

UCSF

UC San Francisco Previously Published Works

Title

Robust, flexible, and scalable tests for Hardy-Weinberg Equilibrium across diverse ancestries

Permalink

<https://escholarship.org/uc/item/7607p308>

Journal

Genetics, 218(1)

ISSN

0016-6731

Authors

Kwong, Alan M
Blackwell, Thomas W
LeFaive, Jonathon
et al.

Publication Date

2021-05-17

DOI

10.1093/genetics/iyab044

Peer reviewed

Robust, flexible, and scalable tests for Hardy–Weinberg equilibrium across diverse ancestries

Alan M. Kwong,¹ Thomas W. Blackwell,¹ Jonathon LeFaive,¹ Mariza de Andrade,² John Barnard,³ Kathleen C. Barnes,⁴ John Blangero,⁵ Eric Boerwinkle,^{6,7} Esteban G. Burchard,^{8,9} Brian E. Cade,^{10,11} Daniel I. Chasman,¹² Han Chen,^{6,13} Matthew P. Conomos,¹⁴ L. Adrienne Cupples,^{15,16} Patrick T. Ellinor,^{17,18} Celeste Eng,⁹ Yan Gao,¹⁹ Xiuqing Guo,²⁰ Marguerite Ryan Irvin,²¹ Tanika N. Kelly,²² Wonji Kim,²³ Charles Kooperberg,²⁴ Steven A. Lubitz,^{17,18} Angel C. Y. Mak,⁹ Ani W. Manichaikul,²⁵ Rasika A. Mathias,²⁶ May E. Montasser,²⁷ Courtney G. Montgomery,²⁸ Solomon Musani,²⁹ Nicholette D. Palmer,³⁰ Gina M. Peloso,¹⁵ Dandi Qiao,²³ Alexander P. Reiner,²⁴ Dan M. Roden,³¹ M. Benjamin Shoemaker,³² Jennifer A. Smith,³³ Nicholas L. Smith,^{34,35,36} Jessica Lasky Su,²³ Hemant K. Tiwari,³⁷ Daniel E. Weeks,³⁸ Scott T. Weiss,²³ NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium, TOPMed Analysis Working Group, Laura J. Scott,¹ Albert V. Smith,¹ Gonçalo R. Abecasis,¹ Michael Boehnke,¹ and Hyun Min Kang^{1,*}

¹Department of Biostatistics, Center for Statistical Genetics, University of Michigan, Ann Arbor, MI 48109, USA

²Mayo Clinic, Rochester, MN 55905, USA

³Department of Quantitative Health Sciences, Lerner Research Institute, Cleveland Clinic, Cleveland, OH 44106, USA

⁴Department of Medicine, Anschutz Medical Campus, University of Colorado, Aurora, CO 80045, USA

⁵Department of Human Genetics, South Texas Diabetes and Obesity Institute, University of Texas Rio Grande Valley School of Medicine, Brownsville, TX 78520, USA

⁶Department of Epidemiology, Human Genetics Center, Human Genetics and Environmental Sciences, School of Public Health, The University of Texas Health Science Center at Houston, Houston, TX 77030, USA

⁷Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX 77030, USA

⁸Department of Bioengineering and Therapeutic Sciences, University of California San Francisco, San Francisco, CA 94143, USA

⁹Department of Medicine, University of California San Francisco, San Francisco, CA 94143, USA

¹⁰Division of Sleep and Circadian Disorders, Brigham and Women's Hospital, Boston, MA 02115, USA

¹¹Division of Sleep Medicine, Harvard Medical School, Boston, MA 02115, USA

¹²Division of Preventive Medicine, Brigham and Women's Hospital, Boston, MA 02215, USA

¹³Center for Precision Health, School of Public Health and School of Biomedical Informatics, The University of Texas Health Science Center at Houston, Houston, TX 77030, USA

¹⁴Department of Biostatistics, University of Washington, Seattle, WA 98195, USA

¹⁵Department of Biostatistics, Boston University School of Public Health, Boston, MA 02118, USA

¹⁶Framingham Heart Study, Framingham, MA 01702, USA

¹⁷Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA 02114, USA

¹⁸Cardiovascular Disease Initiative, The Broad Institute of MIT and Harvard, Cambridge, MA 02124, USA

¹⁹Department of Physiology and Biophysics, University of Mississippi Medical Center, Jackson, MS 39216 USA

²⁰Department of Pediatrics, The Institute for Translational Genomics and Population Sciences, The Lundquist Institute at Harbor-UCLA Medical Center, Torrance, CA 90502, USA

²¹Department of Epidemiology, School of Public Health, University of Alabama at Birmingham, Birmingham, AL 35294, USA

²²Department of Epidemiology, Tulane University, New Orleans, LA 70112, USA

²³Channing Division of Network Medicine, Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA 02115, USA

²⁴Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA

²⁵Department of Public Health Sciences, Center for Public Health Genomics, University of Virginia, Charlottesville, VA 22908, USA

²⁶GeneSTAR Research Program and Division of Allergy and Clinical Immunology, Department of Medicine, Johns Hopkins University, Baltimore, MD 21205, USA

²⁷Division of Endocrinology, Diabetes and Nutrition, Department of Medicine, University of Maryland School of Medicine, Baltimore, MD 21201, USA

²⁸Sarcoidosis Research Unit, Genes and Human Disease Research Program, Oklahoma Medical Research Foundation, Oklahoma City, OK 73104, USA

²⁹Jackson Heart Study, University of Mississippi Medical Center, Jackson, MS 39216, USA

³⁰Department of Biochemistry, Wake Forest School of Medicine, Winston-Salem, NC 27157, USA

³¹Departments of Medicine, Pharmacology, and Biomedical Informatics, Vanderbilt University Medical Center, Nashville, TN 37232, USA

³²Department of Medicine, Vanderbilt University Medical Center, Nashville, TN 37232, USA

³³Department of Epidemiology, School of Public Health, University of Michigan, Ann Arbor, MI 48109, USA

³⁴Department of Epidemiology, University of Washington, Seattle, WA 98195, USA

³⁵Kaiser Permanente Washington Health Research Institute, Kaiser Permanente Washington, Seattle, WA 98101, USA

³⁶Department of Veterans Affairs, Seattle Epidemiologic Research and Information Center, Office of Research and Development, Seattle, WA 98108, USA

³⁷Department of Biostatistics, School of Public Health, University of Alabama at Birmingham, Birmingham, AL 35294, USA

³⁸Departments of Human Genetics and Biostatistics, Graduate School of Public Health, University of Pittsburgh, Pittsburgh, PA 15261, USA

*Corresponding author: Department of Biostatistics, Center for Statistical Genetics, University of Michigan, Ann Arbor, MI 48109, USA. hmkang@umich.edu

Abstract

Traditional Hardy–Weinberg equilibrium (HWE) tests (the χ^2 test and the exact test) have long been used as a metric for evaluating genotype quality, as technical artifacts leading to incorrect genotype calls often can be identified as deviations from HWE. However, in data sets composed of individuals from diverse ancestries, HWE can be violated even without genotyping error, complicating the use of HWE testing to assess genotype data quality. In this manuscript, we present the Robust Unified Test for HWE (RUTH) to test for HWE while accounting for population structure and genotype uncertainty, and to evaluate the impact of population heterogeneity and genotype uncertainty on the standard HWE tests and alternative methods using simulated and real sequence data sets. Our results demonstrate that ignoring population structure or genotype uncertainty in HWE tests can inflate false-positive rates by many orders of magnitude. Our evaluations demonstrate different tradeoffs between false positives and statistical power across the methods, with RUTH consistently among the best across all evaluations. RUTH is implemented as a practical and scalable software tool to rapidly perform HWE tests across millions of markers and hundreds of thousands of individuals while supporting standard VCF/BCF formats. RUTH is publicly available at <https://www.github.com/statgen/ruth>.

Keywords: population structure; principal components analysis; next-generation sequencing; genotype likelihoods

Introduction

Hardy–Weinberg equilibrium (HWE) is a fundamental theorem of population genetics and has been one of the key mathematical principles to understand the characteristics of genetic variation in a population for more than a century (Hardy 1908; Weinberg 1908). Genetic variants in a homogeneous population typically follow HWE except for unusual deviations owing to very strong case–control association and enrichment (Nielsen et al. 1998), sex linkage, or nonrandom sampling (Waples 2015).

HWE tests are often used to assess the quality of microsatellite (Van Oosterhout et al. 2004), SNP-array (Wigginton et al. 2005), and sequence-based (Danecek et al. 2011) genotypes (GTs). Testing for HWE may reveal technical artifacts in sequence or GT data, such as high rates of genotyping error and/or missingness, or sequencing/alignment errors (Nielsen et al. 2011). It can also identify hemizygotes in structural variants which are incorrectly called as homozygotes (McCarroll et al. 2006). Quality control for array- or sequence-based GTs typically includes a HWE test to detect and filter out artifactual or poorly genotyped variants (Laurie et al. 2010; Nielsen et al. 2011).

Although HWE tests are commonly and reliably used for variant quality control in samples from homogeneous populations, applying them to more diverse samples remains challenging. When analyzing individuals from a heterogeneous population, the standard HWE tests may falsely flag real, well-genotyped variants, unnecessarily filtering them out for downstream analyses (Hao and Storey 2019). This problem is important since genetic studies increasingly collect genetic data from heterogeneous populations. In principle, HWE tests in these structured populations can be performed on smaller cohorts with homogenous backgrounds (Bycroft et al. 2018), and the test statistics combined using Fisher’s or Stouffer’s method (Mosteller and Fisher 1948; Stouffer 1949). However, such a procedure requires much more effort than using a single HWE test across all samples. In addition, this approach cannot account for any heterogeneity within each of the smaller cohorts.

Here, we describe RUTH (Robust Unified Test for Hardy–Weinberg Equilibrium) which tests for HWE under heterogeneous population structure. Our primary motivation for developing RUTH is to robustly filter out artifactual or poorly genotyped variants using HWE test statistics. RUTH is (1) computationally efficient, (2) robust against various degrees of population structure, and (3) flexible in accepting key representations of sequence-based GTs including best-guess GTs and genotype likelihoods (GLs). We perform systematic evaluations of RUTH and

alternative methods for HWE testing using simulated and real data to explore the advantages and disadvantages of these methods for samples of diverse ancestries.

Materials and methods

Unadjusted HWE tests

Consider a study of n participants with true (unobserved) GTs g_1, g_2, \dots, g_n at a bi-allelic variant coded as 0 (reference homozygote), 1 (heterozygote), or 2 (alternate homozygote). Represent the best-guess/hard-call (observed) GTs as $\hat{g}_1, \hat{g}_2, \dots, \hat{g}_n$. A simple HWE test uses the χ^2 statistic to compare the expected and observed GT counts, assuming no population structure and no GT uncertainty. The χ^2 HWE test statistic is defined as $T_{\chi^2} = \sum_{j=0}^2 (c_j - \hat{c}_j)^2 / \hat{c}_j$ where $c_j = \sum_{i=0}^n I(\hat{g}_i = j)$ (ignoring missing GTs), $\hat{p} = (c_1 + 2c_2) / 2n$, $\hat{q} = 1 - \hat{p}$, $\hat{c}_0 = n\hat{q}^2$, $\hat{c}_1 = 2n\hat{p}\hat{q}$, and $\hat{c}_2 = n\hat{p}^2$. Under HWE, the asymptotic distribution of T_{χ^2} is assumed to follow χ_1^2 (Rohlf and Weir 2008). An exact test is known to be more accurate for finite samples, particularly for rare variants (Wigginton et al. 2005), and using mid-P-values instead of exact P-values will lead to slightly less conservative estimates (Graffelman and Moreno 2013). HWE tests stratified by case–control status are known to prevent an inflation of Type I errors for disease-associated variants (Li and Li 2008). Widely used software tools such as PLINK (Purcell et al. 2007) and VCFTools (Danecek et al. 2011) implement an exact HWE test based on best-guess GTs. We will refer to the exact test as the unadjusted test.

