1 **Hardy Weinberg Exact Test In Large Scale Variant Calling**

2 **Quality Control**

3

4 Zhuoyi Huang[1] (zhuoyi.huang@bcm.edu), Navin Rustagi[1] (rustagi@bcm.edu),

5 Degui Zhi[2] (dzhi@uab.edu), L. Adrienne Cupples[3] (adrienne@bu.edu), Richard

6 Gibbs[1] (agibbs@bcm.edu), Eric Boerwinkle[1,4] (eric.boerwinkle@uth.tmc.edu), Fuli

7 Yu[1*] (fyu@bcm.edu)

8

9 * Corresponding author

10 [1]Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX.

11 [2]Department of Biostatistics, University of Alabama at Birmingham, Birmingham, AL.

12 [3]Department of Biostatistics, Boston University School of Public Health, Boston, MA.

13 [4]Human Genetics Center, University of Texas Health Science Center at Houston, Houston, TX

14

1 **Abstract**

2 Hardy Weinberg Equilibrium (HWE) test is widely used as a quality control measure to detect

3 sequencing artifacts like mismapping, allelic dropout and biases. However, in the high

4 throughput sequencing era, where the sample size is beyond a thousand scale, the utility of

5 HWE test in reducing the false positive rate remains unclear. In this paper, we demonstrate that

6 HWE test has limited power in identifying sequencing artifacts when the variant allele frequency

7 is lower than 1% in a variant call set produced from more than five thousand whole genome

8 sequenced samples from two homogeneous populations. We develop a novel strategy of

9 implementing HWE filtering in which we incorporate site frequency spectrum information and

10 determine the p-value cutoff which optimizes the tradeoff between sensitivity and specificity.

11 The novel strategy is shown to outperform the exact test of HWE with an empirical constant p-

12 value cutoff regardless of the sequencing sample size. We also present best practice

13 recommendations for identifying possible sources of false positives from large sequencing

14 datasets based on an analysis of intrinsic biases in the variant calling process. Our novel

15 strategy of determining the HWE test p-value cutoff and applying the test to the common

16 variants provides a practical approach for the variant level quality controls in the upcoming

17 sequencing projects with tens to hundreds of thousand of samples.

18

1    **KEYWORDS:** Hardy Weinberg Equilibrium; whole genome sequencing; variant calling quality

2    control; false positive rate; sequencing artifacts; site frequency spectrum; QA/QC

3

4    **Introduction**

5    The decreasing cost of sequencing makes it possible to sequence large cohorts with moderate

6    to low coverage at million sample size scale in the coming years. An analysis of variant callsets

7    produced from large sequencing datasets are also presenting an unprecedented resolution at

8    low frequency variants [1]. The increased role of rare variants in common diseases has also

9    been documented [2], and improving the yield of rare variants in low coverage whole genome

10   datasets has become extremely important. To increase the yield of high fidelity variants, quality

11   control of variant call sets has been a crucial step in the past large scale projects [3, 4] . More

12   specifically, low sensitivity in variant site discovery and low specificity due to false calling can

13   reduce the discovery power of novel and rare variants in low coverage sequencing projects with

14   a large sample size and median coverage less than 10-fold. However quality assessment and

15   quality control (QA/QC) steps often rely on the estimation of the site frequency spectrum (e.g.

16   allele frequency and genotype frequency) which in themselves are challenging in the low

17   frequency range due to low sequencing coverage [5, 6].

