# Selection of X-chromosome Inactivation Model

## Jian Wang[1], Rajesh Talluri[2] and Sanjay Shete[1,3]

[1]Department of Biostatistics–Unit 1411, The University of Texas MD Anderson Cancer Center, Houston, TX, USA. [2]Department of Data Science, The University of Mississippi Medical Center, Jackson, MS, USA. [3]Department of Epidemiology, The University of Texas MD Anderson Cancer Center, Houston, TX, USA.

**ABSTRACT:** To address the complexity of the X-chromosome inactivation (XCI) process, we previously developed a unified approach for the association test for X-chromosomal single-nucleotide polymorphisms (SNPs) and the disease of interest, accounting for different biological possibilities of XCI: random, skewed, and escaping XCI. In the original study, we focused on the SNP-disease association test but did not provide knowledge regarding the underlying XCI models. One can use the highest likelihood ratio (LLR) to select XCI models (max-LLR approach). However, that approach does not formally compare the LLRs corresponding to different XCI models to assess whether the models are distinguishable. Therefore, we propose an LLR comparison procedure (comp-LLR approach), inspired by the Cox test, to formally compare the LLRs of different XCI models to select the most likely XCI model that describes the underlying XCI process. We conduct simulation studies to investigate the max-LLR and comp-LLR approaches. The simulation results show that compared with the max-LLR, the comp-LLR approach has higher probability of identifying the correct underlying XCI model for the scenarios when the underlying XCI process is random XCI, escaping XCI, or skewed XCI to the deleterious allele. We applied both approaches to a head and neck cancer genetic study to investigate the underlying XCI processes for the X-chromosomal genetic variants.

**KEYWORDS:** X-chromosome inactivation, skewness, escaping X-chromosome inactivation, nonnested model comparison, Cox test, likelihood ratio

## Introduction

Analyzing X-chromosomal genetic variants is challenging because of the complexity of the X-chromosome inactivation (XCI) process for female X-chromosome loci, for which one of the 2 copies of the X chromosome is randomly inactivated to achieve the dosage compensation of X-linked genes in men and women.[1–10] The XCI process is in general random, that is, on average, 50% of cells have one of the 2 alleles active. However, an additional model of XCI was found to be skewed XCI,[6,8,9,11–16] which is defined as inactivation of one of the alleles in >75% or 90% of cells.[6,8,9,15,17–26] Previous studies for different complex diseases (eg, mental retardation disorders, rheumatoid arthritis) have shown that the skewed XCI pattern can be more common in affected women compared with unaffected women.[13,20,27–29] For example, Chabchoub et al[20] showed that the significantly skewed XCI pattern was observed in a larger proportion of female patients with rheumatoid arthritis (34%) and patients with autoimmune thyroid disease (26%) than among controls (11%). Therefore, it is important to account for XCI skewness when analyzing X-chromosomal genetic markers.[10] Another complexity of the XCI process is escaping XCI, in which both alleles are active (ie, no dosage

compensation), such as the genes on the pseudo-autosomal regions on the X chromosome.[8,30–33] To analyze X-chromosomal genetic variants, one of the most popular software programs for studying genetic associations, PLINK, uses the process of escaping XCI as the only model, for which the effect of the male deleterious allele is assumed to be the same as the effect of the female heterozygote genotype. In this case, the 3 genotypes of women are coded as 0, 1, or 2, and the 2 allele types of men are coded as 0 or 1. Such a coding strategy accounts for the escaping XCI but ignores other biological models such as random and skewed XCI.[34] In another approach, Clayton[35,36] proposed $\chi^2$ tests to analyze data by accounting for only the random XCI model but ignoring the other XCI models. Neither approach accounts for XCI skewness. To account for these different biological models of XCI, we had developed a novel association test to analyze X-chromosomal single-nucleotide polymorphisms (SNPs) that maximizes the likelihood ratio (LLR) over all such biological possibilities of XCI.[10] In that paper, we showed that the LLR association test had higher power than the existing approaches, including the $\chi^2$ test proposed by Clayton and the PLINK regression approach, in the

scenarios where XCI was skewed, whereas the LLR association test lost a little power in scenarios where XCI was random or escaping XCI occurred.

Even though the originally proposed test allows us to evaluate the association between X-chromosomal SNPs and a disease of interest, it does not provide knowledge regarding the underlying XCI models, which is useful information for understanding the contribution of SNPs to the disease, as well as the inheritance patterns of diseases.[37] In particular, the inheritance pattern describes how a disease is transmitted[38] and helps to predict the recurrence risk of a disease for subsequent generations.[39,40] Also, knowing the mode of disease inheritance helps us calculate the risk ratio that is attributed to the SNP[41] for a relative of an affected individual. Such information enables clinicians to provide genetic counseling to family members regarding the likelihood of developing the disease or passing it on to their children.[42] To identify the XCI model given a SNP, one approach is to select the XCI model that has the highest LLR (denoted as the max-LLR approach). However, such an approach does not formally compare the LLRs that correspond to different XCI models to assess whether these models are statistically equivalent. This is especially relevant when an alternate form of XCI may have an LLR that is very close to that of the form of XCI with the highest LLR and so may be statistically equivalent. Moreover, when the study sample size is small, the highest LLR might not correspond to the true underlying XCI model, and therefore, the max-LLR approach would identify the wrong XCI model in the analysis.

Therefore, in the spirit of the Cox test, which is used for comparing 2 nonnested models,[43–45] herein we propose an LLR comparison procedure (denoted comp-LLR) to formally compare the LLRs of different XCI models using statistical tests to select the most likely XCI model that describes the underlying XCI process. We conduct simulation studies to investigate the performance of the comp-LLR approach and compare it with that of the max-LLR approach. The results show that given an X-chromosomal variant significantly associated with the disease of interest, the comp-LLR approach has higher probability of selecting the correct XCI model for scenarios in which the underlying XCI process is random XCI, escaping XCI, or skewed XCI to the deleterious allele. In our simulations, when the maximum LLR corresponded with the true underlying XCI model, the comp-LLR approach performed similar to the max-LLR approach. However, when the maximum LLR did not correspond to the true simulation XCI model, the comp-LLR more often identified the true model as one of the plausible XCI process models, which implies that this approach has higher probability of selecting the correct XCI model. We applied both approaches, max-LLR and comp-LLR, to data from a genetic association study of head and neck cancers to investigate the underlying XCI processes for the X-chromosomal genetic variants.