Existing HWE tests accounting for structured populations

The unadjusted HWE test assumes a homogeneous population. If a study is composed of a set of discrete structured subpopulations, a straightforward extension of the unadjusted test is to (1) stratify each study participant into exactly one of the subpopulations, (2) perform the unadjusted HWE test for each subpopulation separately, and (3) meta-analyze test statistics across subpopulations to obtain a combined P-value using Stouffer’s method (Stouffer et al. 1949). More specifically, let z_1, z_2, \dots, z_s be the z-scores from HWE test statistics for s distinct subpopulations with sample sizes n_1, n_2, \dots, n_s . A combined meta-analysis HWE test statistic across the subpopulations is $T_{meta} = (\sum_{i=1}^s z_i / \sqrt{n_i}) / \sqrt{\sum_{i=1}^s 1/n_i}$, which asymptotically follows a standard normal distribution when each subpopulation follows HWE.

When the population cannot be easily stratified into distinct subpopulations (e.g., intra-continental diversity or an admixed population), a quantitative representation of genetic ancestry,

such as principal component (PC) co-ordinates or fractional mixture over subpopulations, can be more useful for representing genetic diversity (Rosenberg et al. 2002; Price et al. 2006). HWES takes PCs as additional input to perform HWE tests under population structure with logistic regression (Sha and Zhang 2011), and a similar idea was suggested by Hao et al. (2016). However, existing implementations do not support sequence-based GTs (where GT uncertainty may remain at low or moderate sequencing depth) or other commonly used formats for genetic array data. A recent method PCAngsd estimates PCs from uncertain GTs represented as GLs (Meisner and Albrechtsen 2019) and uses these estimates to perform a likelihood ratio test (LRT) for HWE, similar to the LRT version of RUTH with differences in computational performance (see below).

Robust HWE testing with RUTH

Here we describe RUTH (Robust and Unified Test for Hardy–Weinberg equilibrium) to enable HWE testing under structured populations, which is especially useful for large sequencing studies. We developed RUTH to produce HWE test statistics to allow quality control of sequence-based variant callsets from increasingly diverse samples. RUTH models the uncertainty encoded in sequence-based GTs to robustly distinguish true and artifactual variants in the presence of population structure, and seamlessly scales to millions of individuals and genetic variants.

We assume the observed GT for individual i can be represented as a GL $L_i^{(G)} = \Pr(\text{Data}_i | g_i = G)$, where Data_i represents observed data (e.g., sequence or array), and $g_i \in \{0, 1, 2\}$ the true (unobserved) GT. For example, GLs for sequence-based GTs can be represented as $L_i^{(G)} = \prod_{j=1}^{d_i} \Pr(r_{ij} | g_i = G; q_{ij})$ where d_i is the sequencing depth, r_{ij} is the observed read, and q_{ij} is the corresponding quality score (Ewing and Green 1998; Jun et al. 2012). We model GLs for best-guess GTs \hat{g}_i from SNP arrays as $L_i^{(G)} = (1 - e_i)^2, 2e_i(1 - e_i), e_i^2$ for $\hat{g}_i = 2, 1, 0$ where e_i is the assumed per-allele error rate. Imputed GTs may also be approximately modeled using this framework, but the current implementation requires creating a pseudo-GL to describe this uncertainty (see Discussion section).

Accounting for population structure with individual-specific allele frequencies

We account for population structure by modeling individual-specific allele frequencies from quantitative coordinates of genetic ancestry such as PCs, similar to Hao et al. (2016). For any given variant, instead of assuming that GTs follow HWE with a single universal allele frequency across all individuals, we assume that GTs follow HWE with heterogeneous allele frequencies specific to each individual, modeled as a function of genetic ancestry. Let $\mathbf{x}_i \in \mathbb{R}^k$ represent the genetic ancestry of individual i , where k is the number of PCs used. We estimate individual-specific allele frequency p as a bounded linear function of genetic ancestry

$$p(\mathbf{x}_i; \boldsymbol{\beta}) = \begin{cases} \boldsymbol{\beta}^T \mathbf{x}_i & \varepsilon \leq \boldsymbol{\beta}^T \mathbf{x}_i \leq 1 - \varepsilon \\ \varepsilon & \boldsymbol{\beta}^T \mathbf{x}_i < \varepsilon \\ 1 - \varepsilon & \boldsymbol{\beta}^T \mathbf{x}_i > 1 - \varepsilon, \end{cases}$$

where ε is the minimum frequency threshold. We estimate $\hat{\boldsymbol{\beta}}$ with an Expectation–Maximization (E-M) algorithm. We used $\varepsilon = 1/4n$ in our evaluation. Although we used a linear model for $p(\mathbf{x}_i; \boldsymbol{\beta})$ for computational efficiency, it is straightforward to apply

a logistic model, which is arguably better (Yang et al. 2012; Hao et al. 2016).

Let $p_i = p(\mathbf{x}_i; \boldsymbol{\beta})$ and $q_i = 1 - p_i$ be the individual specific allele frequencies of the nonreference and reference alleles for individual i . Under the null hypothesis of HWE, the frequencies of GTs (0, 1, 2) are $[q_i^2, 2p_i q_i, p_i^2]$. Under the alternative hypothesis, we assume that these frequencies are $[q_i^2 + \theta p_i q_i, 2p_i q_i(1 - \theta), p_i^2 + \theta p_i q_i]$ where θ is the inbreeding coefficient. This model is a straightforward extension of a fully general model where p_i, q_i is identical across all samples. Then the log-likelihood across all study participants is

$$l(\boldsymbol{\beta}, \theta) = \sum_{i=1}^n \log [L_i^{(0)}(q_i^2 + \theta p_i q_i) + L_i^{(1)} 2p_i q_i(1 - \theta) + L_i^{(2)}(p_i^2 + \theta p_i q_i)]$$

Under both the null ($\theta = 0$) and alternative ($\theta \neq 0$) hypotheses, we maximize the log-likelihood using an E-M algorithm (Dempster et al. 1977). As we empirically observed quick convergence within several iterations in most cases, we used a fixed ($n=20$) number of iterations in our implementation (Supplementary Figure S2).

RUTH score test

The score function of the log-likelihood is the derivative of the log-likelihood with respect to θ :

$$\begin{aligned} U(\theta) &= \sum_{i=1}^n \frac{p_i q_i [L_i^{(0)} - 2L_i^{(1)} + L_i^{(2)}]}{L_i^{(0)}(q_i^2 + \theta p_i q_i) + L_i^{(1)} 2p_i q_i(1 - \theta) + L_i^{(2)}(p_i^2 + \theta p_i q_i)} \\ &= \sum_{i=1}^n u_i(\theta) \end{aligned}$$

Since $u_i'(\theta) = -u_i^2(\theta)$, we construct a score test statistic of $H_0 : \theta = 0$ vs $H_1 : \theta \neq 0$ as:

$$T_{\text{score}} = \frac{[U(0)]^2}{I(0)} = \frac{[\sum_{i=1}^n u_i(0)]^2}{\sum_{i=1}^n u_i^2(0)}$$

where $I(0)$ is the Fisher information under the null hypothesis. Under the null, T_{score} has an asymptotic χ^2 distribution with one degree of freedom, i.e., $T_{\text{score}} \sim \chi_1^2$. A detailed algorithm is shown in Supplementary Figure S1.

RUTH likelihood ratio test

The log-likelihood function $l(\boldsymbol{\beta}, \theta)$ can also be used to calculate an LRT statistic:

$$T_{\text{LRT}} = 2[\max_{\boldsymbol{\beta}, \theta} l(\boldsymbol{\beta}, \theta) - \max_{\boldsymbol{\beta}} l(\boldsymbol{\beta}, 0)].$$

Like the score test, we estimate MLE parameters $\boldsymbol{\beta}, \theta$ iteratively using an E-M algorithm to test $H_0 : \theta = 0$ vs $H_1 : \theta \neq 0$. Under the null hypothesis, the asymptotic distribution of T_{LRT} is expected to follow χ_1^2 . This test is very similar to the likelihood-ratio test proposed by PCAngsd (Meisner and Albrechtsen 2019), except PCAngsd does not re-estimate $\boldsymbol{\beta}$ under the alternative hypothesis. In principle, the RUTH LRT should be slightly more powerful owing to this difference; we expect the practical difference in power to be small, as deviations from HWE usually do not change the estimates of $\boldsymbol{\beta}$ substantially.

Simulation of GTs and sequence reads under population structure

We simulated sequence-based GTs under population structure using the following procedure. First, for each variant, we simulated an ancestral allele frequency and population-specific allele frequencies. Second, we sampled unobserved (true) GTs based on these allele frequencies. Third, we sampled sequence reads based on the unobserved GTs. Fourth, we generated GLs and best-guess GTs based on sequence reads. Our goal was to simulate variants such that each subpopulation will have different average allele frequencies from other subpopulations.

To simulate ancestral and population-specific allele frequencies, we followed the procedure of [Balding and Nichols \(1995\)](#), except we sampled ancestral allele frequencies from $p \sim \text{Uniform}(0, 1)$ instead of $p \sim \text{Uniform}(0.1, 0.9)$ to include rare variants. For each of $K \in \{1, 2, 5, 10\}$ populations, we sampled population-specific allele frequencies from $p_k \sim \text{Beta}\left(\frac{p(1-F_{st})}{F_{st}}, \frac{(1-p)(1-F_{st})}{F_{st}}\right)$, where $k \in \{1, \dots, K\}$, and $F_{st} \in \{0.01, 0.02, 0.03, 0.05, 0.10\}$ was the fixation index to quantify the differentiation between populations, as suggested by [Holsinger \(2004\)](#) and implemented in the previously published studies ([Holsinger et al. 2002](#); [Balding 2003](#)). Because p_k no longer follows the uniform distribution, we used rejection sampling to ensure that $-p = \frac{1}{K} \sum_{k=1}^K p_k$ is uniformly distributed across 100 bins across simulations to avoid artifacts caused by systematic differences in allele frequencies.

