1

# A study of analytical strategies to include X-chromosome in variance heterogeneity analysis: evidence for trait–specific polygenic variance structure

4

5

Wei Q. Deng[1], Anette Kalnapenkis[2,3], Shihong Mao[4], Tõnu Esko[2,5], Reedik Mägi[2], Guillaume Paré[4,6], Lei Sun[† 1,7]

8

[1]Department of Statistical Sciences, Faculty of Arts and Science, University of Toronto, Toronto, Ontario, M5S 3G3, Canada;

[2]Estonian Genome Center, University of Tartu, Tartu 51010, Estonia;

[3]Institute of Molecular and Cell Biology, University of Tartu, Tartu 51010, Estonia;

[4]Population Health Research Institute, Hamilton Health Sciences and McMaster University, Hamilton, Ontario, L8L 2X2, Canada;

[5]Program in Medical and Population Genetics, Broad Institute, Cambridge, USA

[6]Department of Pathology and Molecular Medicine, McMaster University, Hamilton, Ontario, L8N 4A6, Canada;

[7]Division of Biostatistics, Dalla Lana School of Public Health, University of Toronto, Toronto, Ontario, M5T 3M7, Canada.

20

[†]*Corresponding author*:

Lei Sun, Department of Statistical Sciences, 100 St George Street, University of Toronto, Toronto, Ontario, M5S 3G3, Canada. E-mail: sun@utstat.toronto.edu.

1 **ABSTRACT**

2 Genetic markers associated with *variance* of quantitative traits are considered promising

3 candidates for follow-up including interaction analyses. However, as in studies of main effects,

4 X-chromosome is routinely excluded from 'whole-genome' scans due to analytical challenges.

5 Specifically, as males carry only one copy of X-chromosome, the inherent sex-genotype

6 dependency could bias the trait-genotype association, through sexual dimorphism in quantitative

7 traits with sex-specific means or variances. Here we investigate phenotypic variance

8 heterogeneity associated with SNPs on X-chromosome and propose robust strategies. Among

9 those, a generalized Levene's test, adjusting for sex and sex-genotype interaction effects, has

10 adequate power and remains robust to sexual dimorphism. An alternative sex-stratified approach

11 via Fisher's method is the most robust at the cost of slightly reduced power. We applied both

12 methods to an Estonian study of gene expression quantitative trait loci (eQTL; n=841), and two

13 complex trait studies of height, hip and waist circumferences, and body mass index (BMI)

14 collected on Caucasians in UK Biobank (UKB; n=132,968) and multi-ethnic study of

15 atherosclerosis (MESA; n=2,073). Consistent with previous eQTL findings on mean, we found

16 some but not conclusive evidence for *cis* regulators being enriched for variance association.

17 Individual SNP rs148191803 was X-chromosome-wide significant for waist circumference

18 ($p$=2.4E-6) and suggestive for BMI ($p$=1.2E-5) in UKB but not in MESA. However, a

19 permutation study based on MESA showed a trait-specific polygenic model whereby multiple X-

20 chromosome loci collectively influence variance of height ($\lambda_{GC}$=1.14, $p$<1/100), calling for

21 developments of methods to examine broad-sense heritability by incorporating variance loci and

22 quantifying their sex-specific contributions.

23

1 **INTRODUCTION**

2 Several recent reports have examined autosomal genetic loci contributing to phenotypic variance

3 (as opposed to mean) for a wide range of complex traits[1–3], and corresponding methodology

4 development remains an active area of research[4–11]. One possible reason for such phenotypic

5 variance and SNP genotype association, or variance heterogeneity, is that genotype-stratified

6 variances of a trait differ in the presence of gene-gene (*G*x*G*) or gene-environment (*G*x*E*)

7 interactions; both referred to as *G*x*E* hereinafter. For example, rs1358030 (*SORCS1*) was shown

8 to interact with treatment type affecting HbA1c levels in Type 1 Diabetes subjects[12]. And

9 indeed, in a proof-of-principle study where the treatment information was intentionally masked,

10 the SNP was then demonstrated to be associated with variance of HbA1c[10]. Conversely, because

11 direct *G*x*E* modeling may not be feasible in an initial whole-genome scan, the question was then

12 raised as to whether SNPs having effects on the variance of a trait make good candidates for

13 follow-up interaction testing[3]. For instance, rs7202116 (*FTO* as the nearest gene) was

14 significantly associated with variance of BMI[2], and at the same locus, rs1121980 (*FTO)* showed

15 evidence for a statistical interaction with physical activity influencing the mean of BMI[13,14]; it is

16 worth noting that un-modeled interaction induces variance heterogeneity, but the causes of

17 variance heterogeneity are multifaceted[8,10,15–18]. In practice, although it is possible that an

18 interacting SNP has a stronger effect on variance than on mean, as in the case of rs12753193

19 (*LEPR*) interacting with BMI in the prediction of CRP levels in the absence of detectable main

20 effect[1], a more powerful approach to selecting association candidates is to jointly evaluate their

21 mean and variance effects[8,10,19,20].

22 Despite enthusiasm to discover SNPs with variance effects and the availability of

23 statistical tests, variance heterogeneity has not been formally explored for SNPs on X-

1   chromosome (XCHR). As in the conventional 'genome-wide' (mean) association studies[21], the

2   reluctance to include XCHR is due to analytical challenges[21,22]. They range from technical

3   difficulties in genotype calling to statistical complexities in imputation and association (e.g.

4   model uncertainty involving random or skewed X-inactivation[23–26] and sex as a potential

5   confounder). Solutions to overcome some of these challenges had been provided, but all in the

6   context of genetic association analysis of main effects[25,27–30].

7        Here we focus on understanding the impact of the inherent sex-genotype dependency on

8   variance heterogeneity association analysis, and when the trait of interest has sex-specific mean

9   or variance values for males and females. In practice, sexual dimorphism is consistently

10  observed. For example, based on the UK Biobank (UKB)[31] and Multi-Ethnic Study of

11  Atherosclerosis[32] (MESA) data, height displays a sex-specific difference in mean, hip

12  circumference differs in variance, while body mass index (BMI) and waist circumference

13  contrast in both mean and variance between males and females (Figure 1). These empirical

14  patterns of sexual dimorphism vary according to the underlying physiology of the trait, which

15  might or might not be related to genes. Thus, association analyses of phenotypic mean or

16  variance with XCHR SNPs could be biased if these potential sex-specific main or variance

17  effects were not appropriately accounted for.

18                           Figure 1 here

19       For an autosomal SNP, evaluating differences in phenotypic variance across the three

20  genotype groups can be readily achieved by the classical Levene's test for variance

21  heterogeneity[33]. SNPs with significant variance association $p$-values are then selected as likely

22  candidates for follow-up interaction studies. However, the same strategy to prioritize SNPs on

23  XCHR can be problematic, because sex-specific mean and variance differences could create

1  spurious variance heterogeneity unrelated to the putative *GxE* interactions. Thus, the correct

2  formulation of variance test is dependent on a proper formulation of sex effect with respect to

3  both mean and variance.