## Methods

### LLR and the association test

In this study, we considered a di-allelic SNP on the X chromosome and denoted the 2 alleles as $A$ and $a$, where $A$ is the deleterious allele. As described in the original study,[10] we used different coding values for different XCI processes. Specifically, for random and skewed XCI, we, respectively, denoted alleles $a$ and $A$ for men using a random variable $X = \{0, 2\}$ and, respectively, denoted genotypes $aa$, $Aa$, and $AA$ for women as $X = \{0, \gamma, 2\}$, where $\gamma \in [0, 2]$. For escaping XCI, for men, we, respectively, denoted alleles $a$ and $A$ as $X = \{0, 1\}$ and for women, we, respectively, denoted genotypes $aa$, $Aa$, and $AA$ as $X = \{0, 1, 2\}$. In our previous study,[10] we proposed to perform a grid search for the $\gamma$ value to obtain the coding value that maximizes the LLRs. However, we also observed from simulations that it is not necessary to use a small-step function for a grid search as it has very little impact on the LLRs, and multiple comparison issues likely lead to a loss of statistical power. Therefore, we considered 3 coding strategies for random and nonrandom XCI. In particular, the $\gamma$ values of 0, 1, or 2, respectively, represent skewed XCI toward the normal allele, random XCI, or skewed XCI toward the deleterious allele. Including the escaping XCI model, we have a total of 4 coding strategies for the different XCI processes.

Denoting the disease of interest as a binary random variable $Y = \{0, 1\}$, where 0 represents individuals without the disease and 1 represents individuals with the disease, the association between an X-chromosomal SNP $X$ and the disease of interest $Y$ can be written as follows:

$$Logit\left(P\left(Y = 1 \mid X, X_{sex}\right)\right) = \beta_0 + \beta_1 X + \beta_s X_{sex}$$

where $\beta_0$, $\beta_1$, and $\beta_s$ are regression coefficients and $X_{sex}$ is the variable for the sex. As described above, the value of $X \in M$ is based on each individual's sex as well as the underlying XCI processes. Therefore, we define the set $M$ as $M = (X^F, X^M)$, where $X^F$ denotes a set of coding for genotypes $aa$, $Aa$, and $AA$ for women, and $X^M$ denotes a set of coding for 2 allele types, $a$ and $A$, for men. Specifically, we evaluate 4 distinct models based on 4 different coding sets: (1) $M_1$ = random XCI in which $X^F \in \{0,1,2\}$ and $X^M \in \{0,2\}$, (2) $M_2$ = escaping XCI in which $X^F \in \{0,1,2\}$ and $X^M \in \{0,1\}$, (3) $M_3$ = skewed XCI to the normal allele in which $X^F \in \{0,0,2\}$ and $X^M \in \{0,2\}$, and (4) $M_4$ = skewed XCI to the deleterious allele in which $X^F \in \{0,2,2\}$ and $X^M \in \{0,2\}$.

Given a sample data set with $N$ individuals, the likelihood is written as $L(Y \mid X, X_{sex}; \beta_0, \beta_1, \beta_s) = \prod_{i=1}^{N} (\exp(\beta_0 + \beta_1 x_i + \beta_s x_{sexi}) / (1 + \exp(\beta_0 + \beta_1 x_i + \beta_s x_{sexi})))^{y_i} (1 / (1 + \exp(\beta_0 + \beta_1 x_i + \beta_s x_{sexi})))^{1-y_i}$ under the alternative hypothesis that the SNP is associated with the disease and $L(Y \mid X_{sex}; \beta_0, \beta_s) = \prod_{i=1}^{N} (\exp(\beta_0 + \beta_s x_{sexi}) / (1 + \exp(\beta_0 + \beta_s x_{sexi})))^{y_i} (1 / (1 + \exp(\beta_0 + \beta_s x_{sexi})))^{1-y_i}$ under the null hypothesis that the SNP is not

associated with the disease, where $x_i$ and $x_{sexi}$ are the observed values of the X-chromosomal SNP and the sex and $y_i$ are the observed disease status. Therefore, the LLR can be written as a function of different coding of $X$:

$$LLR(X) = \frac{L(Y \mid X, X_{sex}; \beta_0, \beta_1, \beta_s)}{L(Y \mid X_{sex}; \beta_0, \beta_s)}, \quad X \in \mathrm{M}$$

For the SNP-disease association, for each $X$ coded according to a specific XCI model, we first estimate the regression coefficients $\beta_0$, $\beta_1$, and $\beta_s$, as well as the corresponding *LLR*. By enumerating all different models, $X \in M$, we obtain the maximum LLR, $LLR^*$, based on the sample data. The significance of the association between a SNP and the disease was assessed by comparing $LLR^*$ with its null distribution obtained using a permutation-based procedure as described in the original study.[10] In this study, we are interested in selecting the most likely XCI model that describes the underlying XCI process. Note that the max-LLR approach selects such an XCI model (ie, $X$) as the one for which the LLR is maximum (ie, $LLR^*$). The details of the comp-LLR approach are described below.

*Procedure for comparisons of LLRs (comp–LLR)*

To compare 2 LLRs, we use a test statistic that is inspired by the Cox test[43–45] for comparing 2 nonnested models. Two models are considered as nonnested if one model cannot be derived from the other model by parametric restriction or a limiting process, such as imposing constraints on the parameters of one of the models.[46,47] This is the case in our study where the LLRs for different XCI processes use different coding values for the X-chromosomal SNPs. When comparing 2 nonnested models, $M_1$ and $M_2$, the basic idea of the Cox test is that if $M_1$ represents the correct model, then a fit of the regressors from $M_2$ to the fitted values from $M_1$ should not have further explanatory value; if it has, it can be concluded that $M_1$ does not represent the correct model.[43–45,48] We use the R function "coxtest" in the "lmtest" package to conduct the test.[49,50] To compare the 4 LLRs based on 4 different XCI processes, including skewed XCI to the normal allele, random XCI, skewed XCI to the deleterious allele, and escaping XCI, we started with the XCI model with the highest LLR as the most likely model that describes the data and compared it with the other XCI models. The formal comparison procedure is described below:

1. Rank the LLRs as $LR_{(1)} \geq LR_{(2)} \geq LR_{(3)} \geq LR_{(4)}$. We denote the XCI models corresponding to the ranked LLRs as $M_{(1)}$, $M_{(2)}$, $M_{(3)}$, and $M_{(4)}$, respectively. Let $X_{(1)}$, $X_{(2)}$, $X_{(3)}$, and $X_{(4)}$ be the sets of regressors associated with the 4 models, respectively. Note that $X_{(1)}$, $X_{(2)}$, $X_{(3)}$, and $X_{(4)}$ are defined as matrices with 3 columns: the first column of 1's represents the intercept and the rest of the columns represent the SNP based on different XCI models and the individual's sex.