The unobserved GT $G_i \in \{0, 1, 2\}$ for individual $i \in \{1, \dots, n_k\}$, belonging to population k with sample size n_k , was simulated from GT frequencies $(q_k^2 + \theta p_k q_k, 2p_k q_k(1 - \theta), p_k^2 + \theta p_k q_k)$, where $q_k = 1 - p_k$ and $\theta \in \left[-\min\left(\frac{q_k}{p_k}, \frac{p_k}{q_k}\right), 1\right]$ quantifies deviation from HWE; $\theta = 0$ represents HWE, whereas $\theta < 0$ and $\theta > 0$ represent excess heterozygosity and homozygosity compared to HWE expectation, respectively. In our experiments, we evaluated $\theta \in \{0, \pm 0.01, \pm 0.05, \pm 0.1, \pm 0.5\}$. When θ was smaller than the minimum possible value for a specific population, we replaced it with the minimum value.

We simulated sequence reads based on unobserved GTs, sequence depths, and base-call error rates. To reflect the variation of sequence depths between individuals, we simulated the mean depth of each sequenced sample as $\mu_i \sim \text{Uniform}(1, 2D - 1)$, where D is the expected depth and $D = 5$ and $D = 30$, representing low-coverage and deep sequencing, respectively. For each sequenced sample and variant site, we sampled the sequence depth from $d_i \sim \text{Poisson}(\mu_i)$. Each sequence read carried either of the possible unobserved (true) alleles $r_{ij} \in \{0, 1\}$, where $j \in \{1, \dots, d_i\}$. Given unobserved GT G_i , we generated $r_{ij} \sim \text{Bernoulli}\left(\frac{G_i}{2}\right)$, with observed allele $o_{ij} = (1 - e_{ij})r_{ij} + e_{ij}(1 - r_{ij})$ flipping to the other allele when a sequencing error occurs with probability $e_{ij} \sim \text{Bernoulli}(\epsilon)$. We used $\epsilon = 0.01$ throughout our simulations (which corresponds to phred-scale base quality of 20) and assumed that all base-calling errors switched between reference and alternate alleles.

We then generated GLs and best-guess GTs from the simulated alleles. Let $t_i = \sum_{j=1}^{d_i} o_{ij}$ be the observed alternate allele count. The GLs for the three possible GTs are $L_i^{(0)} = (1 - \epsilon)^{d_i - t_i} (\epsilon)^{t_i}$, $L_i^{(1)} = 0.5^{d_i}$, $L_i^{(2)} = (\epsilon)^{d_i - t_i} (1 - \epsilon)^{t_i}$. We called best-guess GTs by using the overall ancestral allele frequency $-p$ for a given variant as the prior, then calling the GT corresponding to the highest posterior probability among $(L_i^{(0)}(1 - p)^2, 2L_i^{(1)}-p(1 - p)^2, L_i^{(2)}-p^2)$ for each individual. For each possible combination of F_{st} , K , and θ , we generated

50,000 independent variants across a set of $n = 5000$ samples with per-ancestry sample sizes $n_k = n/K$.

Evaluation of Type I error and statistical power

We used different P -value thresholds, F_{st} values, number of ancestry groups K , and average sequencing depth D to determine the number of variants significantly deviating from HWE. To evaluate Type I error, we simulated sequence reads under HWE ($\theta = 0$) and calculated the proportion of significant variants at each P -value threshold. In RUTH tests, we assumed that PCs were accurately estimated using true GTs unless indicated otherwise. For real data, we summarized ancestral information by projecting PCs estimated from full genomes onto the reference PC space of the Human Genome Diversity Project (HGDP) panel ([Li et al. 2008](#)) using `verifyBamID2` ([Zhang et al. 2020](#)), similar to the procedure for variant calling in the TOPMed Project, which has integrated RUTH as part of its quality control pipeline (https://github.com/statgen/topmed_variant_calling).

In all data sets, we evaluated the tradeoff between Type I Error and power for each method using precision-recall curves (PRCs) and receiver-operator characteristic curves (ROCs). In simulated data, we considered variants with $\theta = 0$ to be true negatives and variants with $\theta = -0.05$ to be true positives. For real data, we labeled high-quality (HQ) variants as negative and low quality (LQ) variants as positive.

Data source

To evaluate our method, we used sequence-based GT data from the 1000 Genomes Project (1000G) ([Auton et al. 2015](#)) and the Trans-Omics Precision Medicine (TOPMed) Project ([Taliun et al. 2021](#)). In both cases, we used subsets of variants from chromosome 20. For 1000G, we started with 1,812,841 variants in 2504 individuals, with an average depth of 7.0 \times . For TOPMed, we started with 12,983,576 variants in 53,831 individuals, with an average depth of 37.2 \times .

Application to 1000 genomes data

To test our method on 1000G data, we first needed to define two sets of variants: one set which is expected to follow HWE, and another set which is expected to deviate from HWE. Unlike simulated data, variants in 1000G are not clearly classified into “true” or “artificial,” so evaluation of false positives and power is less straightforward. We focused on two subsets of variants in chromosome 20 which serve as proxies for these two variant types. We selected non-monomorphic sites found in both the Illumina Infinium Omni2.5 genotyping array and in HapMap3 (The International HapMap Consortium et al. 2010) as HQ variants that mostly follow HWE after controlling for ancestry, ending up with 17,740 variants. We selected variants that displayed high discordance between duplicates or Mendelian inconsistencies within family members in TOPMed as LQ variants which should be enriched for deviations from HWE even after accounting for ancestry, ending up with 10,966 variants. Among 329,699 LQ variants from TOPMed in chromosome 20, we found that only 10,966 overlap with 1000G samples. We suspect that a substantial fraction of these 10,966 LQ variants are true variants since they passed all of the 1000G Project’s quality filters. Nevertheless, we still expect a much larger fraction of these LQ variants to deviate from HWE compared to HQ variants.

We evaluated multiple representations of sequence-based GTs from 1000G. As 1000G samples were sequenced at relatively low-coverage of 7.0 \times on average, best-guess GTs inferred only from sequence reads (raw GT) tend to have poor accuracy.

Therefore, the officially released best-guess GTs in 1000G were estimated by combining GLs, calculated based on sequence reads, with haplotype information from nearby variants through linkage-disequilibrium (LD)-aware GT refinement using SHAPEIT2 (Delaneau et al. 2013). This procedure resulted in more accurate GTs (LD-aware GT), but it implicitly assumed HWE during refinement. As different representations of sequence GTs may result in different performance in HWE tests, we evaluated all three representations—raw GT, LD-aware GT, and GL. In all tests of RUTH using hard GT calls, we assumed the error rate for GT-based GTs to be 0.5%, which is representative of a typical non-reference GT error rate for SNP arrays. We restricted our analyses to biallelic variants. The positions and alleles of 1000G and TOPMed variants were matched using the liftOver software tool (Kuhn et al. 2013).

We evaluated all tests as described above. For meta-analysis with Stouffer's method, we divided the samples into 5 strata, using the five 1000G super population code labels—African (AFR), Admixed American (AMR), East Asian (EAS), European (EUR), and South Asian (SAS). To obtain PC co-ordinates for 1000G samples, we estimated 4 PCs from the aligned sequence reads (BAM) with verifyBamID2 (Zhang et al. 2020), using PCs from 936 samples from the Human Genome Diversity Project (HGDP) panel as reference coordinates. The RUTH score test and LRT used these PCs as inputs, along with GTs in raw GT, LD-aware GT, and GL formats. For PCAnsd, we used GLs from all variants tested as the input. We limited the analysis to a single chromosome due to the heavy computational requirements of PCAnsd.

Application to TOPMed data

We analyzed variants from 53,831 individuals from the TOPMed sequencing study (Taliun et al. 2021). These samples came from multiple studies from a diverse spectrum of ancestries, leading to substantial population structure. Using the same criteria as our 1000G analysis, we identified 17,524 HQ variants and 329,699 low-quality variants across chromosome 20. Since TOPMed genomes were deeply sequenced at $37.2 \times (\pm 4.5 \times)$, LD-aware GT refinement was not necessary to obtain accurate GTs. Therefore, we used two GT representations—raw GT and GL—in our evaluations. This GT data contained no missingness.

Similar to 1000G, for best-guess GTs (raw GT), we used PLINK for the unadjusted test. For meta-analysis, we assigned each sample to one of the five 1000G super populations as follows. First, we summarized the genetic ancestries of aligned sequenced genomes with verifyBamID2 by estimating 4 PCs using HGDP as reference. Second, we used Procrustes analysis (Dryden and Mardia 1998; Wang et al. 2010) to align the PC coordinates of HGDP panels (to account for different genome builds) so that the PC coordinates were compatible between TOPMed and 1000G samples. Third, for each TOPMed sample, we identified the 10 closest corresponding individuals from 1000G using the first 4 PC coordinates with a weighted voting system (assigning the closest individual a score of 10, next closest a score of 9, and so on until the 10th closest individual is assigned a score of 1, then adding up the scores for each super population) to determine the super population code that had the highest sum of scores, and therefore best described that sample. In this way, we classified 15,580 samples as AFR, 4836 as AMR, 29,943 as EUR, 2960 as EAS, and 716 as SAS. Among these samples, 94.5% had the same super population code for all 10 nearest 1000G neighbors. To evaluate the RUTH score test and LRT for both raw GT and GL, we used 4 PCs estimated by verifyBamID2 (Zhang et al. 2020), consistent with the method applied for the 1000G data.

Impact of ancestry estimates on adjusted HWE tests

We examined the effect of changing the number of PCs used as input for RUTH tests by using 2 PCs as opposed to 4 PCs. We also evaluated the impact of using different approaches to classify ancestry when adjusting for population structure with meta-analysis. By default, our analysis classified the 1000G subjects into 5 continental super populations based on published information (Auton et al. 2015). For TOPMed, the best-matching 1000G continental ancestry was carefully determined using the PCA-based matching strategy described above. However, in practice, ancestry classification may be performed with a coarser resolution (Jin et al. 2019). To mimic plausible scenarios in which sample ancestries are not carefully determined, we used k-means clustering on the first 2 PCs of our samples to divide individuals into three distinct groups, roughly corresponding to East Asian, European, and African populations, and performed meta-analyses based on this coarse classification for both 1000G and TOPMed data.

Data availability

RUTH is available at <https://github.com/statgen/ruth>. GT data from 1000G are available from the International Genome Sample Resource at <https://www.internationalgenome.org> (last accessed March 22nd, 2021). TOPMed data are available via a dbGaP application for controlled-access data (see <https://www.nhlbiwgs.org> [last accessed March 22nd, 2021] for details). Supplementary materials have been uploaded to figshare: <https://doi.org/10.25386/genetics.14068970>.

Results

Simulation: effect of GT uncertainty

To evaluate the impact of GT uncertainty, we first compared tests in the absence of population structure (i.e., single ancestry). For the unadjusted test, we used only best-guess GTs. For PCAnsd, we used only GLs. For RUTH score and LRTs, we used both.

Using GLs over GTs substantially reduced Type I errors in HWE tests, especially in low-coverage data (Figure 1, A–C). For example, the standard HWE test based on GTs resulted in a 229-fold inflation (22.9%) at $P < 0.001$ (Figure 1B and Supplementary Table S1), a threshold which allows the evaluation of Type I error with reasonable precision with 50,000 variants (50 expected false positives under the null). GT-based RUTH score and RUTH-LRT tests showed similar inflation. When GLs were used instead of best-guess GTs, RUTH score and RUTH-LRT had Type I errors close to the null expectation (0.001 for RUTH score and 0.0012 for RUTH-LRT). PCAnsd, which also accounts for GT uncertainty (Meisner and Albrechtsen 2019), had similar performance. The severely inflated Type I errors with best-guess GTs can largely be attributed to high uncertainty and bias toward homozygote reference GTs in single-site calls from low-coverage sequence data, resulting in apparent deviations from HWE. For high-coverage sequence data, inflation of Type I error with GTs was substantially attenuated; inflation nearly disappeared when using GLs (0.004 for RUTH score and 0.002 for RUTH-LRT; Figure 1, D–F).