3

1     As a variant site level quality control method, the exact test of Hardy Weinberg Equilibrium

2     (HWE) has been widely used [7] to filter false positive sites which significantly deviate from the

3     random mating hypothesis [4, 8]. Statistical approaches based on deviations from Hardy

4     Weinberg Equilibrium have been frequently used to ascertain genotyping errors, and population

5     stratification in case control studies. The Hardy Weinberg test has a long history in population

6     genetics [9, 10] but the exact test for Hardy-Weinberg has recently gained prominence as a

7     preferred QA/QC method for Genome Wide Association Studies(GWAS) [8, 11, 12]. The test

8     has been particularly useful in scenarios when highly polymorphic sites are targeted, as in SNP

9     array, across several hundred of samples in the same population. Even though a naïve

10     application of HWE filter can result in incorrect inferences for SNP array data [13], HWE test

11     based filters are extensively used in the past, including on datasets consisting of high coverage

12     (>50x) Exome data [14]. While inbreeding events within a population can result in lower

13     prevalence of heterozygotes, significant excess of heterozygotes in a homogeneous population

14     is usually an indicator of sequencing and calling artifacts, such as sequencing errors,

15     misalignment, low coverage allelic dropout, false calling or false imputation [7, 8].. Empirically

16     choosing a constant HWE test p-value between $10^{-2}$ to $10^{-6}$ based on past large scale projects

17     [4, 14], regardless of the sample size, population ethnicity and the variant allele frequency,

18     makes it hard to compare the quality and integrate the variant call sets from different

4

1    sequencing projects. Many standard population genetics tests for estimating important

2    information like diversity, selection and demographic history are direct applications of site

3    frequency spectrum information. The resolution of discovering low frequency variation in

4    populations is unprecedented due to the increasing capacity of sequencing extremely large

5    cohorts, thereby making it intuitive to incorporate allele frequency information into QA/QC

6    procedure.

7

8    In this paper we present a unified picture of QA/QC procedure in which we systematically

9    determine a HWE p-value cutoff optimizing the power of HWE test in different ranges of allele

10   frequency spectrum for datasets with large sequencing sample size. Firstly, we analyze the

11   performance of HWE test in identifying sequencing artifacts at different allele frequency ranges,

12   in a low coverage sequencing project consisting of several thousands of samples from two

13   homogeneous populations. We specially emphasize on the challenges in identifying sources of

14   error using the HWE test at the rare end of the site frequency spectrum. . Secondly, we propose

15   a novel strategy of applying HWE test in an effective allele frequency range, with a p-value

16   cutoff which optimizes the overall sensitivity and specificity. The novel strategy is shown to

17   outperform the HWE test with an empirical constant p-value cutoff independent of the sample

18   size. Thirdly, based on our experiences with an extremely large sequencing cohort, we present

5

1    guidelines for designing a QA/QC procedure which combines the HWE exact test with site

2    frequency information for increased sensitivity and specificity in variant callsets from future large

3    scale studies.

4

5    **Results**

6    The QA/QC analysis presented in this paper has been carried out on the SNV callset of the

7    CHARGES Freeze 3 whole genome dataset (see Methods for more details). This dataset has

8    3396 European Americans (EuAm) samples and 1901 African Americans (AfAm) samples. The

9    gold standard data consists of SNP array genotype information of a subset of 1850 European

10   American samples within the Charge Freeze 3 dataset. Unless otherwise mentioned explicitly,

11   all the results presented in this paper is for the 1850 EuAm samples in the CHARGE dataset for

12   which there is corresponding SNP Array data.

13   **Hardy Weinberg Equilibrium test has limited power for low allele frequency range ($f$<1%).**

14   For large cohort sizes (>1000), Hardy Weinberg Equilibrium (HWE) test has very limited power

15   in filtering out false positives in the rare end of the frequency spectrum ( allele frequency $f$<1%).