4  In this paper, we explicitly model the possible sources of confounding related to sex, and

5  propose two general testing strategies that strike a balance between power and robustness. Using

6  extensive simulations, we demonstrate the danger of directly applying autosomal methods to

7  XCHR that would otherwise be suitable for testing variance heterogeneity, and we conclude that

8  special consideration for sex-genotype dependence must be made for XCHR to maintain correct

9  type I error rates. Application studies include identifying SNPs associated with variances of

10  height, BMI, hip and waist circumference using the UK Biobank and MESA data, as well as

11  detecting loci associated with variance of expression quantitative traits using data from the

12  Estonian Genome Center at the University of Tartu (EGCUT) cohort[34–36].

13

14  **MATERIAL AND METHODS**

15

16  **Notation and model setup**

17  Of interest is a quantitative trait *Y*, assumed to be (approximately) normally distributed or had

18  been inversely transformed to resemble a normal distribution. Without loss of generality,

19  consider the following linear model for the '*true*' association relationship between *Y* and a

20  SNP,

21
$$Y = \beta_0 + \beta_G G + \beta_S S + \beta_{GS} GS + \beta_E E + \beta_{GE} GE + \beta_{SE} SE + \beta_{GSE} GES + \varepsilon, \qquad (1)$$

22  where *G* denotes the SNP genotype (coded additively with respect to the number of the minor

23  allele as 0, 1 and 2 for *bb*, *Bb* and *BB* as in convention[37]), *S* is the sex indicator variable (female

1    $= 0$ and male $= 1$), $E \sim N(0, 1)$ is a standardized continuous covariate following the classical *G-E*

2    independence assumption[38], and the error term $\varepsilon \perp N(0, 1)$ is independent of $G$, $S$ and $E$. The

3    minor allele frequency (MAF) of $G$ is assumed to be the same for male and females; sex-specific

4    MAF affects the naïve methods and we will return to this point in the Discussion section.

5         Under these assumptions, it is possible to identify autosomal SNPs potentially involved

6    in *GxE* or high-order interactions, without having to measure $E$ directly, through detecting

7    phenotypic variance associated with $G$ via the *working* model of $Y \sim G$. Note that the analytical

8    context here is that direct *GxE* (or *GxG*) modeling may not be possible (e.g. $E$ may not be known

9    or measured precisely) or desirable (e.g. due to computational or multiple hypothesis testing

10   concerns for whole-genome *GxG* scans). To see the rationale behind the working model, with

11   the additional assumption of conditional independence between $E$ and $S$ conditional on $G$, one

12   can show that the conditional variance of $Y$ on $G$ is,

$$\text{Var}(Y|G = g) = (\beta_E + \beta_{GE}g)^2 + (\beta_S + \beta_{GS}g)^2 \text{Var}(S|G = g) +$$

13   $$[(\beta_{SE} + \beta_{GSE}g)^2 + 2(\beta_E + \beta_{GE}g)(\beta_{SE} + \beta_{GSE}g)]E(S|G = g) + 1. \quad (2)$$

14   Since sex $S$ is independent of $G$ for an autosomal SNP, $\text{Pr}(S|G = g)$ is constant across $g = 0, 1$

15   and 2, so are *E(S/G=g)* and *Var(S/G=g)*. Thus, if $\beta_{GS} = \beta_{GE} = \beta_{GSE} = 0$, expression (2) can be

16   reduced to a constant with respect to $G$:

17   $$\text{Var}(Y|G = g) = \beta_E{}^2 + \beta_S{}^2 \text{Var}(S|G = g) + [\beta_{SE}{}^2 + 2\beta_E\beta_{SE}]E(S|G = g) + 1 \quad (3)$$

18   $$= \beta_E{}^2 + \beta_S{}^2 \text{Var}(S) + [\beta_{SE}{}^2 + 2\beta_E\beta_{SE}]E(S) + 1. \quad (4)$$

19   Conversely, variation in $\text{Var}(Y/G=g)$ across $G$ suggests that at least some of the (un-modeled)

20   interaction terms involving $G$ (i.e. $\beta_{GS}$, $\beta_{GE}$ and $\beta_{GSE}$) are non-zero. This was precisely the

21   motivation behind the original idea[1] of using Levene's test to identify variance heterogeneity

22   induced by the underlying but un-modeled *GxE* interaction.

23

1   **X-chromosome (XCHR) specific challenges for variance tests**

2   The same approach to draw similar conclusions for XCHR SNPs, however, is questionable,

3   because $\Pr(S|G = g)$ is no longer constant in $G$ and expression (3) cannot be further reduced to

4   (4). For example, under the X-inactivation coding of 0, 1 and 2 for the *bb*, *Bb* and *BB* genotypes

5   in females and 0 and 2 for the *b* and *B* genotypes in males, the *G=1* group contains only females.

6   Similarly, under the no X-inactivation coding of 0, 1 and 2 for females and 0 and 1 for males, the

7   *G=2* group then contains only females. Thus, omitting the sex indicator *S* from the covariates

8   can bias the conclusion through sexual dimorphism as seen in Figure 1.

9   Consider the simplest case of no interaction effects at all ($\beta_{GS} = \beta_{GE} = \beta_{SE} = \beta_{GSE} = 0$) nor

10  environmental main effect ($\beta_E = 0$), but there is a sex main effect ($\beta_S \neq 0$, i.e. the sex-stratified

11  phenotypic means differ between males and females), then expression (3) is reduced to

$$\text{Var}(Y|G = g) = \beta_S{}^2 \text{Var}(S|G = g) + 1. \qquad (5)$$

13  Thus, in the absence of any interactions that involve *G*, there is a spurious phenotypic variance

14  heterogeneity across levels of *G* through a non-zero sex main effect ($\beta_S$), or through a sex-

15  environment interaction effect ($\beta_{SE}$) if present as in equation (3). Severity of the confounding

16  depends on the discrepancy between the two sex-stratified trait distributions (real data in Figure

17  1 and conceptual data in Figures 2 A-D), as well as on the strength of correlation between sex

18  and the observed genotype, which in turn depends on the MAF and proportions of males and

19  females in a sample (details in Supplemental Data).

20  Figure 2 here

21  To avoid spurious variance heterogeneity signals, alternative approaches are needed to

22  quantify variance differences induced by *GxE* or higher order interactions involving *G*. To this

7

1  end, it is important to appropriately define the null hypothesis of variance homogeneity that

2  corresponds to an absence of phenotypic variance associated with genotype while allowing for

3  variance (and mean) to differ between males and females (Figures 2 A-D).

4

5  **X-chromosome (XCHR) variance heterogeneity tests**

6  Here we consider various analytical strategies to assess phenotypic variance associated with

7  genotypes of XCHR SNPs, including naïve methods that directly apply the original Levene's test

8  to different genotype groups, and alternative approaches that utilize a generalized Levene's test

9  derived from a two-stage regression framework[20,39].

10

11  *Naïve methods: apply Levene's test to three or five genotype groups*

12  The original Levene's test for variance heterogeneity treats an autosomal genotype $G$ as a

13  categorical variable[33,39] and examines any variance difference in trait $Y$ amongst the three

14  possible genotype groups. A direct application to XCHR, however, is problematic. Because sex

15  $S$ is inherently correlated with XCHR $G$, so any potential correlation between $S$ and $Y$ (e.g. as

16  observed in human height) would create the classic case of confounding. To see this, consider

17  the null situation where $G$ is *not* associated with *variance* of $Y$ as in the top panel of Figure 2.

18  Assume the X-inactivation coding of $G$ was used (same conclusion for the no X-inactivation

19  coding), the *Bb* group contains only females and its variance would be the same as $\sigma_f^2$, reflected

20  by variance of the orange curve in the figure. In contrast, the other two groups (*bb+b* and

21  *BB+B*) contain both males and females, and their respective variance values involve both the

22  orange and blue curves depending on sex-specific means ($\mu_m$ and $\mu_f$) and variances ($\sigma_m^2$ and $\sigma_f^2$),

23  as well the proportion of males in each group. Thus, in the presence of sexual dimorphism

1. (either in mean-Figure 2B or variance-Figure 2C or both-Figure 2D), there would be spurious

2. variance heterogeneity resulting in increased false positive rates (also confirmed by empirical

3. results).