Next, we evaluated the power to identify variants truly deviating from HWE at various levels of inbreeding (θ). For low-coverage sequence data, we skip interpretation of power of GT-based tests owing to their extremely inflated false-positive rates. All GL-based tests behaved similarly, achieving ~19–21% power at $P < 0.001$ with moderate excess heterozygosity ($\theta = -0.05$) (Figure 2B

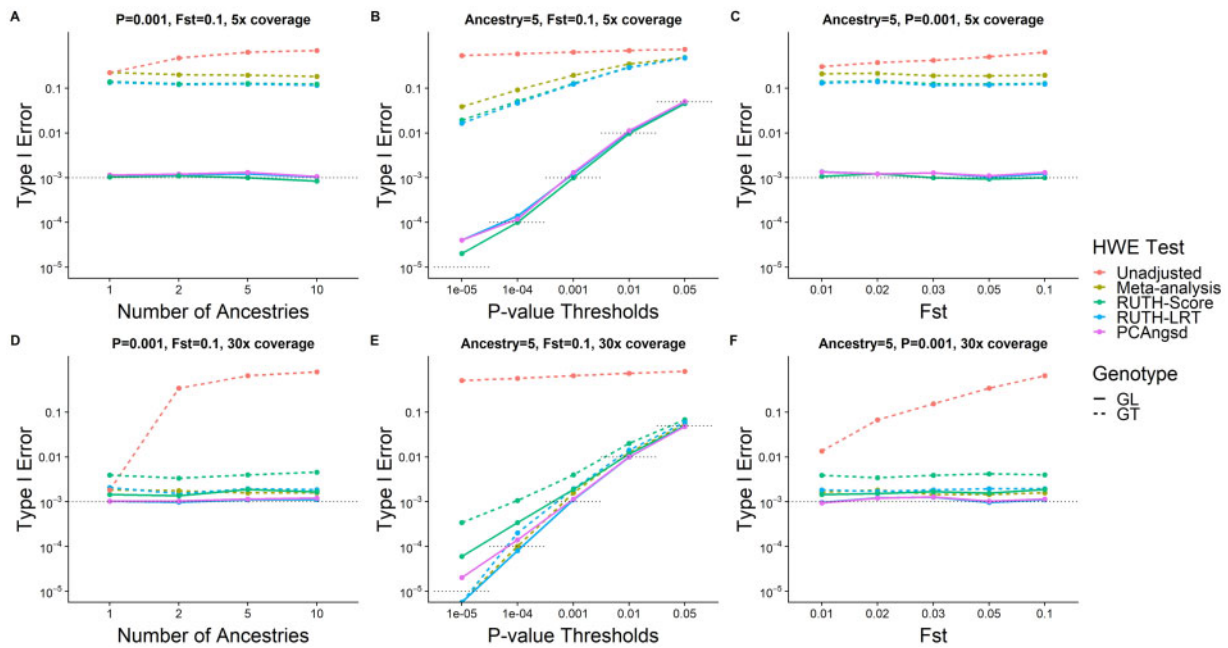


Figure 1 Evaluation of Type I Errors between various HWE tests on simulated GTs. Under each combination of simulation conditions (number of ancestries, sequencing coverage, and fixation index), we simulated 5000 samples with 50,000 variants that follow HWE within each of the subpopulations and determined the Type I error performances of different HWE tests based on the proportion of variants labeled as having significant P-values. Five HWE tests—(1) Unadjusted HWE test (Wigginton *et al.* 2005) implemented in PLINK-1.9 (Purcell *et al.* 2007) using hard GTs, (2) meta-analysis using Stouffer’s method across ancestries using hard GT, (3) RUTH test using hard GTs, (4) RUTH test using phred-scale likelihood (GL) computed from simulated sequence reads, and (5) PCAngsd (Meisner and Albrechtsen 2019)—were tested under HWE with various parameter settings. Gray dotted lines indicate targeted Type I Error rates. Top panels (A–C) represent results from shallow sequencing (5×), and the bottom panels (D–F) represent results from deep sequencing (30×). Using GL-based GTs resulted in Type I Error rates closer to the targeted rate than using GT-based GTs across different numbers of ancestries (A, D), P-value thresholds (B, E), and fixation indices (C, F). The difference is especially large for low-coverage GTs.

and Supplementary Table S1). For high-coverage sequence data, the power of GL-based tests at the same P-value threshold increased to ~56–60%, comparable to corresponding GT-based tests. Interestingly, the unadjusted GT-based test showed much lower power than RUTH and PCAngsd tests under excess heterozygosity ($\theta < 0$) while demonstrating much higher power with excess homozygosity ($\theta > 0$). Upon further investigation, we observed that the tests have lower power than the exact test specifically for rare variants with excess homozygosity owing to the mismatch between the empirical and the asymptotic null distributions (for details, see Discussion section).

We also generated PRC and receiver-operator characteristic (ROC) curves to better understand the tradeoff between the Type I errors and power under moderate excess heterozygosity ($\theta = -0.05$) (Supplementary Figure S3, C and D). Again, accounting for GT uncertainty resulted in better empirical power and Type I error, especially for low-coverage data: at an empirical false-positive rate of 1%, GL-based tests had 41–45% power, as opposed to 4–10% for GT-based tests. For high-coverage data, GL-based tests had 1–2% greater power than GT-based tests at the same false-positive rate. These results suggest that ignoring GT uncertainty in HWE tests is reasonable for high-coverage sequence data.

Simulation: Impact of Population Structure on HWE Test Statistics

As expected, the unadjusted HWE test had substantially inflated Type I errors under population structure based on the model of Balding and Nichols (1995) (Figure 1 and Supplementary Table S1). Even for an intra-continental level of population differentiation ($F_{ST} = 0.01$), the Type I errors at $P < 0.001$ were inflated 13.5-fold even for high-coverage data. With an inter-continental level

of differentiation ($F_{ST} = 0.1$), we observed orders of magnitude more Type I errors across different simulation conditions. This inflation is expected to increase with larger sample sizes, suggesting that adjustment for population structure is important even if a study focuses on a single continental population.

One simple approach to account for population structure is to stratify individuals into distinct subpopulations and apply HWE tests separately, as was done in UK Biobank (Bycroft *et al.* 2018), then meta-analyze the results (Figure 3B). Type I errors were appropriately controlled with this approach in high-coverage but not low-coverage data, likely owing to unmodeled GT uncertainty (Figure 1 and Supplementary Table S1). Instead of classifying individuals into distinct subpopulations, RUTH incorporates PCs to jointly perform HWE tests (Figure 3C). By estimating individual-specific allele frequencies, RUTH was able to adjust for the simulated population structure. For both low- or high-coverage data, GL-based RUTH tests and PCAngsd showed well-controlled Type I errors, whereas GT-based tests showed slight (high-coverage) to severe (low-coverage) inflation.

Although meta-analysis resulted in well-controlled Type I errors for high-coverage data, it was considerably less powerful than RUTH. For example, with moderate excess heterozygosity ($\theta = -0.05$) across five ancestries ($F_{ST} = 0.1$), RUTH tests identified 20–27% more variants as significant at $P < 0.001$ (Figure 2 and Supplementary Table S1) compared to meta-analysis. PRCs also clearly showed better operating characteristics for RUTH and PCAngsd compared to meta-analysis (Supplementary Figure S4). For example, at an empirical false-positive rate of 1%, RUTH showed much greater power (66–68%) than meta-analysis (43%) although the simulation scenario favors meta-analysis because samples were perfectly classified into distinct subpopulations.

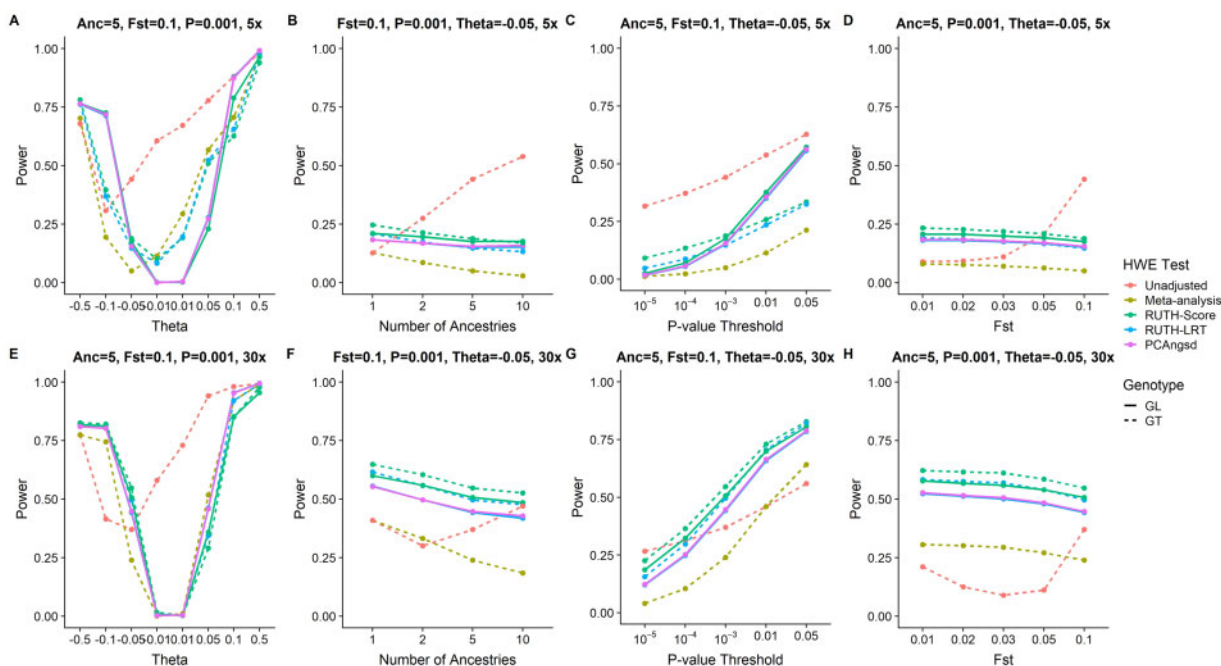


Figure 2 Evaluation of power between different HWE tests on simulated GTs. Under each combination of simulation conditions (number of ancestries, sequencing coverage, fixation index, and deviation from HWE), we simulated 50,000 variants for 5000 samples and evaluated the ability of different HWE tests to find the variants significant. Unless otherwise specified, the default simulation parameters are five ancestries, with $F_{ST} = 0.1$, P-value threshold = 0.001, and $\Theta = -0.05$. Tests that can find a larger proportion of significant variants are considered more powerful. Five HWE tests—(1) Unadjusted HWE test (Wigginton et al. 2005) implemented in PLINK-1.9 using hard GTs, (2) RUTH test using hard GTs, (3) RUTH test using phred-scale likelihood (PL) computed from simulated sequence reads, (4) meta-analysis using Stouffer’s method across ancestries using hard GTs, and (5) PCAngsd (Meisner and Albrechtsen 2019)—were tested for variants deviating from HWE with various parameter settings, for low coverage (A–D) and high coverage (E–H) data. (A, E) Θ controls the degree of deviation from HWE, with negative values indicating excess heterozygosity and positive values, indicating heterozygote depletion. The high Type I Error rates in GT-based tests (Figure 2) lead to those methods appearing to have higher power in some scenarios. The unadjusted test suffers from this problem the most. GL-based methods have slightly lower powers than GT-based methods in exchange for a much better controlled Type I error rate. This pattern mostly holds across different numbers of ancestries (B, F), P-value thresholds (C, G), and fixation indices (D, H). Meta-analysis had the lowest power in the presence of excess heterozygosity.