16   Comparing the variants called in the low coverage whole genome sequencing (WGS) in 1850

17   EuAm samples to the SNP array data, the alternate allele frequency of identified false positive

18   sites in the WGS call set follows a bimodal distribution, aggregating at both the common end

6

1    (allele frequency ($f$>10%) and the rare end ($f$<1%) of the allele frequency spectrum (red, yellow

2    and grey bins in Figure 1). By applying the exact test of HWE to both the true positive and false

3    positive sites with the p-value cutoff($p$)>$10^{-4}$, most common false positives are filtered (grey bins

4    in Figure 1) out. Increasing the p-value cutoff does not result in significant reduction of rare false

5    positives any further. Even for p-value cutoff $p$>0.1, a reduction in false positive sites with $f$<1%

6    is not observed (red bins in Figure 1), while more common true positive sites are filtered (the

7    difference between green and blue curves in Figure 1) out. This suggests that the exact test of

8    HWE should be applied to the common variant sites with allele frequency above a critical limit,

9    denoted as $f_c$. However, filtering mechanisms need to model sequencing artifacts to get higher

10   quality variants at the rare end of the frequency spectrum, as is discussed in more detail in the

11   subsequent paragraphs. .

12   Figure 1: Alternate allele frequency distribution of true positive and false positive variant sites in
13   the CHARGE European American samples taking SNP array data as gold standard. A subset of
14   true positive sites with HWE test p-value cutoff $p$>0.1 is also shown (blue curve), while for false
15   positive sites, subsets with $p$>0.1 (red bins) and $p$>$10^{-4}$ (yellow bins) are overlaid on top of all
16   false positive distributions (grey bins).

17

18   **Distribution of HWE p-values for false positive sites in the allele frequency.** The true

19   positive and false positive sites follow different distributions in the alternate allele frequency and

20   HWE p-value ($f$-$p$) diagram (Figure 2). In the common allele frequency range ($f$>5%), where

7

1    HWE test is effective, the p-value of most false positive sites ($p<10^{-4}$) is significantly lower than

2    the p-value of true positive sites ($p>10^{-3}$). In the range $f<1\%$, most false positive sites and true

3    positive sites have indistinguishable p-values in the HWE test (green and red dots with $p=1$ in

4    Figure 2).

5

6    Figure 2: The distribution of true positive sites (green dots) and false positive sites (red circles)

7    in the diagram of alternate allele frequency and HWE test p-value.

8

9    The variant sites on the false positive trajectory in the $f$-$p$ diagram correspond to sites with lower

10    number of homozygous Alt/Alt samples ($N_{aa}\leq 3$ in Figure 2, see also Supplementary Figure 4). In

11    common allele frequency range (alternate allele frequency $f>5\%$), the extreme deficit of

12    homozygous Alt/Alt samples indicates a significant excess of heterozygous Ref/Alt samples,

13    which causes significant deviation from HWE. On the contrary, in the rare alternate low allele

14    frequency range ($f<1\%$), the deficit of homozygotes Alt/Alt is in agreement with HWE and the

15    test is not informative to distinguish false positive from true positive sites, which may both have

16    excess of heterozygotes. However the effect of sample sizes is reflected in the HWE p value

17    plots as the false positive trajectory shifts towards the lower frequency end with increasing sizes

18    (see Supplementary Figure 4). Other sequencing artifacts like mismapping and allelic dropout

1     can also confound the HWE test and their role in developing effective QA/QC strategies is

2     discussed below.

3

4     **The cause of false positive sites.** To study the origin of false positive trajectory in the $f$-$p$

5     diagram, we take the variant sites called in the CHARGE dataset with SNP array genotypes in a

6     subset of shared samples, and compare the alternate allele frequency and HWE p-value ($f$, $p$) of

7     WGS with that of genotypes in the SNP array ($f_G$, $p_G$). We investigated the sites with a

8     significant (twice or greater) excess or deficit of alternate alleles in the $f$-$p$ diagram (Figure 3).