4.     As an alternative, one may be tempted to treat each genotype and sex combination as one

5. group, resulting in a total of 5 groups. Indeed, this five-group strategy does not induce spurious

6. association in the presence of sex-specific mean effect ($\mu_f \neq \mu_m$ as in Figure 2B). However, it is

7. not difficult to see that the problem remains when there is a sex-specific variance effect ($\sigma_m^2 \neq \sigma_f^2$

8. as in Figures 2C or 2D).

9.

10. *Fisher's method: combine sex-stratified Levene's test*

11. Sex-stratified analysis provides a practical strategy whereby variance heterogeneity is assessed

12. separately in males (two-group Levene's test) and females (three-group Levene's test). Fisher's

13. method can then be used to combine the two *p*-values. Note that Leven's test statistic is

14. asymptotically $\chi^2_{\#groups-1}$ distributed without apparent 'direction of effect', so the traditional

15. meta-analysis that combines the weighted (directional) Z-values for testing mean effect is not

16. applicable here. Though sex-stratified analysis does not allow direct *GxS* modeling, it is robust

17. to various forms of sexual dimorphism as seen in Figures 2B-D, and it does not require the

18. knowledge of X-inactivation status.

19.

20. *Model-based generalized Levene's test: account for sex-specific mean and variance effects via*

21. *two-stage regression models.*

22. Since the null hypothesis is defined in terms of phenotypic variance heterogeneity induced by

23. (un-modeled) *GxE* interactions while allowing for sexual dimorphism (Figures 2B-D), a

1   preferred method should explicitly account for the effect of sex on the phenotype of interest. To

2   this end, we consider the generalized Levene's test that established a flexible two-stage

3   regression framework[20,39]. In essence, stage one regresses $Y$ on $G$ and obtains the absolute

4   residual $d$ (i.e. absolute of deviation between observed and model fitted $Y$ values). Stage two

5   regresses $d$ on $G$ again, and testing the slope was shown to be equivalent to evaluating variance

6   heterogeneity in $Y$ associated with $G$ because the expectation of $d$ linearly depends on variance

7   of $Y$ [20,39].

8       The generalized Levene's test has been used to study autosomal SNPs with more

9   complex data structures including genotype group uncertainty (e.g. imputed SNPs) or sample

10  dependency (e.g. correlated family members)[20]. For XCHR analysis, the implementation

11  requires additional care. For example, it is not immediately clear if $S$ (or $G$x$S$) should be

12  included in both stages. For a comprehensive evaluation, we consider all combinations of the

13  following two-stage models:

14      **Stage One: Mean models,**

15  $$Y \sim \alpha_0 + \alpha_G G \text{ (M1)},$$

16  $$Y \sim \alpha_0 + \alpha_G G + \alpha_s S \text{ (M2)},$$

17  $$Y \sim \alpha_0 + \alpha_G G + \alpha_s S + \alpha_{GS} GS \text{ (M3)}.$$

18      **Stage Two: Variance models,**

19  $$d \sim \gamma_0 + \gamma_G G \text{ (V1)},$$

20  $$d \sim \gamma_0 + \gamma_G G + \gamma_s S \text{ (V2)},$$

21  $$d \sim \gamma_0 + \gamma_G G + \gamma_s S + \gamma_{GS} GS \text{ (V3)}.$$

22      The models in stage one are only used to calculate residuals, using either the traditional

23  ordinary least squares (OLS) or the recommended least absolute deviations (LAD); LAD is more

10

1    robust to data with asymmetric distributions or low genotype counts in a specific group[20,40]. The

2    goal of this stage is to remove any *Mean* effects associated with the covariates included in the

3    model (i.e. *G*, *S* or *G*x*S*), thus denoted as M1, M2 or M3.

4         Test for *Variance* heterogeneity is achieved in stage two (V1, V2 or V3), by testing

5    $H_o: \gamma_G = 0$ or $H_o: \gamma_G = \gamma_{GS} = 0$ via the standard regression *F*-test, where model is fitted using

6    OLS for independent samples or generalized least square for dependent samples.

7         The model-based regression approach includes a total of nine M+V two-stage models,

8    and V3 also allows a two degrees of freedom (d.f.) test (Table S1).  Based on the earlier

9    discussion, it is expected that mean modeling strategies omitting *S* (i.e. M1) would be sensitive

10   to sex-specific mean effect (e.g. Figures 2B or 2D).  Meanwhile, variance testing strategies

11   omitting *S* (i.e. V1) are anticipated to be sensitive to sex-specific variance effect (e.g. Figures 2C

12   or 2D).  For completeness of our empirical validation, we first examined all 12 testing strategies

13   in simulation studies then focused on the robust approaches in applications.

14

15

16   **Simulation studies**

17   A sample of 2,500 females and 2,500 males were simulated, and the MAF was fixed at 0.2; other

18   sample sizes and MAFs led to qualitatively similar results.  Note that although *G* could be coded

19   assuming X-inactivation or no X-inactivation, the two types of coding are generally highly

20   correlated leading to similar association results[29].  Thus, for a more focused study here the

21   genotype was simulated assumed no X-inactivation.  The number of simulated replicates was

22   10,000 so that estimates of the empirical Type I Error (T1E) rates within $\pm 0.5\%$ of the nominal

23   rate of 5% were considered satisfactory.

1    A joint mean and variance test can be more powerful than testing for variance

2    heterogeneity alone, but the power of the joint test depends on the individual components[10].

3    Therefore, here we focus on comparing the different variance-testing strategies as outlined

4    above, recommending the most robust yet powerful method that is also suitable for the joint

5    location-scale test.

6

7    *Simulations for T1E evaluation - design I based on model (1)*

8    The genotype-phenotype relationship was generated according to model (1), where the

9    environmental variable $E \sim N(0, 1)$ was used in generating observed phenotypic values but

10   assumed not being available for the actual association analysis.  The null scenarios were defined

11   by the absence of interaction effects for *GxE* and *GxExS*, so the quantitative trait for each null

12   scenario was generated assuming $\beta_{GE} = \beta_{GES} = 0$ in model (1).  A SNP could have a *G* main

13   effect, but it does not affect the phenotypic variance of interest which is induced by un-modeled

14   $\beta_{GE}$ and $\beta_{GES}$ in the working model, so $\beta_G = 0$ without loss of generality.  Note that the naïve

15   variance methods could also pick up a non-zero *GxS* interaction effect if $\beta_{GS} \neq 0$, but $\beta_{GS}$ itself

16   in fact can be directly tested because gender information is routinely collected (or robustly

17   inferred from the available genotype data).  Thus, $\beta_{GS}$ is not related to the variance heterogeneity

18   of interest here and was set to be zero.  For the remaining parameters, without loss of generality,

19   $\beta_0 = 0$, $\beta_E = 0$ or 0.5, $\beta_S = 0$ or 0.5, and $\beta_{SE} = 0$, -0.25 or 0.25, giving a total of 12 scenarios.

20   They roughly fall into four categories, corresponding to the four conceptual sex-stratified

21   distributions as shown in Figures 2A-D.  For example, sexual dimorphism was introduced via $\beta_S$

22   and $\beta_{SE}$, where a none-zero $\beta_S$ allows for *sex*-specific mean effect (Figures 2B and 2D) while a

23   non-zero $\beta_{SE}$ allows for *sex*-specific variance effect (Figures 2C and 2D).  Note that both $\beta_S$ and

1    $\beta_{SE}$ are independent of the *genotype*-specific variance effect to be identified, which is absent in

2    the null cases.