2. Compare the XCI model corresponding to the second highest LLR (ie, $M_{(2)}$) with the XCI model corresponding to the highest LLR (ie, $M_{(1)}$). We test the hypothesis that $X_{(2)}$ is the set of correct regressors, that is, $H_0: Logit(P(Y = 1 \mid X)) = X_{(2)}\beta$, where $\beta = (\beta_0, \beta_1, \beta_s)'$ is the vector of the model coefficients. The statistic for this test can be written as follows[48]:

$$T = N / 2 \ln\left(s^2_{X_{(1)}} / s^2_{X_{(1)}X_{(2)}}\right)$$

where $N$ is the number of individuals, $s^2_{X_{(1)}}$ is the mean squared residual for model $M_{(1)}$, and $s^2_{X_{(1)}X_{(2)}} = s^2_{X_{(2)}} + \hat{\beta}' X'_{(2)} Q_{X_{(1)}} X_{(2)} \hat{\beta} / N$, where $\hat{\beta} = (X'_{(2)}X_{(2)})^{-1} X'_{(2)}Y$, and $Q_{X_{(1)}} = I - X_{(1)}(X'_{(1)}X_{(1)})^{-1} X'_{(1)}$, where $I$ is the identity matrix. The estimated variance of the test statistic is calculated as follows:

$$\hat{V}(T) = \left(s^2_{X_{(2)}} / s^4_{X_{(1)}X_{(2)}}\right) \hat{\beta}' X'_{(2)} Q_{X_{(1)}} Q_{X_{(2)}} Q_{X_{(1)}} X_{(2)} \hat{\beta},$$

where $s^2_{X_{(2)}}$ is the mean squared residual for $M_{(2)}$ and $Q_{X_{(2)}} = I - X_{(2)}(X'_{(2)}X_{(2)})^{-1} X'_{(2)}$. $T / \sqrt{\hat{V}(T)}$ is asymptotically standard, normally distributed. This test can be conducted by computing regressions and retrieving fitted values and means of squared residuals,[48] the details of which are described in Appendix 1.

If the null hypothesis is rejected, that means model $M_{(2)}$ does not have the correct set of regressors, which implies that model $M_{(1)}$ is the most likely model that describes the XCI process. Note that, in this case, we will not compare $M_{(1)}$ with $M_{(3)}$ and $M_{(4)}$ because model $M_{(2)}$ is rejected and models $M_{(3)}$ and $M_{(4)}$ have lower LLR values than model $M_{(2)}$. If the test does not reject the null hypothesis, we consider that models $M_{(1)}$ and $M_{(2)}$ are equivalent and indistinguishable, and go to step 3.

3. Compare the XCI model corresponding to the third highest LLR (ie, $M_{(3)}$) with the XCI model corresponding to the highest LLR (ie, $M_{(1)}$), using the same test described in step 2. Now the null hypothesis is $H_0: Logit(P(Y = 1 \mid X)) = X_{(3)}\beta$. If the null hypothesis is rejected, that means model $M_{(3)}$ does not have the correct set of regressors, which implies that models $M_{(1)}$ and $M_{(2)}$ remain the most likely models that describe the XCI process. If the test does not reject the null hypothesis, we consider that models $M_{(1)}$, $M_{(2)}$, and $M_{(3)}$ are equivalent and indistinguishable, and go to step 4.

4. Finally, we compare the XCI models $M_{(4)}$ and $M_{(1)}$. If the null hypothesis is rejected, that means model $M_{(4)}$ does not have the correct set of regressors, which implies that models $M_{(1)}$, $M_{(2)}$, and $M_{(3)}$ remain the most likely models to describe the XCI process. If the test does not

reject the null hypothesis, we consider that all the models are equivalent and indistinguishable.

*Simulation approach*

We conducted simulation studies to investigate the performance of the comp-LLR and max-LLR approaches for selecting the XCI models that best describe the XCI process. We investigated 2 di-allele X-chromosomal SNPs that were in high linkage disequilibrium (LD) with the causal variant for the disease of interest, which we denoted as $SNP_1$ and $SNP_2$. To generate the genotypes for SNPs in LD, we employed the R package "HapSim," which is a simulation tool for generating haplotype data based on minor allele frequency (MAF) and LD coefficients.[51,52] Specifically, HapSim generates haplotypes for each of the 2 strands of the chromosome based on the LD pattern and then combines the base pairs at each locus to generate the genotypes.[51,52] For the purpose of simulation, the MAF for the causal SNP was assumed to be 40%, whereas the MAF for both $SNP_1$ and $SNP_2$ was assumed to be either 30% or 40% for different scenarios. We considered different LD values among the SNPs, with the LD coefficient $r^2$ set to be either 0.85 or 0.95 for different scenarios. The sex of each individual was generated on the basis of the prevalence of men in the general population (ie, 50%). Because men are hemizygous, the genotypes of the causal SNP generated from HapSim were recoded conditional on sex according to the different XCI processes as discussed above.