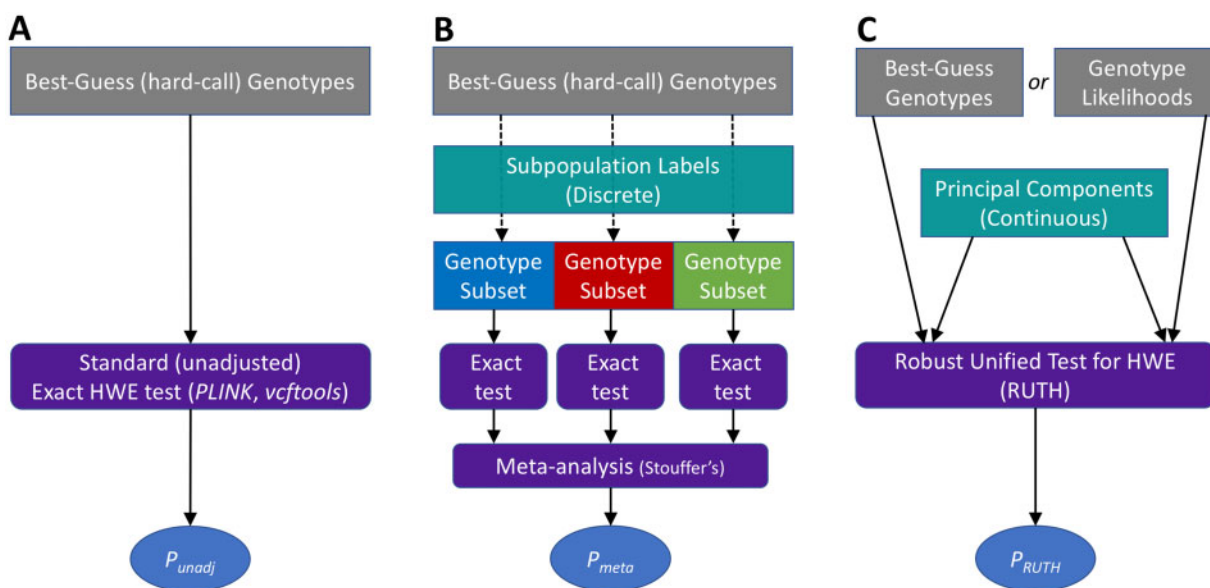


Figure 3 Schematic diagrams of different methods to test HWE under population structure. Three different methods to test HWE under population structure are described. (A) In the standard (unadjusted) HWE test, all samples are tested together using best-guess GTs. This test does not adjust for sample ancestry. (B) In a meta-analysis of stratified HWE tests, the samples must first be categorized into discrete subpopulations, determined *a priori* based on their GTs or self-reported ancestries. Next, standard HWE tests (based on best-guess GTs) are performed on each of these subpopulations. Then, the resulting HWE statistics are converted into Z-scores and combined in a meta-analysis using Stouffer’s method, with the sample sizes of the subpopulations as weights. (C) In our proposed method (RUTH), either best-guess GTs or GLs can be used as input for HWE test. We assume that the genetic ancestries of each sample are estimated *a priori*, typically as PCs. We combine the GTs and PCs to perform either a score test or an LRT to obtain a joint ancestry-adjusted HWE statistic for each variant across all samples.

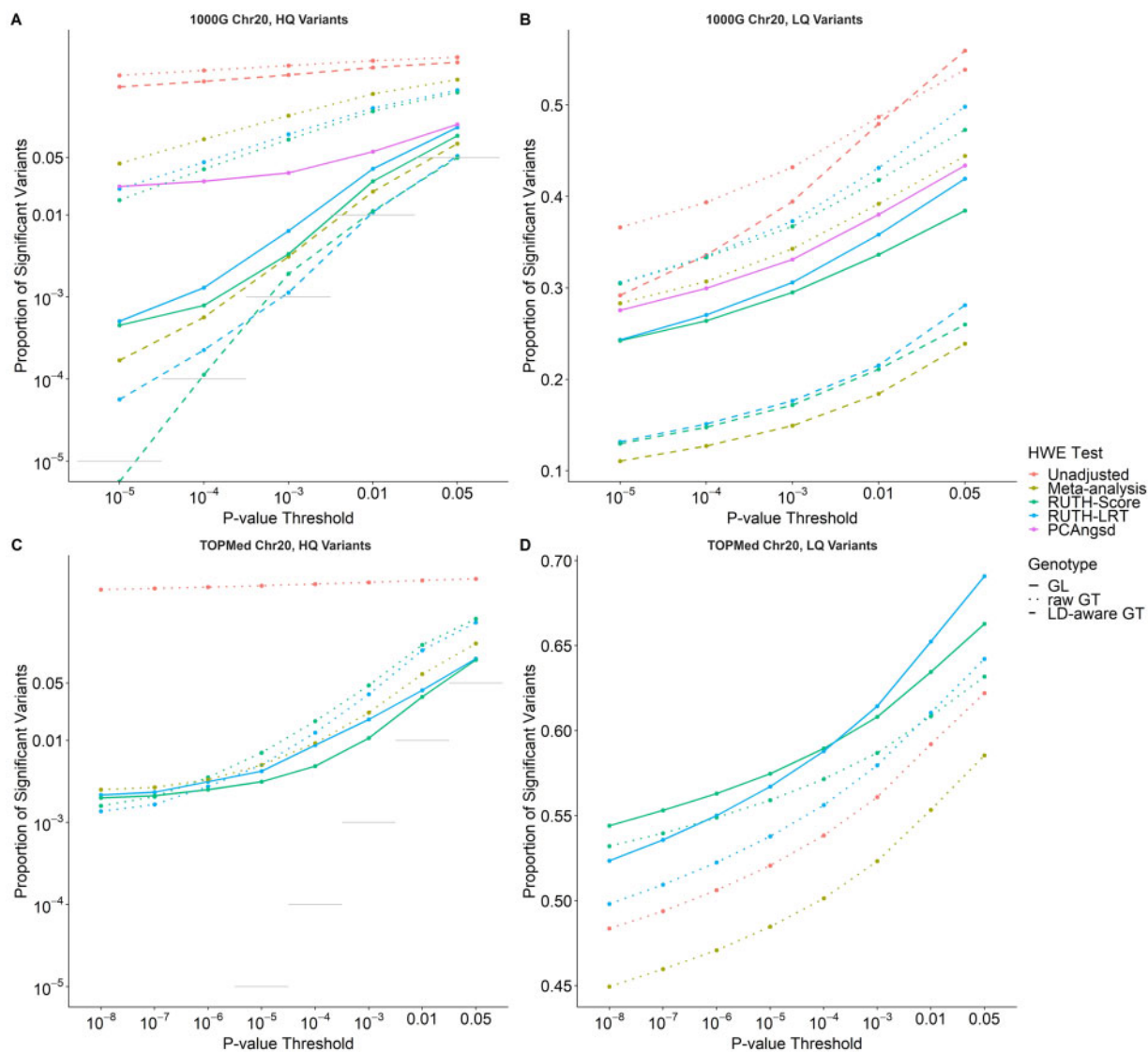


Figure 4 Evaluation of different HWE tests on 1000G and TOPMed variants. In 1000G data (A, B), we identified 17,740 HQ variants and 10,966 LQ variants in chromosome 20. In TOPMed data (C, D), we identified 17,524 HQ variants and 329,699 LQ variants in chromosome 20. A well-behaved HWE test should maximize the proportion of significant LQ variants while controlling the false-positive rate for HQ variants. Dotted gray lines represent targeted Type I error levels if we assume all HQ variants follow HWE. (A) Both the unadjusted test and the PCAngsd found substantially more significant variants than expected in the 1000G HQ variant set, whereas both RUTH and meta-analysis were more conservative. Methods that used raw GTs showed substantial false-positive rates, whereas methods that used GLs and LD-aware GTs had much better control of false positives. (B) In 1000G LQ variants, meta-analysis lagged behind RUTH and the unadjusted test in discovering significant deviation from HWE. RUTH behaved well for HQ variants while having more power to find low-quality variants significantly deviating from HWE. (C) In TOPMed data, the unadjusted test resulted in an excess of false positives. Tests using GL-based GTs outperformed tests using GT-based GTs. (D) Methods using GL-based GTs were able to discover more LQ variants than methods using GT-based GTs, demonstrating the advantage of accounting for GT uncertainty in HWE tests.

When stratified by allele frequency, RUTH showed better operating characteristics for common variants compared to rare variants owing to a difference in power (Supplementary Figure S5).

Application to 1000 Genomes whole-genome sequence data

Next, we evaluated the performance of various HWE tests in low-coverage ($\sim 6\times$) sequence data from the 1000G. We evaluated three representations of GTs: (1) raw GT, (2) LD-aware GT, and (3) GL, as described in Materials and Methods section. Among chromosome 20 variants, we selected 17,740 HQ variants that are polymorphic in GWAS arrays, and 10,966 LQ variants enriched for GT discordance in duplicates and trios. Unlike simulation studies, not all LQ variants are expected to violate HWE, so we

consider the proportion of significant LQ variants as a lower bound for the sensitivity to identify significant variants. Similarly, not all HQ variants are expected to follow HWE, so the proportion of significant HQ variants serves as an upper bound for the false-positive rate.