3

4    *Simulations for T1E evaluation – design II based on sex-stratified mean and variance*

5    The null scenarios based on model (1) may not fully capture the extremes of sexual dimorphism,

6    thus we further simulated trait values directly according to sex-specific distributions using means

7    ($\mu_m$ and $\mu_f$) and variances ($\sigma_m^2$ and $\sigma_f^2$) that mimic the values observed in inverse-normally

8    transformed BMI, height, hip and waist circumference from MESA (Table S2). The simulated

9    traits, generated independent of any genotypes, were then tested for variance association with

10   genotypes of 12,206 XCHR SNPs from the MESA dataset, after filtering by a minimum count of

11   30 observations in the five sex-genotype stratified groups as variance test is sensitive to small

12   group size.

13

14   *Simulations for power study*

15   Only strategies with satisfactory T1E control were considered for power evaluation. We focused

16   on model-based design I where power directly depends on the size of *GxE* and *GxExS*

17   interaction effects and has a clearer genetic interpretation than design II. Under model (1), $\beta_{GE}$

18   was varied from 0 to 0.2 with a 0.025 incremental increase, and combined with a possible three-

19   way interaction $\beta_{GES}$ of 0 or 0.1. Other parameter values were the same as in the null case with

20   the exception that $\beta_S = \beta_E = 0$ were not considered for a more focused study of power. That is,

21   $\beta_0 = 0$, $\beta_G = \beta_{GS} = 0$, $\beta_S = \beta_E = 0.5$, and $\beta_{SE} = 0$, -0.25 or 0.25. In total, there were 54 scenarios.

22

23   **Applications**

13

1  Robust variance testing strategies for X-chromosome SNPs that also had reasonable power

2  performance were then applied to real data. Only reportedly unrelated and ethnically Caucasian

3  individuals were included, and diabetic individuals were excluded based on electronic medical

4  records in the UK Biobank[31], and based on blood glucose level greater than 7 mmol/L in

5  MESA[32]. All quantitative traits were quantile-normally transformed to avoid 'scale-effect'

6  where the variance values tend to be proportional to mean values[1,2].

7  The significance level for discovery was set at a nominal level of 5% with Bonferroni

8  correction, depending on the total number of XCHR SNPs examined in each application.

9  Further, the proportion of truly variance-associated variants was estimated using the method

10  proposed by Storey and Tibshirani[41].

11

12  *The UK Biobank (UKB) data*[31]

13  Available genotyped XCHR SNPs were filtered based on whether they were in pseudo

14  autosomal region and a minimal sample count of 30 across the five sex-genotype groups. In

15  total, 7,344 XCHR SNPs on the Caucasian sample (71,452 females and 61,516 males) were

16  analyzed, and the XCHR-wide significance level was 6.8E-6.

17

18  *The Multi-Ethnic Study of Atherosclerosis* (*MESA*) *data*[32].

19  The genotype data in MESA, available from dbGap (Study accession: phs000209.v10.p2), were

20  filtered similarly as the UKB data. In total, 12,206 XCHR SNPs on the Caucasian sample (1,003

21  females and 1,070 males) were analyzed, and the XCHR-wide significance level was 4.1E-6.

22  We did not perform a multi-ethnic analysis with all ethnicities combined. Instead, we focused on

23  the Caucasian subset and used it to corroborate findings from the UKB data.

14

1

2    *Estonian Genome Center at the University of Tartu (EGCUT) cohort*[34–36]

3    We sought to discover XCHR SNPs influencing the variance of expression traits, as variability

4    of gene expression has been suggested to be associated with genetic variants on autosomes[7]. The

5    recommended strategies were applied to a sample of 413 male and 421 female Estonians across

6    648 gene expression traits that had gone through standard quality control procedures and further

7    inversely normal transformed. After filtering using the same criteria as the UKB data, 4,034

8    XCHR SNPs were analyzed for variance association with each of the 648 gene expression traits,

9    resulting in a total number of 2,614,032 tests and a global significance level of 1.9E-8.

10

11    **RESULTS**

12

13    **Simulation studies**

14    As expected, the naïve Levene's test with either a three-level factor $G$ factor or a five-level $GxS$

15    factor grossly overestimated the number of false positives in almost all scenarios except for when

16    $\beta_E$, $\beta_{GS}$ and $\beta_{SE}$ were all set to zero or the in the absence of any sexual dimorphism (Table 1).

17        For generalized two-stage Levene's tests, good choices of the mean model in stage one

18    and variance test in stage two should remove any effect from sex to avoid inflating the test

19    statistics. Thus, as expected, any strategies involving M1 or V1 had T1E issues, where the

20    degrees of departure from the nominal T1E rate varied according to sizes of the unadjusted sex

21    mean or variance effects (Table 1). The remaining strategies appear to have reasonably

22    controlled T1E rates and are underlined in Table 1. However, when considering design II where

23    sexual dimorphism was more extreme, only M2V3.2 and M3V3.2 have good T1E control and

1    behave quite similarly (Table S3).

2        The sex-stratified approach, as expected, gave correct T1E rates in females and males

3    separately, and subsequently in the combined sample via Fisher's methods, under both design I

4    (Table 1) and design II (Table S3).

5        In terms of statistical power among testing strategies with reasonable T1E control

6    (M2V3.2, M3V3.2 and Fisher), M2V3.2 and M3V3.2 are nearly identical and had slightly better

7    power than Fisher's method across most of the scenarios considered (Figures S1).

8        The reason for performance similarity between M2V3.2 and M3V3.2 is because the

9    model was generated (and correctly modeled) under the assumption of no X-inactivation and $\beta_G$

10   = 0.  Interestingly, when there is a strong genotypic main effect ($\beta_G \neq 0$), an increased variance

11   in the female homozygote *Bb* group could be observed as a result of unknown X-inactivation[42].

12   Indeed, additional simulation studies confirmed that variance heterogeneity *p*-values given by

13   M2V3.2 were influenced by X-inactivation while Fisher's method and M3V3.2 remained

14   consistent (Figure S2, Table S4).  Thus, we recommend the M3V3.2 model-based approach and

15   the complementary sex-stratified Fisher's method, which were then applied to the three

16   application datasets.

17

18   *Applications*

19   XCHR-wide analysis of waist circumference showed that the recommended M3V3.2 and

20   Fisher's tests indeed have good T1E control (Figure 3); similar conclusions were drawn based on

21   results of other traits (Figures S3-5, respectively, for BMI, height and hip circumference).  In

22   UKB (Table S5), rs148191803 (*MED14*; in a region known to escape X-inactivation) was found

23   to be associated with waist circumference, based on the M3V3.2 test ($p = 2.08E-06$) and Fisher's

16

1     method ($p$ = 2.4E-06), at the XCHR-wide significance level ($p$ < 6.8E-06), and the same SNP

2     was associated with variance of BMI (M3V3.2 $p$ = 7.01E-06; Fisher $p$ = 1.20E-05; Table S6).

3     However, no SNPs in the *MED14* locus had waist circumference or BMI variance association $p$-

4     values less than 0.05 in MESA.

5        The association of rs2023750 at *TBL1X* with BMI was suggestive (M3V3.2 $p$ = 5.30E-

6     04; Fisher $p$ = 2.61E-04, Table S5) in UKB, and at the same gene locus, rs2521580 was also

7     suggestively associated with BMI (M3V3.2 $p$ = 2.85E-03; Fisher $p$ = 9.04E-04) and waist

8     circumference (M3V3.2 $p$ = 9.77E-03; Fisher $p$ = 5.75E-04) in MESA (Table S6).

9        Although there were no additional X-chromosome-wide significant SNPs in UKB or

10     MESA, the overall distributions of the $p$-values suggest enrichment of associated variants for

11     some of the traits. In Figure S4 for example, the estimated genomic lambda $\lambda_{GC}$ based on the

12     M3V3.2 variance test for height was 1.028 and 1.194, respectively in UKB and MESA, and it

13     was 1.014 and 1.186 based on Fisher's method. The estimated proportion of truly associated

14     SNPs also suggested signal enrichment (Table S7).

15        To benchmark the observed estimates, a permutation analysis was conducted using the

16     individual-level MESA data available to us. Each quantitative trait under the study was

17     permuted within the two sex strata, independently, 100 times. For each permutated dataset,

18     Fisher's method and M3V3.2 were applied and the corresponding $\lambda_{GC}$ values were then

19     calculated (Figure 4). For height, the permuted values, as expected, centered around $\lambda_{GC}$ = 1.