Finally, given the data of the causal SNP, the disease of interest, denoted as a random variable $Y$, was simulated using the logistic regression model: $Logit(P(Y = 1 | X, X_{sex})) = \beta_0 + \beta_1 X + \beta_s X_{sex}$. In the logistic model, $X$ represents coding associated with the XCI process for the causal SNP and $X_{sex}$ represents the covariate for sex, where we assumed that sex was also associated with the disease. We fixed the regression coefficients at $\beta_1 = 0.53$ and $\beta_s = 0.34$, which correspond to odds ratios of 1.7 and 1.4, respectively. We assumed that the proportion of men in the general population was 50%. The intercept coefficient $\beta_0$ was set as –5. We considered 4 different XCI processes: skewed XCI to the normal allele, skewed XCI to the deleterious allele, random XCI, and escaping XCI. Given different MAFs for the associated SNPs of interest (ie, $SNP_1$ and $SNP_2$), different LD coefficients and different XCI processes, we had a total of 16 scenarios in the simulation studies (Table 1). For each scenario, we simulated 500 replicates, each with 2000 cases and 2000 controls. For the comp-LLR approach, we had up to 3 tests for each genetic variant; therefore, we used a significance level of .017 (ie, .05/3) to control for the multiple comparison issue. In the simulations, we considered common X-chromosomal SNPs (MAFs = 30% or 40%) assuming an effect size of 0.53 for the SNP-disease association, which has been observed in real genetic association studies.[53,54] The prevalence of the disease was assumed to be ~1% in the general population. Because sex was used as a covariate in the simulation models, the simulation scenarios include different proportions of men and women in cases and controls. In particular, across all scenarios, the proportion of men among the controls varied from ~20% to ~40%.

*Application to head and neck cancer data*

We applied the max-LLR and comp-LLR approaches for identifying the underlying XCI process to the X-chromosomal data from a head and neck cancer genome-wide association (GWA) study. The details of the data, including demographical, genotyping, and quality control information, were described in our original study.[10] In particular, this GWA study involved 2 independent phases. In the phase 1 study, 14 169 SNPs were genotyped on the X chromosome for 2718 individuals, where 1161 were patients with head and neck cancer and 1557 were controls. In the phase 2 study, 14 371 SNPs were genotyped on the X chromosome for 3996 individuals, where 1031 were patients and 2965 were controls. For both phases, the patients with head and neck cancer were accrued at The University of Texas MD Anderson Cancer Center (MD Anderson). For the controls in the phase 1 study, 531 individuals were recruited by MD Anderson for the study of head and neck cancers, and 1026 individuals had been previously recruited at MD Anderson for the study of cutaneous melanoma.[55] For the controls in the phase 2 study, 643 individuals were recruited by MD Anderson and data for 2322 individuals were obtained from the Study of Addiction: Genetics and Environment, provided by the National Center for Biotechnology Information, which were downloaded from dbGaP.[56] The study was approved by the institutional review board at MD Anderson, and written informed consent was obtained from all participants.

## Results
### Simulation results

As described in the "Methods" section, using the proposed comp-LLR approach, we may not be able to identify one single model to represent the underlying XCI model as the max-LLR approach does. Instead, the comp-LLR approach may result in multiple possible XCI models. If the true underlying XCI model is included in the resulting multiple equivalent possible XCI models, we consider that the true model is identified by the comp-LLR approach. For the max-LLR, we consider that the true model is identified if the highest LLR corresponds to the true XCI model. Table 1 lists the proportions of the 500 replicates for which the true underlying XCI model was identified for the associated $SNP_1$ and $SNP_2$ under 16 scenarios, using the 2 approaches.

From the results, we observe that compared with the max-LLR approach, the comp-LLR approach has higher or comparable probability of identifying the true underlying XCI model for both SNPs through all the scenarios.

**Table 1.** Proportions of simulated replicates for which the true underlying XCI model was identified by the max-LLR and comp-LLR approaches for $SNP_1$ and $SNP_2$.

| SCENARIO | XCI MODEL | MAF | LD | MAX-LLR[a] | | COMP-LLR[b] | |
|---|---|---|---|---|---|---|---|
| | | | | $SNP_1$ | $SNP_2$ | $SNP_1$ | $SNP_2$ |
| 1 | Skewed XCI to normal allele | 0.3 | 0.95 | 0.828 | 0.840 | 0.828 | 0.840 |
| 2 | Skewed XCI to deleterious allele | 0.3 | 0.95 | 0.830 | 0.832 | 0.910 | 0.918 |
| 3 | Random XCI | 0.3 | 0.95 | 0.942 | 0.944 | 0.990 | 0.992 |
| 4 | Escaping XCI | 0.3 | 0.95 | 0.596 | 0.614 | 0.960 | 0.948 |
| 5 | Skewed XCI to normal allele | 0.4 | 0.95 | 0.970 | 0.974 | 0.970 | 0.974 |
| 6 | Skewed XCI to deleterious allele | 0.4 | 0.95 | 0.932 | 0.912 | 0.954 | 0.932 |
| 7 | Random XCI | 0.4 | 0.95 | 0.962 | 0.970 | 0.994 | 0.998 |
| 8 | Escaping XCI | 0.4 | 0.95 | 0.688 | 0.698 | 0.952 | 0.940 |
| 9 | Skewed XCI to normal allele | 0.3 | 0.85 | 0.616 | 0.660 | 0.618 | 0.660 |
| 10 | Skewed XCI to deleterious allele | 0.3 | 0.85 | 0.696 | 0.678 | 0.876 | 0.844 |
| 11 | Random XCI | 0.3 | 0.85 | 0.926 | 0.940 | 0.982 | 0.982 |
| 12 | Escaping XCI | 0.3 | 0.85 | 0.444 | 0.474 | 0.938 | 0.936 |
| 13 | Skewed XCI to normal allele | 0.4 | 0.85 | 0.860 | 0.834 | 0.864 | 0.834 |
| 14 | Skewed XCI to deleterious allele | 0.4 | 0.85 | 0.768 | 0.756 | 0.808 | 0.804 |
| 15 | Random XCI | 0.4 | 0.85 | 0.986 | 0.978 | 0.994 | 0.990 |
| 16 | Escaping XCI | 0.4 | 0.85 | 0.500 | 0.496 | 0.926 | 0.934 |

Abbreviations: LD: linkage disequilibrium ($r^2$); LLR, likelihood ratio; MAF, minor allele frequency for the associated SNPs; SNP, single-nucleotide polymorphism; XCI, X-chromosome inactivation.
The proportions were calculated based on 500 replicates, each with 2000 cases and 2000 controls.
[a]If the maximum LLR corresponds to the true XCI model, we consider that the true model is identified using the max-LLR approach.
[b]If the true underlying XCI model is included in the resulting one or multiple equivalent possible XCI models, we consider that the true model is identified by the comp-LLR approach.