Consistent with simulation results, all tests based on raw GTs generated from low-coverage sequence data had severe inflations of false positives (Figure 4A and Table 1). This was true even for HQ variants, presumably owing to genotyping errors and bias in raw GTs. Standard HWE tests, which model neither GT uncertainty nor population structure, showed the highest inflation of false positives at 44% for $P < 10^{-6}$, a threshold commonly used for HWE testing in large genetic studies (Locke et al. 2015; Fritsche et al. 2016). Modeling population structure substantially reduced

Table 1 Performance of the unadjusted test, meta-analysis, RUTH, and PCAngsd on 1000G chromosome 20 variants

| Variant category | GT format | HWE test | Proportion of significant variants | | | | | Total variant count | |
|------------------|-------------|---------------|------------------------------------|----------------------|----------------------|----------------------|----------------------|---------------------|--------|
| | | | $P < 10^{-2}$ | $P < 10^{-3}$ | $P < 10^{-4}$ | $P < 10^{-5}$ | $P < 10^{-6}$ | | |
| LQ variants | Raw GT | Unadjusted | 0.487 | 0.432 | 0.394 | 0.366 | 0.339 | 10,966 | |
| | | Meta-analysis | 0.392 | 0.343 | 0.307 | 0.283 | 0.262 | 10,966 | |
| | | RUTH score | 0.418 | 0.367 | 0.333 | 0.305 | 0.284 | 10,966 | |
| | LD-aware GT | RUTH-LRT | 0.431 | 0.373 | 0.335 | 0.305 | 0.280 | 10,966 | |
| | | Unadjusted | 0.479 | 0.395 | 0.336 | 0.292 | 0.259 | 10,966 | |
| | | Meta-analysis | 0.184 | 0.149 | 0.127 | 0.111 | 0.098 | 10,966 | |
| | | RUTH score | 0.211 | 0.172 | 0.147 | 0.130 | 0.112 | 10,966 | |
| | | RUTH-LRT | 0.215 | 0.177 | 0.151 | 0.131 | 0.115 | 10,966 | |
| | | GL | RUTH score | 0.336 | 0.295 | 0.264 | 0.242 | 0.223 | 10,966 |
| | HQ variants | Raw GT | RUTH-LRT | 0.358 | 0.306 | 0.270 | 0.243 | 0.225 | 10,966 |
| | | | PCAngsd | 0.380 | 0.331 | 0.300 | 0.275 | 0.255 | 10,920 |
| | | | Unadjusted | 0.755 | 0.657 | 0.573 | 0.501 | 0.443 | 17,740 |
| LD-aware GT | | Meta-analysis | 0.298 | 0.161 | 0.084 | 0.042 | 0.020 | 17,740 | |
| | | RUTH score | 0.183 | 0.083 | 0.036 | 0.015 | 7.4×10^{-3} | 17,740 | |
| | | RUTH-LRT | 0.200 | 0.095 | 0.044 | 0.021 | 0.010 | 17,740 | |
| | | Unadjusted | 0.623 | 0.507 | 0.422 | 0.361 | 0.311 | 17,740 | |
| | | Meta-analysis | 0.019 | 3.1×10^{-3} | 5.6×10^{-4} | 1.7×10^{-4} | 1.1×10^{-4} | 17,740 | |
| | | RUTH score | 0.011 | 1.9×10^{-3} | 1.1×10^{-4} | 0 | 0 | 17,740 | |
| GL | | RUTH-LRT | 0.011 | 1.1×10^{-3} | 2.3×10^{-4} | 5.6×10^{-5} | 0 | 17,740 | |
| | | RUTH score | 0.026 | 3.3×10^{-3} | 7.9×10^{-4} | 4.5×10^{-4} | 3.4×10^{-4} | 17,740 | |
| | | RUTH-LRT | 0.036 | 6.4×10^{-3} | 1.3×10^{-3} | 5.1×10^{-4} | 3.4×10^{-4} | 17,740 | |
| | PCAngsd | 0.059 | 0.032 | 0.026 | 0.022 | 0.021 | 17,740 | | |

The numbers within cells represent the proportions of significant variants under the corresponding testing conditions at the given P-value threshold. We expect our LQ variants to violate HWE at a higher rate than our HQ variants. A well-behaved test is expected to find a high proportion of LQ variants to be significant while maintaining the targeted Type I Error rate in HQ variants. The unadjusted test consistently shows the highest false-positive rate among all the tests. HWE tests that rely on raw GTs also show much higher false-positive rates than tests that use other GT representations. RUTH tests were the best at controlling false positives while still maintaining comparable power to the other methods. PCAngsd had a much higher false-positive rate than RUTH-based methods, especially at more stringent P-value thresholds.

inflation, with RUTH tests showing fewer false positives (0.7–1.0% at $P < 10^{-6}$) than meta-analysis (2.0% at $P < 10^{-6}$). False positives were inflated across all methods when using raw GTs.

Similarly, GL-based RUTH tests further reduced false positives (0.034% at $P < 10^{-6}$). In contrast to our simulations, however, PCAngsd demonstrated considerably higher false positives than RUTH (2.1% at $P < 10^{-6}$) because PCAngsd estimates PCs from the input data without the ability to use externally provided PCs (Discussion section). The sensitivity for detecting significant LQ variants was also consistent with our simulations (Figure 4B and Table 1). GL-based tests, which showed better control of false positives, identified 22–25% of LQ variants as significant at $P < 10^{-6}$.

Strikingly, while using LD-aware GTs reduced false positives with adjusted tests, it was at the expense of substantially reduced sensitivity to detect LQ variants. The false-positive rates of any adjusted test with LD-aware GTs were uniformly lower than those of any GL- and raw GT-based tests across all P-value thresholds (Figure 4A). However, sensitivity was also substantially reduced with LD-aware GTs (Figure 4B). For example, at $P < 10^{-6}$, GL-based RUTH tests identified 22–23% of LQ variants as significant, while using LD-aware GTs halved the proportions. Meta-analysis with LD-aware GTs had even lower sensitivity, likely because the implicit HWE assumption in LD-aware GT refinement altered the LD-aware GTs to conform to HWE, further reducing both false positives and sensitivity.

We evaluated PRCs between HQ and LQ variants to further evaluate this tradeoff. The results clearly demonstrated that HWE tests using LD-aware GTs are substantially less robust than tests using other GT representations (Supplementary Table S2 and Figure S6A). For example, for the RUTH score test, when LD-aware GTs identified 0.1% of HQ variants as significant, 17% of LQ variants were identified as significant. However, with raw GT and GL, 24–27% were identified as significant at the same

threshold. Even fewer were significant in meta-analysis with LD-aware GTs (13%). Similar trends were observed across all thresholds, suggesting that using LD-aware GTs results in substantially poorer operating characteristics. As more accurate genotyping in LD-aware GT refinement is expected to improve the performance of QC metrics compared to raw GTs, these results are quite striking, and highlight a potential oversight in using LD-aware GTs in various QC metrics for sequence-based GTs. It should also be noted that the significance threshold we used can be subjective (Discussion section), but the relative trends between the methods largely remained similar (Table 1).

Application to TOPMed Deep whole-genome sequence data

We evaluated the various HWE tests on a subset of the Freeze 5 variant calls from high-coverage (~37×) whole-genome sequence data in the TOPMed Project (Taliun et al. 2021). We identified 17,524 HQ variants and 329,699 LQ variants using the same criteria used for 1000G variants and evaluated raw GTs and GLs. We did not evaluate PCAngsd owing to excessive computational time (see Computational cost section).

We first evaluated the false-positive rates of different HWE tests indirectly by using HQ variants. With a >20-fold larger sample size than 1000G, we identified more significant HQ variants, whereas the false-positive rates were still reasonable with adjusted tests. At $P < 10^{-6}$, 74% of HQ variants were significant with unadjusted tests, whereas the adjusted GL-based tests identified ~0.3% at $P < 10^{-6}$ (Figure 4, C and D; Table 2). Adjusted GT-based tests had only slightly higher levels of false positives at $P < 10^{-6}$. However, inflation was more noticeable at less stringent P-value thresholds, suggesting that GL-based tests may be needed for larger sample sizes.

Next, we evaluated the proportions of LQ variants found to be significant by different tests to indirectly evaluate their statistical

Table 2 Performance of the unadjusted test, meta-analysis, and RUTH on TOPMed freeze 5 chromosome 20 variants

| Variant set | GT format | HWE test | Proportion of significant variants | | | | | Total variant count |
|-------------|-----------|---------------|------------------------------------|---------------|----------------------|----------------------|----------------------|---------------------|
| | | | $P < 10^{-2}$ | $P < 10^{-3}$ | $P < 10^{-4}$ | $P < 10^{-5}$ | $P < 10^{-6}$ | |
| LQ Variants | Raw GT | Unadjusted | 0.592 | 0.561 | 0.539 | 0.521 | 0.506 | 329,699 |
| | Raw GT | Meta-analysis | 0.554 | 0.524 | 0.502 | 0.485 | 0.471 | 329,699 |
| | Raw GT | RUTH score | 0.608 | 0.587 | 0.572 | 0.559 | 0.549 | 329,699 |
| | GL | RUTH score | 0.635 | 0.608 | 0.590 | 0.575 | 0.563 | 329,699 |
| | Raw GT | RUTH-LRT | 0.610 | 0.580 | 0.556 | 0.538 | 0.522 | 329,699 |
| | GL | RUTH-LRT | 0.653 | 0.615 | 0.588 | 0.567 | 0.550 | 329,699 |
| HQ Variants | Raw GT | Unadjusted | 0.890 | 0.842 | 0.800 | 0.766 | 0.736 | 17,524 |
| | Raw GT | Meta-analysis | 0.065 | 0.022 | 9.0×10^{-3} | 4.8×10^{-3} | 3.3×10^{-3} | 17,524 |
| | raw GT | RUTH score | 0.145 | 0.047 | 0.172 | 7.1×10^{-3} | 3.5×10^{-3} | 17,524 |
| | GL | RUTH score | 0.034 | 0.011 | 4.9×10^{-3} | 3.1×10^{-3} | 2.5×10^{-3} | 17,524 |
| | raw GT | RUTH-LRT | 0.125 | 0.036 | 0.012 | 5.0×10^{-3} | 2.7×10^{-3} | 17,524 |
| | GL | RUTH-LRT | 0.041 | 0.018 | 8.5×10^{-3} | 4.3×10^{-3} | 3.1×10^{-3} | 17,524 |

The numbers within cells represent the proportions of significant variants under the corresponding testing conditions at the given P -value threshold. These results are based on the tests that used likelihood-based GT representations as input. A well-behaved test should reduce the number of significant HQ variants while increasing the number of significant low-quality (LQ) variants. The unadjusted test had a greatly inflated false-positive rate for HQ variants while showing a lower true positive rate for LQ variants. While meta-analysis performed better for HQ variants, it had reduced power to find LQ variants to be significant. RUTH performed the best, with fewer false positives (significant HQ variants) compared to both the unadjusted test and the meta-analysis, while at the same time finding more true positives (significant LQ variants).