20     The observed $\lambda_{GC}$ value, based on either the M3V3.2 or Fisher's method was larger than all 100

21     null values (Figure 4), supporting the apparent enrichment of XCHR SNPs associated with

22     variance of height. The same conclusion holds, but to a lesser extent for waist circumference.

23     However, for hip circumference and BMI, the observed estimates were not visibly different from

17

1   the null estimates obtained from the permuted datasets.  Similar observations were made based

2   on the estimated proportion of truly associated SNPs (Figure S6).

3          For the eQTL analysis, we observed various forms of sexual dimorphism in expression

4   traits.  In total, 182 out of the 648 expression traits had $p$ <0.05 based on either a $t$-test of

5   equality of means or an $F$-test for equality of variance between the two sexes (Figure S7).

6   Among the eQTLs, the top five variance-associated SNPs belonged to three genes, *FTX*, *PLAC1*

7   and *TEX11* (Figures S8-10), but no SNPs passed the strict Bonferroni correction at $p < 1.9E-8$

8   (Table S8).  There was no apparent enrichment of association globally over all SNP-expression

9   2,614,032 (= 648*4,034) $p$-values (Figures S11-12).  However, upon further investigation based

10  on stratifying SNPs and gene expression pairs according to whether they were *cis* or *trans* acting

11  (using a physical distance of 5Mbps from the start and the end of the gene for each expression

12  trait), we found that the estimated proportion of truly associated SNP-expression pairs appear to

13  be slightly higher for SNPs in *cis* (0.006 and 0.013 based on $p$-values of Fisher's and M3V3.2

14  methods, respectively), as compared to those in *trans* (0 based on either).  The estimated lambda

15  control was 0.996 and 1.003 for *cis*-acting pairs using Fisher's and M3V3.2 methods,

16  respectively, and both at 0.99 for SNP-expression pairs in *trans*, suggesting additional studies are

17  needed to establish convincing evidence for enrichment of variance-associated eQTLs.

18

19  **DISCUSSION**

20  This work was motivated by the recent call to include X-chromosome (XCHR) in 'whole-

21  genome' scans[21], as well as the recent development of identifying autosomal SNPs associated

22  with phenotypic variance[1,20].  To pave the way for future XCHR-wide study of variance

23  heterogeneity and subsequent joint location-scale test[10], we examined a catalogue of analytical

18

1    strategies and recommended two robust and power approaches. We emphasize the importance of

2    recognizing sex as an inherent confounder in analyzing XCHR variants that contribute to

3    phenotypic variance heterogeneity, particularly for traits displaying sexual dimorphism with

4    either sex-specific means or variances, or both; this also holds for the traditional association

5    analysis of XCHR variants studying their effects on phenotypic mean[22].

6         Between the two recommended strategies, Fisher's method to combine sex-specific

7    Levene's *p*-values is intuitive and the most robust, but at the cost of slightly reduced power.

8    Through exploiting the recently proposed generalized Levene's test based on a two-stage

9    regression approach, the model-based M3V3.2 test can directly account for sex main effect as

10    well as *G*x*S* interaction effect. The model-based regression testing strategy also allows flexible

11    adjustments for other potential confounder such as principal components[43]. Thus, we

12    recommend in practice to apply both methods to complement each other.

13         The naïve strategies that directly test for variance heterogeneity across either the classical

14    three genotype groups or the sex-stratified five groups are inadequate with grossed inflated T1E

15    rates in the presence of any sexual dimorphism (Table 1). Note that if the status of X-

16    inactivation were known *a priori*, a non-additive variance model (VNA) in stage two may be

17    considered:

18 $$d \sim \gamma_0 + \gamma_{G1} G1 + \gamma_{G2} G2 \; (\textbf{NAV1}),$$

19 $$d \sim \gamma_0 + \gamma_{G1} G1 + \gamma_{G2} G2 + \gamma_s S \; (\textbf{NAV2}),$$

20 $$d \sim \gamma_0 + \gamma_{G1} G1 + \gamma_{G2} G2 + \gamma_s S + \gamma_{G1S} G1S \; (\textbf{NAV3}),$$

21    where $G1$ and $G2$ are indicator variables for the *Bb* and *BB+B* groups under X-inactivation,

22    respectively, or alternatively for the *Bb+B* and *BB* groups without X-inactivation. Under this

23    representation, Lev3 is equivalent to M1VNA1 in which the main effect of sex is not account for

1 in stage one; while Lev5 is equivalent to the M3VNA3 model-based testing of $\gamma_{G1} = \gamma_{G2} = \gamma_S =$

2 $\gamma_{G1S} = 0$ in which variance heterogeneity due to sex $\gamma_S$ is erroneously being tested. These

3 observations clearly revealed the source of bias inherent in the naïve methods.

4       The additive coding for $G$ in stage one is believed to sufficiently capture the genetic main

5 effect[37], while it may not be the case for analysis of variance. The ambiguous genotype grouping

6 under unknown X-inactivation status adds another layer of complexity for non-additive variance

7 models (Figure S13). In fact, the choice of reference allele coding matters in XCHR when the

8 mean association model does not include the $G$x$S$ interaction terms as shown recently[29]. This

9 has direct consequences for XCHR variance testing since a variance difference in the female

10 homozygote group could be observed as a result of a strong marginal effect coupled with

11 unknown X-inactivation[42]. Additional simulation results under model (1) with a non-zero

12 genetic main effect suggested that $p$-values derived from M2V3.2 varied according to the

13 underlying X-inactivation status (Table S4). Though in applications, the genetic main effect

14 would have to be extremely large for the M2V3.2 $p$-values to lead to different conclusions. We

15 also note that X-inactivation and no X-inactivation lead to similar (mean) association results if

16 the $G$x$S$ interaction term is included in the model, which explains the consistent performance of

17 M3V3.2 irrespective of the X-inactivation status.

18       Since variance testing requires larger sample size than mean testing, detecting individual

19 variance signals that are significant at the XCHR-wide or genome-wide level requires studies of

20 very large size that might only be viable through meta-analysis. Meta-analyses of variance

21 heterogeneity[4] for XCHR variants can be conducted in parallel to that of a single study

22 incorporating the analytical strategies proposed for autosomal variants[4].

1        Similar to a polygenic model proposed for association studies of main effects, it is

2    possible that a large proportion of genetic variants, though not individually detectable, could

3    collectively contribute to variance heterogeneity in certain complex traits[44,45]. The permutated

4    genomic control values showed a clear enrichment and point to a possible XCHR polygenic

5    inheritance model for variance of height, which suggests that height could be potentially

6    enriched for gene-environment interactions. Some have suggested that X-linked genes

7    contribute to the sex-specific architecture of complex traits[46], yet the amount of contribution

8    from XCHR SNPs involved in possible *GxE* or higher-order interactions is unclear. On that

9    note, it is interest to point out that a sex-stratified approach could reveal sex-specific enrichment

10    of XCHR SNPs associated with phenotypic variance of, for example, waist circumference in

11    males (Figure 4). It has been reported genetic loci can have gender specific marginal effects on

12    traits such as height[47] and HDL[48], i.e. *GxS* interactions. Results from this study call for new

13    developments of the broad-sense heritability estimation methods that can incorporate variance

14    loci as well as quantify their contributions to sex-specific heritability.