When the true underlying XCI model is skewed toward the deleterious allele (scenarios 2, 6, 10, and 14), the comp-LLR approach has higher probability of identifying the true model. For example, in scenario 2 (MAF = 0.3 and $r^2$ = 0.95), the proportions of replicates for which the true XCI model was identified are, respectively, 0.830 and 0.832 for the 2 SNPs when using the max-LLR approach, whereas when using the comp-LLR approach, the proportions are, respectively, 0.910 and 0.918 for the 2 SNPs.

Similarly, when the true underlying XCI model is random (scenarios 3, 7, 11, and 15), we observe a pattern similar to that for the scenarios with XCI skewed toward the deleterious allele. For example, in scenario 11 (MAF = 0.3 and $r^2$ = 0.85), the proportions are, respectively, 0.926 and 0.940 for the 2 SNPs when using the max-LLR approach, whereas when using the comp-LLR approach, the proportions are, respectively, 0.982 and 0.982 for the 2 SNPs.

When the true underlying XCI model is escaping XCI (scenarios 4, 8, 12, and 16), the proposed comp-LLR approach provides substantial gains in identifying the true XCI model compared with that of the max-LLR approach. For example, in

scenario 16 (MAF = 0.4 and $r^2$ = 0.85), the proportions are, respectively, 0.500 and 0.496 for the 2 SNPs when using the max-LLR approach, whereas when using the comp-LLR approach, the proportions are, respectively, 0.926 and 0.934 for the 2 SNPs.

When the true underlying XCI model for the simulations is skewed toward the normal allele (scenarios 1, 5, 9, and 13), the 2 approaches perform similarly, with almost identical proportions for all the scenarios.

*Effect of sample size*

Through simulations, we also investigated the effects of sample sizes on selecting the true underlying XCI model. In particular, we considered scenario 8, in which the true underlying model was escaping XCI. We simulated a data set with 2000 cases and 2000 controls. The highest LLR corresponded with the random XCI model, so the max-LLR approach led to incorrect model selection for the XCI process. In contrast, the comp-LLR approach concluded that the random XCI and the true escaping XCI models were equivalent and indistinguishable.

To further assess the sample size issue, we increased the sample size to 10 000 cases and 10 000 controls. The highest LLR corresponded with the escaping XCI model, so the max-LLR approach led to correct model selection for the XCI process. Similarly, the comp-LLR approach concluded that the true escaping XCI model was the most likely XCI model.

*XCI models for X-chromosomal genetic variants using head and neck cancer data*

The purpose of the analysis of real data is to identify the potential underlying XCI models for the X-chromosomal SNPs reported in our original study, using both approaches, max-LLR and comp-LLR, based on the phase 1 data, phase 2 data, and the combined data merged from both phases of the head and neck cancer GWA study. In our original study,[10] we analyzed the top 33 X-chromosomal SNPs using phase 1 and phase 2 data separately and then performed the meta-analysis based on the results from both phases using the fixed and random effects models, as well as the Fisher method, in which the natural logarithm of *P* values from 2 phases were summed and then multiplied by –2 to create a test statistic, which has a $\chi^2$ distribution when the null hypotheses are true and the *P* values are independent.[57,58] We focused on the 4 top SNPs ranked by the *P* values obtained using the Fisher method. The results for the XCI models identified for the 4 top SNPs are reported in Table S1 of the Supplementary Material. In particular, we were interested in the SNP rs12388803 (meta-analysis *P* = 2.04 × $10^{-6}$ using Fisher method), which is the only SNP that reached the chromosome-wide significance threshold in our original study. When using the max-LLR approach, the SNP rs12388803 was found to have an XCI model that skewed toward the deleterious allele based on phase 1, phase 2, and the combined data. When using the comp-LLR approach, rs12388803 was found to have an XCI model that skewed toward the deleterious allele based on phase 1 data and the combined data but to have equivalent models of random XCI and XCI skewed toward the deleterious allele based on phase 2 data. Of note, the other 3 SNPs were statistically not significant at the chromosome-wide significance level and had LLRs for different XCI models that were very close; consequently, the comp-LLR approach was not able to distinguish different XCI models and resulted in multiple equivalent models.

## Discussion

In this article, we extended the approach we developed in the original study,[10] which was a novel approach to analyze X-chromosomal SNPs, and proposed a comp-LLR approach to select the most likely XCI model that describes the underlying XCI biological process. We performed simulation studies to investigate the performance of the proposed comp-LLR. Our results show that if the SNP is significantly associated with the disease of interest, the comp-LLR approach has higher probability of identifying the true XCI model for scenarios where the

random XCI, skewed XCI to the deleterious allele, and escaping XCI were used as the underlying models for the simulations. When the underlying XCI model is skewed toward the normal allele, both approaches lead to similar conclusions. Note that when the SNP is not associated with the disease of interest, there is no associated mode of inheritance or XCI model; therefore, there is no need to conduct model comparisons for such a SNP. In practice, for each SNP, we first check the SNP-disease association using the association test proposed in our original study, and the comparison procedures proposed in this article are only needed when the SNP-disease association is statistically significant. Note that the proposed comp-LLR focuses on identifying the XCI models; thus, it will not lead to an inflated type 1 error rate for the SNP-disease association test as the significance of the X-chromosomal SNP has already been checked. Although logistic regression was used to formulate the problem and demonstrate the proposed approach, the comp-LLR approach is applicable to other types of data (eg, continuous) or different sampling schemes (eg, case-control study based on frequency matching for covariates such as sex). Note that the unconditional logistic regression, used in this article, is commonly used for frequency-matched case-control studies. In some situations, the likelihood function may need to be modified accordingly (eg, linear regression for continuous phenotype, conditional logistic regression for a pairwise matched sampling design). In genetic association studies, one may observe different MAFs and effect sizes for men and women, as well as different proportions of male:female among the cases and controls.[3,4] To account for such differences between men and women, we recommend always including sex as a covariate in the study of X-chromosomal genetic variants using the proposed approach. Our method does not assume Hardy-Weinberg proportions in women.

The approach proposed here has some assumptions. For example, we assumed that the XCI model corresponding to the highest LLR, $M_{(1)}$, is the most likely XCI model that describes the underlying XCI process and compared this model with the other XCI models. For each comparison (steps 2-4), we also assumed that if the null model (eg, $M_{(2)}$ in step 2) is rejected, then other models (eg, $M_{(3)}$ and $M_{(4)}$) are not compared with $M_{(1)}$ because of the lower LLRs for the other models. In addition, we considered only the analysis of a single locus but not a multifactorial disease model. We also assumed that there are only 4 XCI models represented by 4 different coding strategies.