Table 3 Runtimes for RUTH and PCAnsd on simulated data

| Sample size | Wall Time (s) | | | User Time (s) | | |
|-------------|---------------|------------|------------|---------------|------------|-----------|
| | RUTH-LRT | RUTH score | PCAnsd | RUTH-LRT | RUTH score | PCAnsd |
| 1,000 | 16.21 | 27.24 | 173.11 | 16.16 | 27.09 | 172.37 |
| 2,000 | 32.19 | 54.63 | 347.10 | 31.94 | 54.51 | 345.58 |
| 5,000 | 82.80 | 136.44 | 1,124.83 | 81.81 | 136.20 | 1,102.85 |
| 10,000 | 165.48 | 273.67 | 7,396.00 | 163.88 | 273.27 | 7,235.91 |
| 20,000 | 336.75 | 553.92 | 38,807.67 | 332.06 | 553.05 | 37,338.69 |
| 50,000 | 902.81 | 1,438.32 | 461,971.33 | 886.67 | 1,435.87 | 403,296.5 |

We simulated 10,000 GL-based variants for varying numbers of samples. Wall time indicates total runtime, whereas user time is the amount of time the CPUs spent running each program. All programs were run in single-threaded mode. System processes make up the difference between the two values, with a majority consisting of file I/O. We used VCF files with GL fields in RUTH and converted them to Beagle3 format for PCAnsd. The RUTH LRT was the fastest method, with the score test about 60% slower. PCAnsd was about 10 times slower than RUTH-LRT with the smallest sample sizes and over 400 times slower with our largest tested size of 50,000 samples.

power. GT- and GL-based RUTH tests showed similar power, whereas meta-analysis showed considerably lower power. For example, at $P < 10^{-6}$, meta-analysis identified 47% of LQ variants as significant, whereas RUTH tests identified 54–58%. This pattern was similar across different P -value thresholds (Figure 4, C and D) or choices of LQ variants (Supplementary Table S3 and Figure S7). Our results suggest that GL-based RUTH tests are suitable for testing HWE for tens of thousands of deeply sequenced genomes with diverse ancestries, and that using raw GTs will also result in a comparable performance at typically used HWE P -value thresholds (e.g., $P < 10^{-6}$).

We used PRCs to evaluate the tradeoff between empirical false-positive rates and power. Consistent with the previous results, the GL-based RUTH test showed the best tradeoff between false positives and power, whereas the GT-based RUTH test and meta-analysis were slightly less robust but largely comparable (Supplementary Figure S6). Notably, when we evaluated the different methods at an empirical false-positive rate of 0.1%, RUTH score tests had ~4% higher power than RUTH LRT for both raw GTs and GLs (Supplementary Figures S8 and S9).

Impact of ancestry estimation accuracy on Hardy–Weinberg equilibrium tests

So far, our evaluations relied on genetic ancestry estimates carefully determined with sophisticated methods (Materials and Methods section). However, using simpler approaches instead

during the variant QC step may affect the performance of adjusted HWE tests. We evaluated whether the number of PC coordinates affected the performance of RUTH tests by comparing the use of 2 vs 4 PCs (default). The results from both simulated and real data sets consistently demonstrated that using 4 PCs led to substantially reduced Type I errors compared to using 2 PCs at a similar level of power (Supplementary Tables S2 and S4, Figure S10). PRCs also clearly showed that using 4 PCs was more robust against population structure across both simulated and real data sets (Supplementary Figure S11).

We also evaluated whether the classification accuracy of subpopulations affected the performance of meta-analysis. Instead of assigning 1000G individuals into five continental populations, we used the k -means algorithm on those samples' top 2 PCs to classify them into three crude subpopulations (Supplementary Figure S12). This led to a much higher false-positive rate with virtually no increase in true positives (Supplementary Figure S13 and Table S2). We saw the same pattern in simulated data (Supplementary Figure S11 and Table S5).

Computational cost

We compared the computational costs of RUTH and PCAnsd for simulated and real data. RUTH has linear time complexity to sample size, whereas PCAnsd appears to have quadratic time complexity owing to joint estimation of PCs (Table 3 and Supplementary Table S6). RUTH also has low memory

requirements compared to PCAnsd (e.g., 14 MB vs 2 GB for 1000G data). Extrapolating our results to the whole-genome scale, analyzing 1000G (i.e., 80 million variants) is expected to take 120 CPU-hours for RUTH and 3200 CPU-hours for PCAnsd (with >1 TB memory consumption). Additionally, RUTH can be parallelized into smaller regions in a straightforward manner.

Discussion

RUTH is a unified, flexible, and robust approach to incorporate genetic ancestry and GT uncertainty for testing Hardy–Weinberg Equilibrium capable of handling large amounts of GT data with structured populations. Sha and Zhang (2011) proposed HWES, an HWE test for structured populations, to address some of these challenges, but it has not been widely used owing to the lack of an implementation that supports popular GT data formats (e.g., PED, BED, VCF, or BCF) and inability to handle imputed or uncertain GTs. Hao et al. (2016) proposed sHWE which can only handle best-guess (hard call) GTs (i.e., 0, 1, or 2 for biallelic variants) and does not account for GT uncertainty. Meisner and Albrechtsen (2019) proposed PCAnsd to address some of these issues, but it does not support the standard VCF/BCF formats for sequence-based GTs, and its current implementation scales poorly with genome-wide analyses of large samples.

Similar to previous studies (Sha and Zhang 2011; Hao et al. 2016), our proposed framework uses individual-specific allele frequencies rather than allele frequencies pooled across all samples to systematically account for population structure in HWE tests. Unlike those previous studies, we model GT uncertainty in sequence-based GTs using a likelihood-based framework. We implemented two RUTH tests—a score test and a LRT—to test for HWE under population structure for GTs with uncertainty. While RUTH LRT is similar to the independently developed PCAnsd, the software implementation of RUTH is more flexible, scales much better to large studies, and supports the standard VCF format.

We provide a comprehensive evaluation of various approaches for testing HWE using simulated and real data. Our results demonstrate that modeling population stratification is necessary for HWE tests on heterogeneous populations. We showed that accounting for GT uncertainty via GT likelihoods performs substantially better than using best-guess GTs, especially for low-coverage sequenced genomes. Importantly, we included evaluations for an unpublished but commonly used approach—meta-analysis across stratified subpopulations, cohorts, or batches. Our results demonstrate that while meta-analysis may be effective in reducing false positives, it does so at the expense of substantially reduced power compared to RUTH.

We observed that the current implementation of PCAnsd does not scale well to large-scale sequencing data, though in principle, it can be implemented more efficiently, because the underlying HWE test itself is similar to RUTH LRT. PCAnsd requires loading all GTs into memory, which is often infeasible for large sequencing studies. For example, loading all of 1000G will require ~4.8 TB of memory. In our evaluation of 1000G chromosome 20 variants, the inability of PCAnsd to estimate PCs from the whole genome may have contributed to the observed difference in results from RUTH compared to our simulation studies. Moreover, PCAnsd does not offer an option to externally provide PCs or exclude false-positive variants when calculating PCs, so it performs poorly when false-positive variants confound PC estimation as demonstrated in the 1000G examples.

Although our 1000G experiments demonstrated the unexpected result that using raw GTs had better sensitivity than using LD-aware GTs at the same empirical false-positive rates for low-coverage data, we do not advocate using raw GTs for low-coverage sequence data. First, the results for raw GTs were still consistently less robust than GL-based RUTH tests. Moreover, it would be tricky to determine an appropriate P -value threshold when false positives are severely inflated. Therefore, we strongly advocate using GL-based RUTH tests for robust HWE tests with low-coverage sequence data. For the now more typical high-coverage sequence data, GL-based tests are still preferred, but GT-based RUTH tests should be acceptable for cases in which GLs are unavailable.

Our experiment compared using 2 vs 4 PCs only because the *verifyBamID2* software tool estimated up to 4 PCs projected onto the HGDP panel by default (Zhang et al. 2020). Because our method focuses on testing HWE during the QC steps in sequence-based variant calls, a curated version of PCs, estimated from the sequenced cohort themselves, may not be readily available. However, it is possible to use a larger number of PCs (e.g., >10 PCs) if available at the time of HWE test. We expect that a larger number of PCs will account for finer-grained population structure and may improve the performance of HWE tests, but additional experiments are needed to quantify the effect.

Our results demonstrate that RUTH score and LRT tests perform similarly in simulated and experimental data sets. Overall, the RUTH-LRT was slightly more powerful than the RUTH score test at the expense of slightly greater false-positive rates although this tendency was not consistent. We observed that the RUTH tests tended to be slightly more powerful in identifying deviation from HWE in the direction of excess heterozygosity than excess homozygosity when compared to adjusted meta-analysis. These results might be caused by the difference between our model-based asymptotic tests compared to the exact test used in meta-analysis.

We did not evaluate our methods on imputed GTs in this manuscript. Because imputed GTs implicitly assume HWE, we suspect that HWE tests based on imputed GTs may have reduced power compared to directly genotyped variants. It is possible to use approximate GLs instead of best-guess GTs for imputed GTs, but this requires GT probabilities, not just GT dosages. If GT probabilities $\Pr(g_i = G | \text{Data}_i)$ are available, they can be converted to GLs $L_i^{(G)} = \Pr(\text{Data}_i | g_i = G)$ using Bayes' rule by modeling $\Pr(g_i = G)$ as a binomial distribution based on allele frequencies (which implicitly assumes HWE). However, similar to LD-aware GTs in low-coverage sequencing, the power of HWE tests with imputed GTs may be poor. Further evaluation is needed to understand the effect of using imputed GTs on the behavior of HWE tests.

As described in our results, we observed that the current implementations of RUTH (and PCAnsd) tests relying on asymptotic distributions do not work more robustly than the exact test when testing for excess homozygosity ($\theta > 0$). This is mainly because the empirical null distribution becomes increasingly asymmetric between the two directions of effects for rarer variants, but the asymptotic approximation assumes symmetry between them, causing loss of power for excess homozygosity. Using RUTH score test will further reduce power because score tests are known to have reduced power than LRT when θ strongly deviates from zero, which happens in rare variants with excess homozygosity. Applying Saddlepoint approximation (Dey et al., 2017) or similar techniques may help address this issue.

In practice, when we examined LQ variants, determined by high Mendelian errors, the vast majority (65% for 1000G, 82% for TOPMed) of them deviated from HWE toward excess heterozygosity ($\theta < 0$) as opposed to excess homozygosity ($\theta > 0$) when we examined the direction of deviation from HWE regardless of its significance. On the other hand, the majority of HQ variants (77% for 1000G, 64% for TOPMed) mildly deviated from HWE toward excess homozygosity ($\theta > 0$), presumably owing to residual population structure and cryptic relatedness. These observations suggest that detecting excess heterozygosity is practically more important for variant QC, on which RUTH tests are expected to perform well.

Our methods have room for further improvement. First, we used a truncated linear model for individual-specific allele frequencies for computational efficiency. Although such an approximation was demonstrated to be effective in practice (Zhang *et al.* 2020), applying a logistic model or some other more sophisticated model may be more effective in improving the precision and recall of RUTH tests. Second, we did not attempt to model or evaluate the effect of admixture in our method. Because HWE is reached in two generations with random mating, accounting for admixed individuals may only have a marginal impact. On the other hand, admixture can lead to higher observed heterozygosity. It may be possible to improve RUTH by explicitly modeling and adjusting for the effect of admixture on individual-specific allele frequencies. Systematic evaluations focusing on admixed populations are needed to evaluate whether an admixture adjustment is necessary. Third, RUTH tests do not account for family structure or individual-level inbreeding. We suspect that the apparent inflation of Type I error for the TOPMed data was partially owing to sample relatedness. Accounting for family structure or individual-level inbreeding in other ways, for example using variance components models, will require much longer computational times and may not be feasible for large-scale data sets. Fourth, RUTH currently does not directly support imputed GTs or GT dosages. In principle, it is possible to convert posterior probabilities for imputed GTs into GT likelihoods to account for GT uncertainty (by using individual-specific allele frequencies). However, because most GT imputation methods implicitly assume HWE, we suspect that HWE tests on imputed GTs will be underpowered, similar to our observations with LD-aware GTs in the 1000G data set, even though explicitly modeling posterior probabilities may slightly mitigate this reduction in power.

