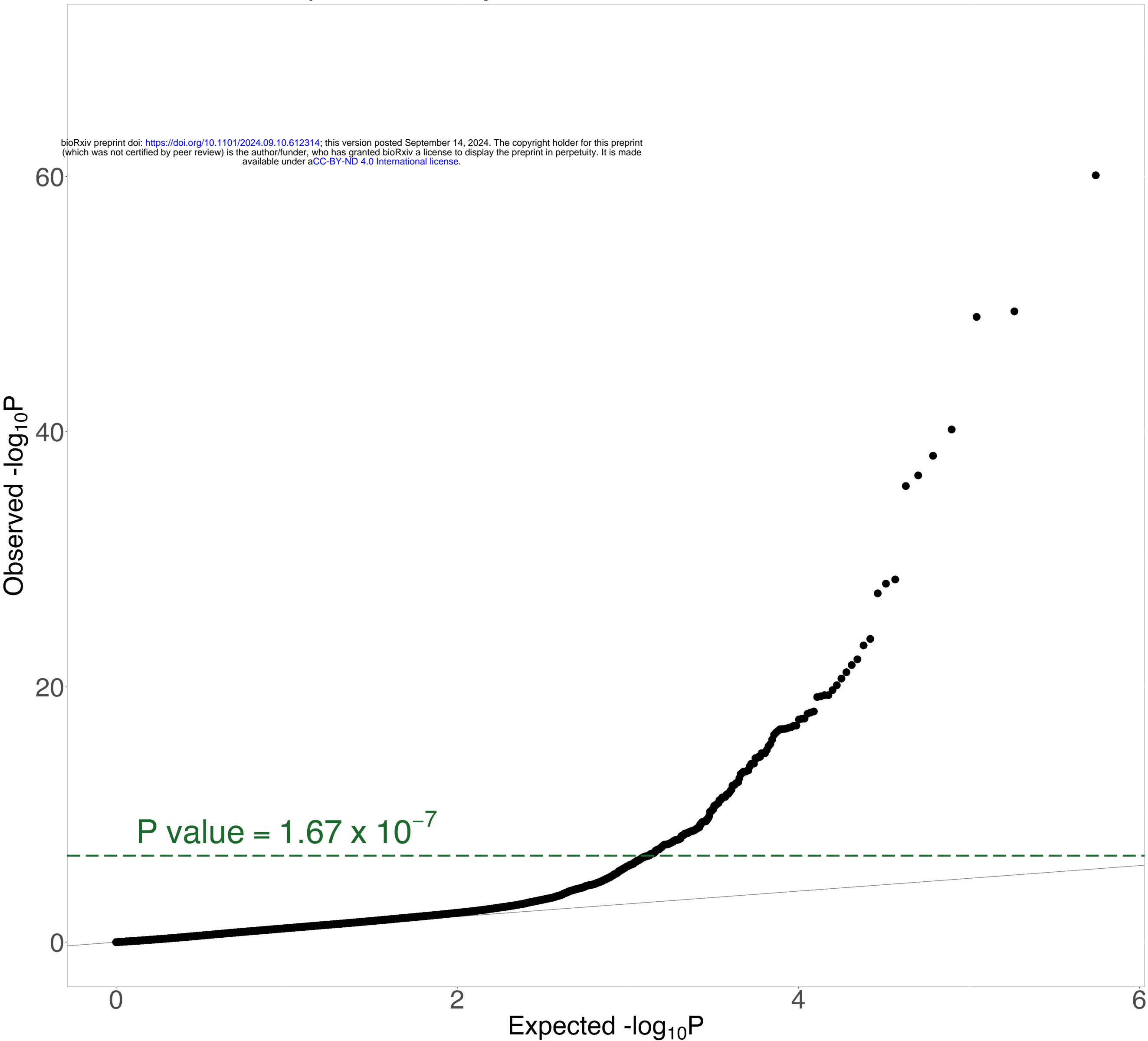


$\sigma^2 = 0.5$

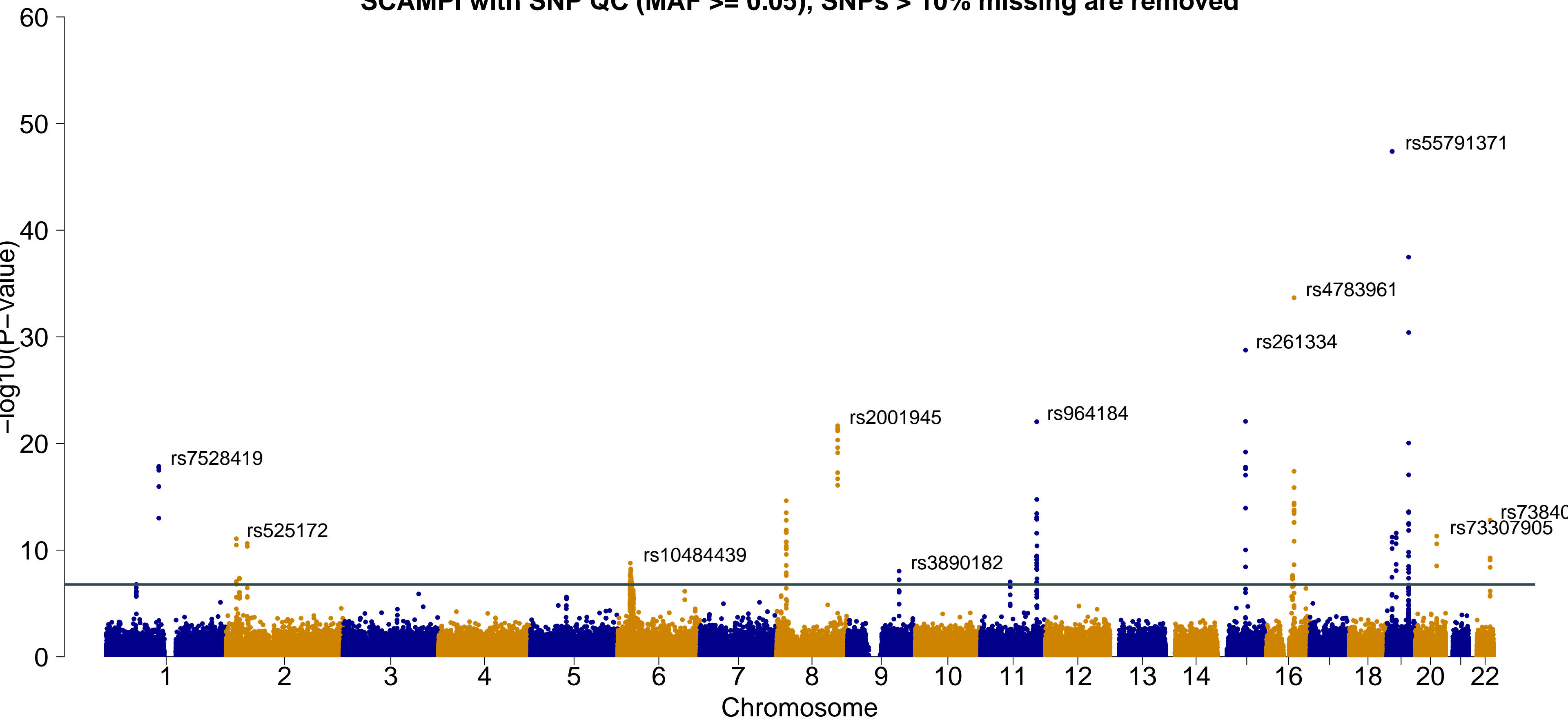


SCAMPI with SNP QC (MAF ≥ 0.05),
SNPs $> 10\%$ missing are removed
N = 288,709 Independent Subjects in UKBB

bioRxiv preprint doi: <https://doi.org/10.1101/2024.09.10.612314>; this version posted September 14, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.



N = 288,709 Independent Subjects in UKBB
SCAMPI with SNP QC (MAF \geq 0.05), SNPs $>$ 10% missing are removed



SCAMPI with SNP QC (MAF ≥ 0.05),
SNPs $> 10\%$ missing are removed
N = 241,167 Independent Subjects in UKBB

bioRxiv preprint doi: <https://doi.org/10.1101/2024.09.10.612314>; this version posted September 14, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