The sample size has a significant impact on the performance of the 2 approaches. When the sample size is large, both approaches can identify the correct underlying model. However, when the sample size is relatively small, as is commonly observed in many real data sets, the max-LLR approach is more likely to identify the wrong model, whereas the comp-LLR approach still provides a solution with multiple equivalent possible XCI models, including the true underlying model.

Alternative approaches are available to compare 2 non-nested models (eg, $M_{(1)}$ and $M_{(2)}$), such as the encompassing

approach and the $J$ test proposed by Davidson and MacKinnon.[45,48] The encompassing approach creates a "supermodel," which contains all regressors from both models, $M_{(1)}$ and $M_{(2)}$, for comparison. In this case, both models are nested within the "supermodel," and each model is compared with the "supermodel." However, this approach does not specifically distinguish between $M_{(1)}$ and $M_{(2)}$, rather it distinguishes between the "supermodel" and $M_{(1)}$ or $M_{(2)}$[48]; therefore, we did not consider this approach in our study. However, the basic idea of the $J$ test is that if $M_{(2)}$ contains the correct regressors, then inclusion of the fitted values of $M_{(1)}$ into $M_{(2)}$ as a regressor will not significantly improve the fitted values as compared with fitted values based on $M_{(2)}$.[45,48] The $J$ test is also available in the R "lmtest" package.[49,50] We applied the $J$ test to the simulated data and obtained results similar to those obtained from the Cox test. For example, for scenario 4 (MAF = 0.3 and $r^2$ = 0.95), using the $J$ test for comparing LLRs, the proportions of replicates for which the true underlying XCI model was identified were 0.966 and 0.954, respectively, for the associated $SNP_1$ and $SNP_2$. Comparatively, these proportions using the Cox test were 0.960 and 0.948, respectively.

This study has 2-fold improvement over the original study. First, the proposed comp-LLR approach, which compares different XCI models, will help us better understand the underlying XCI models and provide useful information about the contribution of SNPs to diseases and the inheritance model of disease. Second, in this study, we observed that the SNPs in high LD ($r^2$ > 0.8) tend to have the same XCI models. This finding potentially improves the power of the original association test by reducing the number of multiple comparisons. For each X-chromosomal SNP, the original association test needs to conduct 4 tests, each corresponding to one of the 4 XCI biological processes. In this case, the power of the test might be impaired due to the multiple comparison issue. For example, given that $SNP_1$ and $SNP_2$ are in high LD ($r^2$ > 0.8), if the comp-LLR shows that both the random XCI model and the skewed XCI to the normal allele model can describe $SNP_1$, then for $SNP_2$, we do not need to check all 4 XCI models, but instead can investigate only 2 models when testing the association. Thus, when testing tens of thousands of X-chromosomal variants, such a strategy potentially reduces the number of multiple comparisons.

In summary, we have proposed a comp-LLR approach to select the most likely XCI model that describes the underlying XCI biological process, which has higher probability of identifying the true XCI model compared with the naïve approach that uses the highest LLR.

## Author Contributions

Conceived and designed the experiments: JW, SS. Conducted simulation and analyzed the data: JW, RT, SS. Wrote the first draft of the manuscript: JW, SS. Contributed to the writing of the manuscript: JW, SS. Agree with manuscript results and conclusions: JW, RT, SS. Jointly developed the structure and arguments for the paper: JW, SS. Made critical revisions and approved final version: JW, SS. All authors reviewed and approved of the final manuscript.

## REFERENCES

1. Chow JC, Yen Z, Ziesche SM, Brown CJ. Silencing of the mammalian X chromosome. *Annu Rev Genomics Hum Genet*. 2005;6:69–92.
2. Gendrel AV, Heard E. Fifty years of X-inactivation research. *Development*. 2011;138:5049–5055.
3. Hickey PF, Bahlo M. X chromosome association testing in genome wide association studies. *Genet Epidemiol*. 2011;35:664–670.
4. Loley C, Ziegler A, Konig IR. Association tests for X-chromosomal markers—a comparison of different test statistics. *Hum Hered*. 2011;71:23–36.
5. Lyon MF. Gene action in the X-chromosome of the mouse (*Mus musculus* L.). *Nature*. 1961;190:372–373.
6. Minks J, Robinson WP, Brown CJ. A skewed view of X chromosome inactivation. *J Clin Invest*. 2008;118:20–23.
7. Starmer J, Magnuson T. A new model for random X chromosome inactivation. *Development*. 2009;136:1–10.
8. Willard HF. The sex chromosomes and X chromosome inactivation. In: Scriver CR, Beaudet al, Sly WS, Valle D, Childs B, Vogelstein B, eds. *The Metabolic and Molecular Bases of Inherited Disease*. New York, NY: McGraw-Hill; 2000:1191–1221.
9. Wong CC, Caspi A, Williams B, Houts R, Craig IW, Mill J. A longitudinal twin study of skewed X chromosome-inactivation. *PLoS ONE*. 2011;6:e17873.
10. Wang J, Yu R, Shete S. X-chromosome genetic association test accounting for X-inactivation, skewed X-inactivation, and escape from X-inactivation. *Genet Epidemiol*. 2014;38:483–493.
11. Busque L, Paquette Y, Provost S, et al. Skewing of X-inactivation ratios in blood cells of aging women is confirmed by independent methodologies. *Blood*. 2009;113:3472–3474.
12. Chagnon P, Provost S, Belisle C, Bolduc V, Gingras M, Busque L. Age-associated skewing of X-inactivation ratios of blood cells in normal females: a candidate-gene analysis approach. *Exp Hematol*. 2005;33:1209–1214.
13. Plenge RM, Stevenson RA, Lubs HA, Schwartz CE, Willard HF. Skewed X-chromosome inactivation is a common feature of X-linked mental retardation disorders. *Am J Hum Genet*. 2002;71:168–173.
14. Struewing JP, Pineda MA, Sherman ME, et al. Skewed X chromosome inactivation and early-onset breast cancer. *J Med Genet*. 2006;43:48–53.
15. Amos-Landgraf JM, Cottle A, Plenge RM, et al. X chromosome-inactivation patterns of 1,005 phenotypically unaffected females. *Am J Hum Genet*. 2006;79:493–499.
16. Belmont JW. Genetic control of X inactivation and processes leading to X-inactivation skewing. *Am J Hum Genet*. 1996;58:1101–1108.
17. Abkowitz JL, Taboada M, Shelton GH, Catlin SN, Guttorp P, Kiklevich JV. An X chromosome gene regulates hematopoietic stem cell kinetics. *Proc Natl Acad Sci U S A*. 1998;95:3862–3866.
18. Naumova AK, Olien L, Bird LM, et al. Genetic mapping of X-linked loci involved in skewing of X chromosome inactivation in the human. *Eur J Hum Genet*. 1998;6:552–562.
19. Renault NK, Pritchett SM, Howell RE, et al. Human X-chromosome inactivation pattern distributions fit a model of genetically influenced choice better than models of completely random choice. *Eur J Hum Genet*. 2013;21:1396–1402.
20. Chabchoub G, Uz E, Maalej A, et al. Analysis of skewed X-chromosome inactivation in females with rheumatoid arthritis and autoimmune thyroid diseases. *Arthritis Res Ther*. 2009;11:R106.
21. Sharp A, Robinson D, Jacobs P. Age- and tissue-specific variation of X chromosome inactivation ratios in normal women. *Hum Genet*. 2000;107:343–349.
22. Busque L, Mio R, Mattioli J, et al. Nonrandom X-inactivation patterns in normal females: lyonization ratios vary with age. *Blood*. 1996;88:59–65.
23. Champion KM, Gilbert JG, Asimakopoulos FA, Hinshelwood S, Green AR. Clonal haemopoiesis in normal elderly women: implications for the myeloproliferative disorders and myelodysplastic syndromes. *Br J Haematol*. 1997;97:920–926.
24. Gale RE, Fielding AK, Harrison CN, Linch DC. Acquired skewing of X-chromosome inactivation patterns in myeloid cells of the elderly suggests stochastic clonal loss with age. *Br J Haematol*. 1997;98:512–519.
25. Hatakeyama C, Anderson CL, Beever CL, Penaherrera MS, Brown CJ, Robinson WP. The dynamics of X-inactivation skewing as women age. *Clin Genet*. 2004;66:327–332.
26. Tonon L, Bergamaschi G, Dellavecchia C, et al. Unbalanced X-chromosome inactivation in haemopoietic cells from normal women. *Br J Haematol*. 1998;102: 996–1003.