9     The sites with deficient alternate alleles $f_G > 2f$ (Figure 3a) aggregate in the low allele frequency

10    range ($f < 5\%$) with a generally high p-value $p \geq 0.01$ (see also histogram in Supplementary Figure

11    7). These are the sites with one or both of the alternate alleles dropped-out due to low

12    sequencing coverage. On the other hand, sites with excess of alternate alleles, $f > 2f_G$ (Figure

13    3b), fall in two regions. In the common allele frequency range ($f > 5\%$), they are located at the

14    false positive trajectory in $f$-$p$ diagram, deviating clearly from the HWE. The large offset in the $f$-$p$

15    diagram is contributed by possible misalignment and sequencing errors (see Supplementary

16    Figure 9-10), which results in excess of heterozygotes, or contributed by false imputation, which

17    introduces false alternate allele at sites with no coverage. The second region where sites with

18    excess of alternate allele aggregate is the low allele frequency range, where the cause of offset

1    is the same but high p-values in HWE test alone is not effective in separating them from the

2    sites with deficit of alternate alleles (as in Figure 3a).

3

4    Figure 3: The difference of the allele frequency and HWE p-value between the WGS call set

5    (filled points) and the SNP array data (open circles), at variant sites with deficit of alternate

6    alleles (3a) and variants sites with excess of alternate alleles (3b).

7

8    We further investigate the genotyping error at sites with large offsets in $f$-$p$ diagram. The overall

9    genotype concordance at the true positive sites is high, on average (99.957%, 98.063%,

10    98.451%) for three genotypes (Ref/Ref, Ref/Alt, Alt/Alt) (Supplementary Table 1). At sites where

11    variant allele are correctly called but per sample genotyping error causes excess of alternate

12    alleles (0.92% of all true positive sites), the average proportion of genotypes (Ref/Ref, Ref/Alt,

13    Alt/Alt) in SNP array data is (99.887%, 0.090%,0.023%), while in the WGS call set, it is

14    (99.463%, 0.535%,0.002%) (Table 1). At these sites, 90% of the Alt/Alt samples and some

15    Ref/Ref samples are incorrectly genotyped as Ref/Alt. On the other hand, for sites with deficit of

16    alternate alleles (5.14% of all true positive sites), the average proportion of genotypes of both

17    Ref/Alt and Alt/Alt drops from 0.255% and 0.262%, in the SNP array, to 0.110% and 0.006%,

18    respectively, in the low coverage WGS call set (Table 1), as a result of low sequencing

19    coverage.

20

Table 1. The difference in the average genotype proportion of the variant sites between SNP array and WGS call set.

| | All TPs | | | TPs with excess of Alt | | | TPs with deficit of Alt | | |
|---|---|---|---|---|---|---|---|---|---|
| | Ref/Ref | Ref/Alt | Alt/Alt | Ref/Ref | Ref/Alt | Alt/Alt | Ref/Ref | Ref/Alt | Alt/Alt |
| SNP array | 88.34% | 10.08% | 6.57% | 99.89% | 0.09% | 0.02% | 99.48% | 0.26% | 0.26% |
| Called | 83.51% | 9.97% | 6.52% | 99.46% | 0.54% | 0.00% | 99.89% | 0.11% | 0.01% |

1

2     At false positive sites, besides the genotyping error as mentioned above, all samples are

3     genotyped against as alternate allele which is falsely called due to similar reasons. We find that

4     most of the sites with these genotyping errors are rare with $f$<1% (Supplementary Figure 7), is

5     consistent with the excess of false positive sites in the low frequency range where the HWE

6     filtering power is limited.

7

8     **The strategy to determine optimal p-value cutoff.** Within the allele frequency range where

9     the HWE test is informative, we can determine the optimal p-value cutoff. We take the difference

10     between the true positive rate (TPR) and false positive rate (FPR) as $\delta$=TPR-FPR. It is an

11     evaluator of the HWE filtering performance, similar to the receiver operating characteristic

12     (ROC) analysis. Since the 1+$\delta$ is the sum of sensitivity and specificity, we define the cutoff

13     which maximizes $\delta$ as the optimal tradeoff between the two. Applying HWE to the full allele

11

1    frequency spectrum only gives maximum $\delta=0.08$ at p-value cutoff $p>10^{-3}$ (Figure 4 red dash

2    line). Since the p-value of true positive and false positive sites is as high as 1 below certain