15

16    **SUPPLEMENTAL DATA**

17    Supplemental data include 13 figures, 8 tables, and theoretical derivations.

18

1 cohort. This research was funded by the Canadian Institutes of Health Research (CIHR,

2 201309MOP-310732-G-CEAA-117978) and the Natural Sciences and Engineering Research

3 Council of Canada (NSERC, 250053-2013) to LS. WQD is supported by NSERC Alexander

4 Graham Bell Canada Graduate Scholarship and Ontario Graduate Scholarship.

5

6 **WEB RESOURCES**

7 An implementation of our methods is provided as an open-source and user-friendly R package

8 available on github (https://github.com/WeiAkaneDeng/Xvarhet).

9

10 **REFERENCES**

11 1. Paré, G., Cook, N.R., Ridker, P.M., and Chasman, D.I. (2010). On the use of variance per

12 genotype as a tool to identify quantitative trait interaction effects: a report from the Women's

13 Genome Health Study. PLoS Genet. *6*, e1000981.

14 2. Yang, J., Loos, R.J.F., Powell, J.E., Medland, S.E., Speliotes, E.K., Chasman, D.I., Rose,

15 L.M., Thorleifsson, G., Steinthorsdottir, V., Mägi, R., et al. (2012). FTO genotype is associated

16 with phenotypic variability of body mass index. Nature.

17 3. Shungin, D., Deng, W.Q., Varga, T. V, Luan, J., Mihailov, E., Metspalu, A., Morris, A.P.,

18 Forouhi, N.G., Lindgren, C., Magnusson, P.K.E., et al. (2017). Ranking and characterization of

19 established BMI and lipid associated loci as candidates for gene-environment interactions. PLOS

20 Genet. *13*, e1006812.

21 4. Struchalin, M. V, Dehghan, A., Witteman, J.C., van Duijn, C., and Aulchenko, Y.S. (2010).

22 Variance heterogeneity analysis for detection of potentially interacting genetic loci: method and

23 its limitations. BMC Genet. *11*, 92.

5. Deng, W.Q., and Paré, G. (2011). A fast algorithm to optimize SNP prioritization for gene-gene and gene-environment interactions. Genet. Epidemiol. 1–10.

6. Aschard, H., Chen, J., Cornelis, M.C., Chibnik, L.B., Karlson, E.W., and Kraft, P. (2012). Inclusion of gene-gene and gene-environment interactions unlikely to dramatically improve risk prediction for complex diseases. Am. J. Hum. Genet. *90*, 962–972.

7. Hulse, A.M., and Cai, J.J. (2013). Genetic variants contribute to gene expression variability in humans. Genetics *193*, 95–108.

8. Cao, Y., Wei, P., Bailey, M., Kauwe, J.S.K., and Maxwell, T.J. (2014). A versatile omnibus test for detecting mean and variance heterogeneity. Genet. Epidemiol. *38*, 51–59.

9. Deng, W.Q., Asma, S., and Paré, G. (2014). Meta-analysis of SNPs involved in variance heterogeneity using Levene's test for equal variances. Eur. J. Hum. Genet. *22*, 427–430.

10. Soave, D., Corvol, H., Panjwani, N., Gong, J., Li, W., Boëlle, P.Y., Durie, P.R., Paterson, A.D., Rommens, J.M., Strug, L.J., et al. (2015). A Joint Location-Scale Test Improves Power to Detect Associated SNPs, Gene Sets, and Pathways. Am. J. Hum. Genet. *97*, 125–138.

11. Hong, C., Ning, Y., Wei, P., Cao, Y., and Chen, Y. (2016). A semiparametric model for vQTL mapping. Biometrics.

12. Paterson, A.D., Waggott, D., Boright, A.P., Hosseini, S.M., Shen, E., Sylvestre, M.P., Wong, I., Bharaj, B., Cleary, P.A., Lachin, J.M., et al. (2010). A genome-wide association study identifies a novel major locus for glycemic control in type 1 diabetes, as measured by both A1C and glucose. Diabetes *59*, 539–549.

13. Qi, L., Brage, S., Sharp, S.J., Sonestedt, E., Demerath, E., Ahmad, T., Mora, S., Kaakinen, M., Sandholt, C.H., Holzapfel, C., et al. (2011). Physical Activity Attenuates the Influence of FTO Variants on Obesity Risk□: A Meta-Analysis of 218 , 166. PLoS Med. *8*,.

1   14. Ahmad, S., Rukh, G., Varga, T. V., Ali, A., Kurbasic, A., Shungin, D., Ericson, U., Koivula,

2   R.W., Chu, A.Y., Rose, L.M., et al. (2013). Gene x Physical Activity Interactions in Obesity:

3   Combined Analysis of 111,421 Individuals of European Ancestry. PLoS Genet. *9*, 1–9.

4   15. Paré, G., Mehta, S.R., Yusuf, S., Anand, S.S., Connolly, S.J., Hirsh, J., Simonsen, K., Bhatt,

5   D.L., Fox, K. a a, and Eikelboom, J.W. (2010). Effects of CYP2C19 genotype on outcomes of

6   clopidogrel treatment. N. Engl. J. Med. *363*, 1704–1714.

7   16. Sun, X., Elston, R., Morris, N., and Zhu, X. (2013). What is the significance of difference in

8   phenotypic variability across SNP genotypes? Am. J. Hum. Genet. *93*, 390–397.

9   17. Dudbridge, F., and Fletcher, O. (2014). Gene-Environment Dependence Creates Spurious

10  Gene-Environment Interaction. Am. J. Hum. Genet. *95*, 301–307.

11  18. Wood, A.R., Tuke, M.A., Nalls, M.A., Hernandez, D.G., Bandinelli, S., Singleton, A.B.,

12  Melzer, D., Ferrucci, L., Frayling, T.M., and Weedon, M.N. (2014). Another explanation for

13  apparent epistasis. Nature *514*, E3–E5.

14  19. Aschard, H., Hancock, D.B., London, S.J., and Kraft, P. (2011). Genome-wide meta-analysis

15  of joint tests for genetic and gene-environment interaction effects. Hum. Hered. *70*, 292–300.

16  20. Soave, D., and Sun, L. (2017). A generalized Levene's scale test for variance heterogeneity

17  in the presence of sample correlation and group uncertainty. Biometrics *73*, 960–971.

18  21. Wise, A.L., Gyi, L., and Manolio, T.A. (2013). EXclusion: Toward integrating the X

19  chromosome in genome-wide association analyses. Am. J. Hum. Genet. *92*, 643–647.

20  22. König, I.R., Loley, C., Erdmann, J., and Ziegler, A. (2014). How to Include Chromosome X

21  in Your Genome-Wide Association Study. Genet. Epidemiol. *38*, 97–103.

22  23. Carrel, L., and Willard, H.F. (2005). X-inactivation profile reveals extensive variability in X-

23  linked gene expression in females. Nature *434*, 400–404.

1   24. Ross, M.T., Grafham, D. V, Coffey, A.J., Scherer, S., McLay, K., Muzny, D., Platzer, M.,

2   Howell, G.R., Burrows, C., Bird, C.P., et al. (2005). The DNA sequence of the human X

3   chromosome. Nature *434*, 325–337.

4   25. Wang, J., Yu, R., and Shete, S. (2014). X-Chromosome Genetic Association Test Accounting

5   for X-Inactivation, Skewed X-Inactivation, and Escape from X-Inactivation. Genet. Epidemiol.

6   *38*, 483–493.

7   26. Tukiainen, T., Villani, A.-C., Yen, A., Rivas, M.A., Marshall, J.L., Satija, R., Aguirre, M.,

8   Gauthier, L., Fleharty, M., Kirby, A., et al. (2017). Landscape of X chromosome inactivation

9   across human tissues. Nature *550*, 244–248.