27. Buller RE, Sood AK, Lallas T, Buekers T, Skilling JS. Association between nonrandom X-chromosome inactivation and BRCA1 mutation in germline DNA of patients with ovarian cancer. *J Natl Cancer Inst*. 1999;91:339–346.

28. Kristiansen M, Langerod A, Knudsen GP, Weber BL, Borresen-Dale AL, Orstavik KH. High frequency of skewed X inactivation in young breast cancer patients. *J Med Genet*. 2002;39:30–33.

29. Talebizadeh Z, Bittel DC, Veatch OJ, Kibiryeva N, Butler MG. Brief report: non-random X chromosome inactivation in females with autism. *J Autism Dev Disord*. 2005;35:675–681.

30. Carrel L, Willard HF. X-inactivation profile reveals extensive variability in X-linked gene expression in females. *Nature*. 2005;434:400–404.

31. Brown CJ, Carrel L, Willard HF. Expression of genes from the human active and inactive X chromosomes. *Am J Hum Genet*. 1997;60:1333–1343.

32. Carrel L, Park C, Tyekucheva S, Dunn J, Chiaromonte F, Makova KD. Genomic environment predicts expression patterns on the human inactive X chromosome. *PLoS Genet*. 2006;2:e151.

33. Miller AP, Willard HF. Chromosomal basis of X chromosome inactivation: identification of a multigene domain in Xp11.21-p11.22 that escapes X inactivation. *Proc Natl Acad Sci U S A*. 1998;95:8709–8714.

34. Purcell S, Neale B, Todd-Brown K, et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet*. 2007;81:559–575.

35. Clayton D. Testing for association on the X chromosome. *Biostatistics*. 2008;9:593–600.

36. Clayton D. *snpStats: SnpMatrix and XSnpMatrix classes and methods. R package version 1.2.1*; Vienna, Austria: R Foundation for Statistical Computing; 2011.

37. Nussbaum RL, McInnes RR, Willard HF. *Genetics in Medicine*. Philadelphia, PA: Saunders; 2007.

38. John P. Hussman Institute for Human Genomics. Genetics Basics—lesson 3: modes of Inheritance. http://hihg.med.miami.edu/code/http/modules/education/Design/Print.asp?CourseNum=1&;LessonNum=3. Published 2017.

39. *Understanding Genetics: A District of Columbia Guide for Patients and Health Professionals*. Washington, DC: Genetic Alliance; 2010.

40. Van Regemorter N, Smith C. The importance of determining the mode of inheritance for the estimation of recurrence risks. *J Genet Hum*. 1976;24:49–60.

41. Risch N. Linkage strategies for genetically complex traits. I. Multilocus models. *Am J Hum Genet*. 1990;46:222–228.

42. Patterns of Inheritance. www.bogari.net/Bogari/Principle_files/2-1%20Patterens%20of%20Inheritance.pdf. Published 2017.

43. Cox DR. Tests of separate families of hypotheses. *Proc Fourth Berkeley Symp Math Statist Prob*. 1961;1:105–123.

44. Cox DR. A return to an old paper: "Tests of separate families of hypotheses." *J R Stat Soc B*. 2013;75:207–215.

45. Davidson R, Mackinnon JG. Several tests for model-specification in the presence of alternative hypotheses. *Econometrica*. 1981;49:781–793.

46. Pesaran MH, Dupleich Ulloa MR. Non-nested hypotheses. In: Durlauf SN, Blume LE, eds. *The New Palgrave Dictionary of Economics*. Basingstoke, UK: Palgrave MacMillan; 2008;6:4619–4626.

47. Clarke KA. Testing nonnested models of international relations: reevaluating realism. *Am J Polit Sci*. 2001;45:724–744.