3    allele frequency, we can choose it as a limiting allele frequency $f_c$, and only apply HWE test to

4    the common sites with alternate allele frequency $f>f_c$ (see Methods). For instance, when the

5    sample size is 1850, and the number of homozygotes $N_{aa}=0$, the limiting allele frequency is

6    $f_c=0.0239$. Above this allele frequency, the optimal p-value cutoff gives $\delta=0.8$ (Figure 4 blue

7    dash line), which is increased by 10-fold compared to a constant p-value applied to all allele

8    frequency spectrum. The filtering performance is significantly improved due to two factors. First,

9    when $f > f_c$, the true positives and false positives follow very different p-value distribution as is

10   mentioned above. Second, towards the low allele frequency end, where the HWE test is non-

11   informative, the number of variant sites is exponentially increasing, therefore excluding the low

12   frequency range in the test significantly increases the performance. With the optimal p-value

13   cutoff, the overall sensitivity, specificity and FDR of the WGS call set is 85.709%, 97.727% and

14   2.388%, respectively, taking SNP array data as the gold standard (Supplementary Table 1).

15

16   Figure 4: The difference between true positive rate (TPR) and false positive rate (FPR) as a

17   function of HWE p-value cutoff. The cutoff is applied in the allele frequency range where the

18   HWE test is effective (blue curve) and to the full allele frequency spectrum (red curve).

19

1    **Best practices of applying Hardy Weinberg Equilibrium test in the large scale variant**

2    **calling quality control.** In Text Box 1, we summarize the recommended site level quality

3    control (QC) procedures to apply the HWE test in the variant calling with extremely large sample

4    size.

---

**Text Box 1**

**Best practice recommendation for the QC procedure:**

- Apply joint calling and genotyping across all samples (see Methods). (For variant callers which already make use of HWE information, we recommend keeping all sites and HWE p-value in the annotation as input of QC);

- Subset samples by population and in each population, apply the exact test of Hardy Weinberg Equilibrium at each variant site;

- Take a subset of samples from each population with shared samples in gold standard dataset (e.g. SNP array), recalculate HWE p-value for the subset of samples in each population, and identify true positive and false positive sites (see Discussion for the choice of golden data sets, and for cases when golden data set is not available);

- Based on the sample size in population $i,$ find the limit allele frequency $f^i_c$ where HWE is informative at $f > f^i_c$ (see Methods);

- Find the optimal p-value cutoff $p^i_c$ which maximizes TPR-FPR at $f^i > f^i_c$;

- Apply HWE p-value cutoff and filter variant sites with p-value $p^i < p^i_c$ at $f^i > f^i_c$ in the full call set in population $i$.

---

13

1

## Discussion

3    The effectiveness of HWE test for detecting sequencing artifacts in large sequencing cohorts is

4    undocumented in previous literature, because of limitations on samples sizes. In this paper, we

5    took a large cohort of whole genome sequencing samples with 6x-8x low coverage and

6    revealed the limited power of HWE test at the rare end of the site frequency spectrum. As we

7    showed as part of our results, the HWE test is very sensitive to the excess or deficit of

8    heterozygotes samples at common false positive sites, but at the rare end ($f$<1%), false positive

9    sites with significant excess of heterozygotes samples cannot be distinguished from the true

10    positives using the test. As was discussed in Results section, misalignment in local repetitive

11    regions, false imputation and sequencing error, all plays a role in genotyping errors at the lower

12    end of the frequency spectrum. We present a novel strategy to determine the cutoff on HWE

13    test p-value, which effectively filters the false positive sites in the allele frequency range where

14    HWE test is informative and achieves a better tradeoff between the sensitivity and specificity,

15    compared to the filtering result with a constant p-value cutoff for all allele frequencies.