10  27. Clayton, D.G. (2009). Sex chromosomes and genetic association studies. Genome Med. *1*,

11  110.

12  28. Clayton, D. (2008). Testing for association on the X chromosome. Biostatistics *9*, 593–600.

13  29. Chen, B., Radu, Craiu, V., and Sun, L. (2017). Bayesian Model Averaging for the X-

14  Chromosome Inactivation Dilemma in Genetic Association Study. arXiv.

15  30. Peter F Hickey, and Bahlo, M. (2011). X Chromosome Association Testing in Genome Wide

16  Association Studies. Genet. Epidemiol. *670*, 664–670.

17  31. Sudlow, C., Gallacher, J., Allen, N., Beral, V., Burton, P., Danesh, J., Downey, P., Elliott, P.,

18  Green, J., Landray, M., et al. (2015). UK Biobank: An Open Access Resource for Identifying the

19  Causes of a Wide Range of Complex Diseases of Middle and Old Age. PLoS Med. *12*, 1–10.

20  32. Bild, D.E., Bluemke, D.A., Burke, G.L., Detrano, R., Diez Roux, A. V., Folsom, A.R.,

21  Greenland, P., Jacobs, D.R., Kronmal, R., Liu, K., et al. (2002). Multi-Ethnic Study of

22  Atherosclerosis: Objectives and design. Am. J. Epidemiol. *156*, 871–881.

23  33. Levene, H. (1960). Robust tests for equality of variances. In Contributions to Probability and

Statistics: Essays in Honor of Harold Hotelling. In In Contributions to Probability and Statistics: Essays in Honor of Harold Hotelling., I. Olkin, ed. (Stanford University Press), pp. 278–292.

34. Leitsalu, L., Haller, T., Esko, T., Tammesoo, M.L., Alavere, H., Snieder, H., Perola, M., Ng, P.C., Mägi, R., Milani, L., et al. (2015). Cohort profile: Estonian biobank of the Estonian genome center, university of Tartu. Int. J. Epidemiol. *44*, 1137–1147.

35. Westra, H.J., Peters, M.J., Esko, T., Yaghootkar, H., Schurmann, C., Kettunen, J., Christiansen, M.W., Fairfax, B.P., Schramm, K., Powell, J.E., et al. (2013). Systematic identification of trans eQTLs as putative drivers of known disease associations. Nat. Genet. *45*, 1238–1243.

36. Metspalu, A. (2004). The Estonian Genome Project. Drug Dev. Res. *62*, 97–101.

37. Hill, W.G., Goddard, M.E., and Visscher, P.M. (2008). Data and theory point to mainly additive genetic variance for complex traits. PLoS Genet. *4*,.

38. Lindström, S., Yen, Y.-C., Spiegelman, D., and Kraft, P. (2009). The impact of gene-environment dependence and misclassification in genetic association studies incorporating gene-environment interactions. Hum. Hered. *68*, 171–181.

39. Gastwirth, J.L., Gel, Y.R., and Miao, W. (2009). The Impact of Levene's Test of Equality of Variances on Statistical Theory and Practice. Stat. Sci. *24*, 343–360.

40. Hines, W.G., and Hines, R.J. (2000). Increased power with modified forms of the Levene (Med) test for heterogeneity of variance. Biometrics *56*, 451–454.

41. Storey, J.D., and Tibshirani, R. (2003). Statistical significance for genomewide studies. Proc. Natl. Acad. Sci. U. S. A. *100*, 9440–9445.

42. Ma, L., Hoffman, G., and Keinan, A. (2015). X-inactivation informs variance-based testing for X-linked association of a quantitative trait. BMC Genomics *16*, 241.

1   43. Price, A.L., Patterson, N.J., Plenge, R.M., Weinblatt, M.E., Shadick, N. a, and Reich, D.

2   (2006). Principal components analysis corrects for stratification in genome-wide association

3   studies. Nat. Genet. *38*, 904–909.

4   44. Yang, J., Weedon, M.N., Purcell, S., Lettre, G., Estrada, K., Willer, C.J., Smith, A. V,

5   Ingelsson, E., Connell, J.R.O., Mangino, M., et al. (2011). Genomic inflation factors under

6   polygenic inheritance. Eur. J. Hum. Genet. 807–812.

7   45. Purcell, S.M., Wray, N.R., Stone, J.L., Visscher, P.M., O'Donovan, M.C., Sullivan, P.F.,

8   Ruderfer, D.M., McQuillin, A., Morris, D.W., Oĝdushlaine, C.T., et al. (2009). Common

9   polygenic variation contributes to risk of schizophrenia and bipolar disorder. Nature *460*, 748–

10  752.

11  46. Weiss, L.A., Pan, L., Abney, M., and Ober, C. (2006). The sex-specific genetic architecture

12  of quantitative traits in humans. Nat Genet *38*, 218–222.

13  47. Randall, J.C., Winkler, T.W., Kutalik, Z., Berndt, S.I., Jackson, A.U., Monda, K.L.,

14  Kilpelainen, T.O., Esko, T., Magi, R., Li, S., et al. (2013). Sex-stratified Genome-wide

15  Association Studies Including 270,000 Individuals Show Sexual Dimorphism in Genetic Loci for

16  Anthropometric Traits. PLoS Genet. *9*,.

17  48. Ober, C., Loisel, D. a, and Gilad, Y. (2008). Sex-specific genetic architecture of human

18  disease. Nat. Genet. *9*, 911–922.

19

20

21

22

23

1 **Table 1. Empirical type I error (T1E) rates of X-chromosome variance heterogeneity tests**
2 **under Simulation Design I.**

3 A quantitative trait was simulated according to simulation design I based on linear regression
4 model (1) with coefficient values specified above such that CONDITION 1 captures the null
5 scenario of no sexual dimorphism, i.e. no sex-specific mean nor variance differences as depicted
6 in Figure 2A; CONDITION 2 corresponds to the conceptual null scenario in Figure 2B with the
7 presence of sex-specific means via a non-zero $\beta_S$; CONDITIONS 3A and 3B correspond to
8 Figure 2C, representing a sex-specific variance difference through a non-zero $\beta_{SE}$, the $SxE$
9 interaction effect, where the environmental effect $\beta_E$ takes a value of either 0 (CONDITION 3A)
10 or 0.5 (CONDITION 3B).  Similarly, CONDITIONS 4A and 4B correspond to Figure 2D with
11 sexual dimorphism in both means and variances, with the absence and presence of environmental
12 effect $\beta_E$, respectively.  The total sample size was 5,000 with 2,500 females and 2,500 males,
13 and the MAF was 0.2.  The nominal T1E rate was set to 0.05 and the empirical T1E rates were
14 calculated based on 10,000 simulated replicates.  Those empirical T1E rates exceeding 5%±0.5%
15 were in bold.  Testing strategies that showed satisfactory T1E controls were underlined, and
16 details of the testing strategies are provided in the text and summarized in Table S1.