48. Greene WH. *Econometric Analysis*. Upper Saddle River, NJ: Prentice Hall; 2003.

49. Zeileis A, Hothorn T. Diagnostic checking in regression relationships. *R News*. 2002;2:7–10.

50. R Development Core Team. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing; 2014.

51. Montana G. HapSim: a simulation tool for generating haplotype data with pre-specified allele frequencies and LD coefficients. *Bioinformatics*. 2005;21:4309–4311.

52. Talluri R, Shete S. A linkage disequilibrium-based approach to selecting disease-associated rare variants. *PLoS ONE* 2013;8:e69226.

53. Frayling TM, Timpson NJ, Weedon MN, et al. A common variant in the FTO gene is associated with body mass index and predisposes to childhood and adult obesity. *Science*. 2007;16:889–894.

54. Winkelmann J, Schormair B, Lichtner P, et al. Genome-wide association study of restless legs syndrome identifies common variants in three genomic regions. *Nat Genet*. 2007;39:1000–1006.

55. Amos CI, Wang LE, Lee JE, et al. Genome-wide association study identifies novel loci predisposing to cutaneous melanoma. *Hum Mol Gen*. 2011;20:5012–5023.

56. Mailman MD, Feolo M, Jin Y, et al. The NCBI dbGaP database of genotypes and phenotypes. *Nat Genet*. 2007;39:1181–1186.

57. Rosenthal R. Combining results of independent studies. *Psychol Bull*. 1978;85:185–193.

58. Sutton AJ, Abrams KR, Jones DR, Sheldon TA, Song F. *Methods for Meta-Analysis in Medical Research*. Chichester, UK: Wiley; 2000.

## Appendix 1

When comparing 2 XCI models, we follow the Cox test presented by Greene.[48] The details of the test are described in the 5 steps shown below[48] using the comparison of $M_{(2)}$ with $M_{(1)}$ as an example; that is, to test the hypothesis that $X_{(2)}$ is the set of correct regressors (ie, $H_0 : Logit(P(Y=1 \mid X)) = X_{(2)}\beta$). Other tests for the comparisons of $M_{(3)}$ or $M_{(4)}$ with $M_{(1)}$ can be conducted similarly following these same steps. In the calculation, $X_{(1)}$ and $X_{(2)}$ are defined as matrices with 3 columns: the first column of 1's represents the intercept and the rest of the columns represent the SNP based on different XCI models and the individual's sex. The test is conducted using the R function "coxtest" in the "lmtest" package[49,50]:

*Step 1.* Regress $Y$ on $X_{(2)}$ to obtain the estimations of the coefficients $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_s)'$. The fitted values can be calculated as $\hat{Y}_{X_{(2)}} = X_{(2)}\hat{\beta}$, and the residuals can be calculated as $\hat{r}_{X_{(2)}} = Y - X_{(2)}\hat{\beta}$. Thus, $s_{X_{(2)}}^2 = \hat{r}'_{X_{(2)}}\hat{r}_{X_{(2)}} / N$.

*Step 2.* Regress $Y$ on $X_{(1)}$ to obtain the estimations of the coefficients $\hat{\alpha} = (\hat{\alpha}_0, \hat{\alpha}_1, \hat{\alpha}_s)'$. The fitted values can be calculated as $\hat{Y}_{X_{(1)}} = X_{(1)}\hat{\alpha}$, and the residuals can be calculated as $\hat{r}_{X_{(1)}} = Y - X_{(1)}\hat{\alpha}$. Thus, $s_{X_{(1)}}^2 = \hat{r}'_{X_{(1)}}\hat{r}_{X_{(1)}} / N$.

*Step 3.* Regress $\hat{Y}_{X_{(2)}}$ on $X_{(1)}$ to obtain the estimations of the coefficients $\hat{\beta}_{X_{(2)}} = (\hat{\beta}_{X_{(2)}0}, \hat{\beta}_{X_{(2)}1}, \hat{\beta}_{X_{(2)}s})'$. The residuals can be calculated as $\hat{r}_{X_{(1)}X_{(2)}} = \hat{Y}_{X_{(2)}} - X_{(1)}\hat{\beta}_{X_{(2)}} = Q_{X_{(1)}}X_{(2)}\hat{\beta}$, where $Q_{X_{(1)}} = I - X_{(1)}(X_1'X_{(1)})^{-1}X_1'$. Thus, $\hat{r}'_{X_{(1)}X_{(2)}}\hat{r}_{X_{(1)}X_{(2)}} = \hat{\beta}'X_{(2)}'Q_{X_{(1)}}X_{(2)}\hat{\beta}$.

*Step 4.* Regress $\hat{r}_{X_{(1)}X_{(2)}}$ on $X_{(2)}$ to obtain the estimations of the coefficients $\hat{\alpha}_{X_{(1)}} = (\hat{\alpha}_{X_{(1)}0}, \hat{\alpha}_{X_{(1)}1}, \hat{\alpha}_{X_{(1)}s})'$. The residuals can be calculated as $\hat{r}_{X_{(2)}X_{(1)}X_{(2)}} = Q_{X_{(2)}}Q_{X_{(1)}}X_{(2)}\hat{\beta}$, where $Q_{X_{(2)}} = I - X_{(2)}(X_2'X_{(2)})^{-1}X_2'$. Thus, $\hat{r}'_{X_{(2)}X_{(1)}X_{(2)}}\hat{r}_{X_{(2)}X_{(1)}X_{(2)}} = \hat{\beta}'X_{(2)}'Q_{X_{(1)}}Q_{X_{(2)}}Q_{X_{(1)}}X_{(2)}\hat{\beta}$.

*Step 5.* Compute $s_{X_{(1)}X_{(2)}}^2 = s_{X_{(2)}}^2 + \hat{r}'_{X_{(1)}X_{(2)}}\hat{r}_{X_{(1)}X_{(2)}}$. Compute $T = N\ln(s_{X_{(1)}}^2 / s_{X_{(1)}X_{(2)}}^2)/2$ and $\hat{V}(T) = (s_{X_{(2)}}^2 \hat{r}'_{X_{(2)}X_{(1)}X_{(2)}}\hat{r}_{X_{(2)}X_{(1)}X_{(2)}} / s_{X_{(1)}X_{(2)}}^4)$. Compare the value of $T/\sqrt{\hat{V}(T)}$ to the critical value of the standard normal distribution to assess the significance (2-sided test).