16

17    We recommend using high quality genotype data set as the quality control gold standard, which

18    includes large number of samples commonly shared between the WGS call set. When selecting

14

1    the golden data sets, we recommended considering following three factors. Firstly, the number

2    of common samples in the golden data set should be large enough to avoid under-sampling of

3    the rare false positive sites, and the samples should be within the same population so that the

4    HWE p-value is not biased by any population stratification or significant population

5    substructures. For instance, when the variant calling sample size is 10,000 and the HWE test

6    limiting allele frequency is about 1% (Supplemental Figure 4), at least 500 samples in the same

7    population should be present in the golden data set in order to assess the site quality with allele

8    frequency at least an order of magnitude lower than the limiting allele frequency. Secondly, to

9    assess the calling quality across the whole genome in a uniform manner, the golden data sets

10   should adequately cover both the coding and non-coding regions, both the highly polymorphic

11   sites and the rare variant sites. Using either the SNP array data containing only the common

12   variants or the high coverage target sequencing call set (e.g. whole Exome sequencing, or

13   WES) may not be sufficient. When large number of samples are available in quality control,

14   combining different control data sets may provide balanced genomic coverage. Thirdly, the

15   golden data set should include only the sites with high quality genotypes. The allele frequency

16   dependent genotyping error in the golden data set will introduce noise in the quality

17   assessment, and consequently bias the determination of HWE p-value cutoff. In case of SNP

18   array data, using a consensus of different arrays can effectively avoid the array specific

15

1    genotyping bias. In particular, any site with missing genotype in golden data set should be

2    excluded from the negative controls, while sites with large fraction (e.g. 5%) of samples missing

3    genotype, should be excluded from the positive controls. If the high coverage WES variant

4    calling result is chosen as the golden data set, stringent site level genotype quality filter must be

5    first applied. This is because in WES variant calling, the allele frequency dependency of

6    genotyping error is not negligible and rare false positives may still be present due to the

7    misalignment (see Supplemental Figure 3,6-7).

8

9    If the golden data set with shared samples is not available, we still recommend estimating the

10    limit allele frequency $f_c$ according to the sample size, and applying the HWE test to the sites with

11    allele frequency $f > f_c$. In this case, an empirical p-value cutoff needs to be determined from the

12    optimal p-value cutoff with similar sample size in other large scale projects, but the tradeoff

13    between sensitivity and specificity may not be optimal.

14

15    Since the HWE test has limited power in filtering the rare false positive variant sites, additional

16    orthogonal filters may be used to improve the quality control in the rare end of the site frequency

17    spectrum. For example, false positives due to local repetitive sequences can be filtered using

18    low complexity regions (LCR) filter [10], while the sites with higher switching errors in the

16

1    imputation and phasing may cluster in recombination hot spots regions in the genome. Since

2    these regional or loci-specific filters are generally independent of the allele frequency of variant

3    sites, they provide an orthogonal diagnosis to identify part of the rare false positive sites.

4    Furthermore, per sample coverage information, when available, may also provide independent

5    information to filter rare false positives with insufficient read depth at sample level. Improving the

6    genotyping accuracy in the low coverage variant calling by making use of position specific

7    genotype likelihood prior in a known population also helps to reduce the false positive rate.

8    Besides, sample level quality control can be useful to identify rare false positives contributed by

9    some outlying samples with extremely low or high coverage, or sequencing mishaps. Ultimately,

10   given the sample size, increasing the sequencing coverage is still the most effective way to

11   reduce the false positive calls in the rare end of site frequency spectrum.

12

13   **Conclusions**

14   As the sequencing sample size continues to grow in the high throughput sequencing era, more

15   rare variants are being discovered and quality control procedure to effectively remove the rare

16   false positives is becoming increasingly important. Using the whole genome sequencing variant

17   call set with more than 5000 samples, we revealed the limited power of exact test of HWE in

18   filtering the rare false positive variant sites with allele frequency less than 1%. By analyzing the

1    genotype discordance, we discussed the different causes of false calling and genotyping error in

2    the low coverage sequencing. Given the availability of gold standard genotype data, we propose

3    a novel strategy to ascertain the allele frequency range where the HWE test is informative and

4    to determine the p-value cutoff which optimizes the overall sensitivity and specificity in the entire

5    allele frequency spectrum. This strategy outperforms the conventional HWE filtering approaches

6    with only an empirical constant p-value cutoff regardless the sequencing sample size. This can

7    serve as a practical consideration for the variant level quality controls in the upcoming

8    sequencing projects with tens or hundreds of thousand of samples.