The choice of a P-value threshold to indicate deviation from HWE remains an open question. In previous studies, stringent P-value thresholds were used to prevent high-quality variants from being filtered due to population structure. Adjusting for population structure with RUTH helps mitigate this problem, allowing the use of less stringent thresholds to improve test performance, but the choice of P-value threshold remains subjective, based on the trade-off between sensitivity and specificity. Future development of more robust methods to determine significance thresholds would help further improve the use of HWE tests for variant quality control.

In summary, we have developed and implemented robust and rapid methods and software tools to enable HWE tests that account for population structure and GT uncertainty. We comprehensively evaluated both our methods and alternative approaches. Our tools can be used to evaluate variant quality in very large-scale genetic data sets, with the ability to handle standard VCF formats for storing sequence-based GTs. Our software tools are publicly available at <http://github.com/statgen/ruth>.

Acknowledgments

TOPMed source studies and sample counts are listed in [Supplementary Table S7](#). Acknowledgements for TOPMed omics support are detailed in [Supplementary Table S8](#). Full TOPMed study acknowledgements are listed in [Supplementary File S1](#).

Funding

This work was supported by National Institutes of Health grants HL137182 (from National Heart, Lung, and Blood Institute), HG009976 (from National Human Genome Research Institute), HG007022 (from National Human Genome Research Institute), DA037904 (from National Institute of Drug Abuse), HL117626-05-S2 (from National Heart, Lung, and Blood Institute), and MH105653 (from National Institute of Mental Health). Molecular data for the Trans-Omics in Precision Medicine (TOPMed) program were supported by the National Heart, Lung and Blood Institute. Core support including centralized genomic read mapping and GT calling, along with variant quality metrics and filtering were provided by the TOPMed Informatics Research Center (3R01HL-117626-02S1; contract HHSN268201800002I). Core support including phenotype harmonization, data management, sample-identity QC, and general program coordination were provided by the TOPMed Data Coordinating Center (R01HL-120393; U01HL-120393; contract HHSN268201800001I). We gratefully acknowledge the studies and participants who provided biological samples and data for TOPMed.

Conflicts of interest

K.C.B. receives funds from the NIH and receives royalties from UpToDate. E.G.B. has received funds from the following: National Institute of Health, Lung, Blood Institute, National Institute of Health, General Medical Sciences, National Institute on Minority Health and Health Disparities, The Tobacco-Related Disease Research Program, Food and Drug Administration, and The Sandler Family Foundation. L.A.C. spends part of her time consulting for Dyslipidemia Foundation, a nonprofit company, as a statistical consultant. P.T.E. is supported by a grant from Bayer to the Broad Institute focused on the genetics and therapeutics of cardiovascular diseases. P.T.E. has also served on advisory boards or consulted for Quest Diagnostics and Novartis. S.A.L. receives sponsored research support from Bristol Myers Squibb/Pfizer, Bayer, Boehringer Ingelheim, and Fitbit, has consulted for Bristol Myers Squibb/Pfizer and Bayer, and participates in a research collaboration with IBM. M.E.M. is an inventor on a patent that was published by the United States Patent and Trademark Office on December 6, 2018 under Publication Number US 2018-0346888, and an international patent application that was published on December 13, 2018 under Publication Number WO-2018/226560 regarding B4GALT1 Variants and uses thereof. S.T.W. receives royalties from UpToDate. G.R.A. and H.M.K. are employees of Regeneron Pharmaceuticals, they own stock and stock options for Regeneron Pharmaceuticals.

Literature cited

- Altshuler DM, Gibbs RA, Peltonen L, Altshuler DM, The International HapMap Consortium, *et al.* 2010. Integrating common and rare genetic variation in diverse human populations. *Nature*. 467:52–58.
- Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, The 1000 Genomes Project Consortium, *et al.* 2015. A global reference for human genetic variation. *Nature*. 526:68–74.

- Balding DJ, Nichols RA. 1995. A method for quantifying differentiation between populations at multi-allelic loci and its implications for investigating identity and paternity. *Genetica*. 96:3–12.
- Balding DJ. 2003. Likelihood-based inference for genetic correlation coefficients. *Theor Popul Biol*. 63:221–230.
- Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, et al. 2018. The UK Biobank resource with deep phenotyping and genomic data. *Nature*. 562:203–209.
- Danecek P, Auton A, Abecasis G, Albers CA, Banks E, 1000 Genomes Project Analysis Group, et al. 2011. The variant call format and VCFtools. *Bioinformatics*. 27:2156–2158.
- Delaneau O, Zagury JF, Marchini J. 2013. Improved whole-chromosome phasing for disease and population genetic studies. *Nat Methods*. 10:5–6.
- Dempster AP, Laird NM, Rubin DB. 1977. Maximum likelihood from incomplete data via the EM algorithm. *J R Stat Soc B*. 39:1–22.
- Dey R, Schmidt EM, Abecasis GR, Lee S. 2017. A Fast and Accurate Algorithm to Test for Binary Phenotypes and Its Application to PheWAS. *The American Journal of Human Genetics*. 101:37–49. 10.1016/j.ajhg.2017.05.014
- Dryden IL, Mardia KV. 1998. *Statistical Shape Analysis*. Chichester; New York: John Wiley & Sons.
- Ewing B, Green P. 1998. Base-calling of automated sequencer traces using phred. II. Error probabilities. *Genome Res*. 8:186–194.
- Fritsche LG, Igl W, Bailey JN, Grassmann F, Sengupta S, et al. 2016. A large genome-wide association study of age-related macular degeneration highlights contributions of rare and common variants. *Nat Genet*. 48:134–143.
- Graffelman J, Moreno V. 2013. The mid p-value in exact tests for Hardy-Weinberg equilibrium. *Stat Appl Genet Mol Biol*. 12:433–448.
- Hao W, Song M, Storey JD. 2016. Probabilistic models of genetic variation in structured populations applied to global human studies. *Bioinformatics*. 32:713–721.
- Hao W, Storey JD. 2019. Extending tests of Hardy-Weinberg equilibrium to structured populations. *Genetics*. 213:759–770.
- Hardy GH. 1908. Mendelian proportions in a mixed population. *Science*. 28:49–50.
- Holsinger KE, Lewis PO, Dey DK. 2002. A Bayesian approach to inferring population structure from dominant markers. *Mol Ecol*. 11:1157–1164.
- Holsinger KE. 2004. Analysis of genetic diversity in geographically structured populations: a Bayesian perspective. *Hereditas*. 130:245–255.
- Jin Y, Schaffer AA, Feolo M, Holmes JB, Kattman BL. 2019. GRAF-pop: a fast distance-based method to infer subject ancestry from multiple genotype datasets without principal components analysis. *G3*. 9:2447–2461.
- Jun G, Flickinger M, Hetrick KN, Romm JM, Doheny KF, et al. 2012. Detecting and estimating contamination of human DNA samples in sequencing and array-based genotype data. *Am J Hum Genet*. 91:839–848.
- Kuhn RM, Haussler D, Kent WJ. 2013. The UCSC genome browser and associated tools. *Brief Bioinform*. 14:144–161.
- Laurie CC, Doheny KF, Mirel DB, Pugh EW, Bierut LJ, GENEVA Investigators, et al. 2010. Quality control and quality assurance in genotypic data for genome-wide association studies. *Genet Epidemiol*. 34:591–602.
- Li JZ, Absher DM, Tang H, Southwick AM, Casto AM, et al. 2008. Worldwide human relationships inferred from genome-wide patterns of variation. *Science*. 319:1100–1104.
- Li M, Li C. 2008. Assessing departure from Hardy-Weinberg equilibrium in the presence of disease association. *Genet Epidemiol*. 32:589–599.
- Locke AE, Kahali B, Berndt SI, Justice AE, Pers TH, The Lifelines Cohort Study, et al. 2015. Genetic studies of body mass index yield new insights for obesity biology. *Nature*. 518:197–206.
- McCarroll SA, Hadnott TN, Perry GH, Sabeti PC, Zody MC, International HapMap Consortium, et al. 2006. Common deletion polymorphisms in the human genome. *Nat Genet*. 38:86–92.
- Meisner J, Albrechtsen A. 2019. Testing for Hardy-Weinberg equilibrium in structured populations using genotype or low-depth NGS data. *Mol Ecol Res*. 19:1144–1152.
- Mosteller F, Fisher RA. 1948. Questions and answers. *Am Statist*. 2:16–31.
- Nielsen DM, Ehm MG, Weir BS. 1998. Detecting marker-disease association by testing for Hardy-Weinberg disequilibrium at a marker locus. *Am J Hum Genet*. 63:1531–1540.
- Nielsen R, Paul JS, Albrechtsen A, Song YS. 2011. Genotype and SNP calling from next-generation sequencing data. *Nat Rev Genet*. 12:443–451.
- Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, et al. 2006. Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet*. 38:904–909.
- Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, et al. 2007. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet*. 81:559–575.
- Rohlf RV, Weir BS. 2008. Distributions of Hardy-Weinberg equilibrium test statistics. *Genetics*. 180:1609–1616.
- Rosenberg NA, Pritchard JK, Weber JL, Cann HM, Kidd KK, et al. 2002. Genetic structure of human populations. *Science*. 298:2381–2385.
- Sha Q, Zhang S. 2011. A test of Hardy-Weinberg equilibrium in structured populations. *Genet Epidemiol*. 35:671–678.
- Stouffer SA, Suchman EA, DeVinney LC, Star SA, Williams RM, Jr. 1949. *The American soldier: adjustment during army life*. Stud Soc Psychol World War II. Princeton, NJ: Princeton University Press.
- Stouffer SA. 1949. *The American Soldier*. Princeton: Princeton University Press.
- Taliun D, Harris DN, Kessler MD, Carlson J, Szpiech ZA, et al. 2021. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature*. 590:290–299.
- Van Oosterhout C, Hutchinson WF, Wills DPM, Shipley P. 2004. MICRO-CHECKER: software for identifying and correcting genotyping errors in microsatellite data. *Mol Ecol Notes*. 4:535–538.
- Wang C, Szpiech ZA, Degnan JH, Jakobsson M, Pemberton TJ, et al. 2010. Comparing spatial maps of human population-genetic variation using Procrustes analysis. *Stat Appl Genet Mol Biol*. 9:13.
- Waples RS. 2015. Testing for Hardy-Weinberg proportions: have we lost the plot? *J Hered*. 106:1–19.
- Weinberg W. 1908. Über den nachweis der vererbung beim menschen. *Jh Ver vaterl Naturk Wurttemb*. 64:369–382.
- Wigginton JE, Cutler DJ, Abecasis GR. 2005. A note on exact tests of Hardy-Weinberg equilibrium. *Am J Hum Genet*. 76:887–893.
- Yang WY, Novembre J, Eskin E, Halperin E. 2012. A model-based approach for analysis of spatial structure in genetic data. *Nat Genet*. 44:725–731.
- Zhang F, Flickinger M, Taliun SAG, In PPGC, Abecasis GR, InPSyght Psychiatric Genetics Consortium, et al. 2020. Ancestry-agnostic estimation of DNA sample contamination from sequence reads. *Genome Res*. 30:185–194.