17

| | CONDITION 1 | | CONDITION 2 | | CONDITION 3A | | CONDITION 3B | | CONDITION 4A | | CONDITION 4B | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\beta_E$ | 0 | 0.5 | 0 | 0.5 | 0 | 0 | 0.5 | 0.5 | 0 | 0 | 0.5 | 0.5 |
| $\beta_S$ | 0 | 0 | 0.5 | 0.5 | 0 | 0 | 0 | 0 | 0.5 | 0.5 | 0.5 | 0.5 |
| $\beta_{SE}$ | 0 | 0 | 0 | 0 | 0.25 | -0.25 | 0.25 | -0.25 | 0.25 | -0.25 | 0.25 | -0.25 |
| LEV3 | 0.0544 | 0.0513 | **0.0620** | 0.0465 | **0.0957** | **0.0754** | **0.0778** | **0.0589** | 0.0552 | **0.0683** | **0.1217** | **0.0794** |
| LEV5 | 0.0472 | 0.0504 | 0.0475 | 0.0482 | **0.8584** | **0.8734** | **0.9983** | **0.1807** | **0.8932** | **0.9053** | **0.9989** | **0.1738** |
| <u>FEMALE</u> | 0.0525 | 0.0525 | 0.0502 | 0.0476 | 0.0528 | 0.0444 | 0.0253 | 0.0514 | 0.0508 | 0.0434 | 0.0484 | 0.0478 |
| <u>MALE</u> | 0.0513 | 0.0485 | 0.0496 | 0.0509 | 0.0496 | 0.0494 | 0.0395 | 0.0518 | 0.0501 | 0.0505 | 0.0508 | 0.0535 |
| <u>FISHER</u> | 0.0517 | 0.0505 | 0.0502 | 0.0478 | 0.0506 | 0.0472 | 0.0286 | 0.0522 | 0.0487 | 0.0454 | 0.0472 | 0.0480 |
| M1V1 | 0.0540 | 0.0460 | 0.0526 | 0.0456 | **0.1130** | **0.0941** | **0.119** | **0.0614** | **0.0639** | **0.0884** | **0.1660** | **0.0709** |
| M1V2 | 0.0535 | 0.0461 | **0.0575** | 0.0424 | 0.0489 | 0.0548 | 0.0375 | 0.0518 | **0.0662** | **0.0566** | 0.0276 | 0.0501 |
| M1V3 | 0.0492 | 0.0472 | **0.1579** | **0.1234** | 0.0459 | 0.0438 | 0.0280 | 0.0538 | **0.1311** | **0.1292** | 0.0292 | **0.1695** |
| M1V3.2 | 0.0508 | 0.0490 | **0.1474** | **0.0990** | 0.0473 | 0.0501 | 0.0320 | 0.0536 | **0.1465** | **0.1269** | 0.0234 | **0.1461** |
| M2V1 | 0.0536 | 0.0455 | 0.0497 | 0.0479 | **0.1136** | **0.0948** | **0.1184** | **0.061** | **0.0998** | **0.1142** | **0.1914** | **0.0610** |
| <u>M2V2</u> | 0.0539 | 0.0459 | 0.0510 | 0.0465 | 0.0484 | 0.0539 | 0.0378 | 0.0526 | 0.0510 | 0.0507 | 0.0318 | 0.0487 |
| <u>M2V3</u> | 0.0485 | 0.0474 | 0.0494 | 0.0499 | 0.0458 | 0.0440 | 0.0278 | 0.0548 | 0.0448 | 0.0433 | 0.0440 | 0.0537 |
| <u>M2V3.2</u> | 0.0514 | 0.0492 | 0.0504 | 0.0492 | 0.0479 | 0.0490 | 0.0315 | 0.0534 | 0.0486 | 0.0471 | 0.0351 | 0.0525 |
| M3V1 | 0.0533 | 0.0456 | 0.0499 | 0.0482 | **0.1109** | **0.092** | **0.1152** | **0.0619** | **0.0971** | **0.1106** | **0.1886** | **0.0617** |
| <u>M3V2</u> | 0.0545 | 0.0458 | 0.0501 | 0.0458 | 0.0489 | 0.0544 | 0.0384 | 0.0528 | 0.0523 | 0.0505 | 0.0313 | 0.0485 |
| <u>M3V3</u> | 0.0483 | 0.0466 | 0.0494 | 0.0511 | 0.0461 | 0.0441 | 0.0267 | 0.0545 | 0.0444 | 0.0428 | 0.0455 | 0.0532 |
| <u>M3V3.2</u> | 0.0516 | 0.0492 | 0.0508 | 0.0493 | 0.0477 | 0.0487 | 0.0314 | 0.0522 | 0.0487 | 0.0477 | 0.0353 | 0.0528 |

18 **Figure 1. Empirical examples of sexual dimorphism: quantitative trait distribution**
19 **stratified by sex.**

28

Phenotype data from UK Biobank (UKB, top row) and the Multi-Ethnic Study of Atherosclerosis (MESA, bottom row) were used to illustrate the possible types of sexual dimorphism as characterized by a location shift in mean of height, a scale difference in the variance of hip circumference, and changes in both mean and variance of waist circumference and BMI. Each trait (*Y*) was inversely normal transformed so the overall distribution (solid curve) is normal with mean 0 and variance 1. Areas under the sex-stratified distributions (dashed curves) are colored by blue for male and orange for female, respectively.


**Figure 2. Defining null and alternative hypotheses for X-chromosome variance heterogeneity test allowing for sexual dimorphism.**

Upper panel: Figures A-D showcase the different types of conceptual *null* distributions, where the variance of a quantitative trait does not vary across the different genotype groups, but is subjected to a possible sex-specific difference in either mean (B), variance (C) or both (D). The black curve is for the overall distribution, and without loss of generality the orange curve is for female and the blue curve is for male as in Figure 1. Lower panel: Figures E-H represent the respective *alternative* distributions. The different genotype groups are marked by different line types and visible only under the alternative conditions when there is phenotypic variance heterogeneity among the genotype groups.


**Figure 3. XCHR-wide variance heterogeneity test results for waist circumference using the UK Biobank (upper panel) and MESA (lower panel) data.**

For each XCHR SNP, the variance heterogeneity *p*-value was calculated using a sex-stratified Fisher's method (grey color) or the model-based M3V3.2 approach (orange color). Manhattan plots (A and D), quantile-quantile plots (B and E), and histograms (C and F) of the *p*-values using data from the UK Biobank are shown on the top row, and on the bottom row for the Multi-Ethnic Study of Atherosclerosis (MESA). In Figure 3A, SNP rs148191803 (M3V3.2 test *p* = 2.08E-06 and Fisher's method *p* = 2.4E-06) in the *MED14* locus was annotated for passing the XCHR-wide significance at 6.8E-06 in the UK Biobank data; but with *p*>0.05 in the MESA data. Results for other traits are in Figures S2-S4 and Tables S5-6.


**Figure 4. Estimated genomic control lambda values using the observed and permuted individual MESA data.**

The permuted dataset was obtained by permuting each quantitative trait stratified by sex, independently, 100 times. The genomic lambda $\lambda_{GC}$ was computed for each trait in a permuted dataset and is shown as a black dot under each test. The red line represents the genomic lambda estimated for the original observed data, and the black horizontal line represents the reference line at of 1. The corresponding estimates of the proportion of truly associated variants are shown in Figure S6.

1

30

**A)**

$H_o$: no genetic variance effect

$\mu_F = \mu_M, \sigma_F = \sigma_M$

**B)**

$H_o$: no genetic variance effect

$\mu_F \neq \mu_M, \sigma_F = \sigma_M$

**C)**

$H_o$: no genetic variance effect

$\mu_F = \mu_M, \sigma_F \neq \sigma_M$

**D)**

$H_o$: no genetic variance effect

$\mu_F \neq \mu_M, \sigma_F \neq \sigma_M$

**E)**

$H_1$: genetic variance effect

$\mu_F = \mu_M, \sigma_F = \sigma_M$

**F)**

$H_1$: genetic variance effect

$\mu_F \neq \mu_M, \sigma_F = \sigma_M$

**G)**

$H_1$: genetic variance effect

$\mu_F = \mu_M, \sigma_F \neq \sigma_M$

**H)**

$H_1$: genetic variance effect

$\mu_F \neq \mu_M, \sigma_F \neq \sigma_M$

SEX
F
M

GENO
0
1
2