9

10   **Materials and Methods**

11   **Sequencing samples and variant calling.** We took 5297 whole genome sequenced samples

12   from the CHARGE project (Cohort for Heart and Aging Research in Genomic Epidemiology

13   [15]). 3396 samples are European Americans (EuAm) and 1901 African Americans (AfAm). The

14   samples were sequenced using Illumina HiSeq 2000 with an average coverage 6x-8x. The

15   alignment was done using BWA [16] integrated in the Mercury pipeline [17]. We called biallelic

16   SNPs across all 5297 samples, taking an ensemble variant calling approach goSNAP [18],

17   which employs four variant callers GATK-HaplotypeCaller [19, 20] with gVCF option, GATK-

18   UnifiedGenotyper [19, 20], SNPTools [21] and GotCloud [22], each enforced in a joint calling

1    mode. To ensure a high quality variant call set, we applied a consensus filtering and selected

2    72,945,834 variant sites which were called at least in 3 out of all 4 callers. The genotype

3    likelihood of each sample at each variant site was calculated using BAM specific Binomial

4    Mixture Model (BBMM) algorithm implemented in SNPTools [21]. Imputation and phasing was

5    done using SNPTools reference panel independent imputation engine. After imputation, we

6    obtain phased genotypes of 5297 samples at 72,762,406 biallelic variant sites (52,116,900 in

7    AfAm samples and 46,201,314 in EuAm samples). The site frequency spectrum of both

8    populations follows similar patterns as in 1000 Genomes Project [4] (Supplementary Figure 1).

9    We use the call set of both populations to evaluate the performance of Hardy Weinberg

10   Equilibrium test.

11

12   **Golden data set for quality assessment.** We use CHARGE SNP array data as the golden

13   data set to assess the low coverage whole genome variant calling quality. Among 5297

14   CHARGE WGS sequencing samples, there are 3533 samples with SNP array genotype, 1683

15   EuAm samples and 1850 AfAm samples. The total number of autosomal variant sites in SNP

16   array data is 228,963 across all 3533 samples. To ensure the quality of control data set, we

17   further remove the control sites according to following criteria. In each population, we remove

18   the sites called in WGS samples but have more than 5% of samples missing genotype in the

1    array. We also remove sites which are not called in WGS samples and have any sample with

2    missing genotype in the array. After filtering, the number of positive and negative control sites is

3    107,339 and 99,034, respectively, in EuAm samples (132,344 and 75,163 in AfAm samples).

4

5    Alternatively, we compare the WGS variant calling quality to the WES variant call set, which

6    includes 4612 samples in common with CHARGE WGS, 1782 AfAm samples and 2830 EuAm

7    samples. The sequencing coverage of WES is 80-100x. The number of control sites in AfAm

8    and EuAm samples is in Supplemental Table 1.

9

10   **The exact test of Hardy Weinberg Equilibrium.** We apply the exact test of Hardy Weinberg

11   Equilibrium [6] to 52,116,900 variant sites called in 1901 AfAm WGS samples, and 46,201,314

12   variant sites called in 3396 EuAm samples. The exact HWE test is also applied to the filtered

13   control sites in a subset of 1850 AfAm samples and 1683 EuAm samples, which are shared

14   between the WGS and SNP array. Besides the p-value of exact HWE test, in each subset of

15   samples and at each biallelic variant site, we also extract the number of samples with three

16   genotypes, homozygous Ref/Ref ($N_{rr}$), heterozygous Ref/Alt ($N_{ra}$) and homozygous Alt/Alt ($N_{aa}$).

17   The alternate allele frequency is derived as $f = (N_{ra}+2N_{aa})/((N_{rr}+N_{ra}+N_{aa})\times2)$.

18

20

1     **Calculate the limiting allele frequency of HWE test.** To determine the p-value cutoff in HWE

2     test, we first calculate the limiting allele frequency $f_c$, above which the HWE is effective, based

3     on the number of samples $N$ in the same population. We keep the number of Alt/Alt samples

4     $N_{aa}=0$, and increase the number of heterozygous samples $N_{ra}$ from 0, and $N_{rr} = N - N_{ra}$. For each

5     configuration of $(N_{rr}, N_{ra}, N_{aa})$, we apply the exact test of HWE. When $N_{ra}$ is small, the HWE p-

6     value is always $p=1$. The HWE test is informative when the p-value starts to be less than 1,

7     where we take number of heterozygotes $N'_{ra}$ of this configuration and calculate the limiting allele

8     frequency as $f_c = N'_{ra} / 2N$. We apply the HWE test only to variant sites with alternate allele

9     frequency $f > f_c$.

10

11     **Determine optimal p-value cutoff in the HWE test.** By comparing to the golden data set

12     among a subset of shared samples, we classify variant sites in WGS call set as true positive

13     (shared in golden data set) and false positive sites (unique in WGS call set). For $T$ true positive

14     sites and $F$ false positive sites with alternate allele frequency (derived based on genotypes in

15     the WGS call set) higher than the limit allele frequency ($f > f_c$), we apply the exact test of HWE

16     and calculate the p-value. By changing p-value cutoff $p_c$, we obtain $T'$ ($T' < T$) true positive sites

17     and $F'$ ($F' < F$) false positive sites with $p > p_c$. We calculate the true positive rate $TPR(p_c) = T'/T$ and

1    false positive rate FPR($p_c$)=$F^7/F$. The optimal HWE p-value cutoff is determined when TPR-FPR

2    is maximized (Figure 4).

3    **Data**

4    CHARGE dataset is accessible on dbGap with following study accession numbers, FHS:

5    phs000651, ARIC: phs000668, CHS: phs000667. The permission is required to access the

6    data. Participant consent is not required to access the data.

7

8    **References**

9

10   1.      Lek M, Karczewski K, Minikel E, Samocha K, Banks E, Fennell T, et al. Analysis of

11   protein-coding genetic variation in 60,706 humans. BioRxiv. 2016:030338.

12   2.      Lupski JR, Belmont JW, Boerwinkle E, Gibbs RA. Clan genomics and the complex

13   architecture of human disease. Cell. 2011;147(1):32-43.

14   3.      Genomes Project C. A map of human genome variation from population-scale

15   sequencing. Nature. 2010;467(7319):1061-73.

16   4.      Consortium GP. A global reference for human genetic variation. Nature.

17   2015;526(7571):68-74.

1  5.      Vieira FG, Fumagalli M, Albrechtsen A, Nielsen R. Estimating inbreeding coefficients

2  from NGS data: impact on genotype calling and allele frequency estimation. Genome

3  research. 2013;23(11):1852-61.

4  6.      Kim SY, Lohmueller KE, Albrechtsen A, Li Y, Korneliussen T, Tian G, et al. Estimation

5  of allele frequency and association mapping using next-generation sequencing data. BMC

6  bioinformatics. 2011;12(1):1.

7  7.      Nielsen R, Paul JS, Albrechtsen A, Song YS. Genotype and SNP calling from next-

8  generation sequencing data. Nature Reviews Genetics. 2011;12(6):443-51.

9  8.      Wigginton JE, Cutler DJ, Abecasis GR. A note on exact tests of Hardy-Weinberg

10  equilibrium. The American Journal of Human Genetics. 2005;76(5):887-93.

11  9.      Haldane JBS. An exact test for randomness of mating. Journal of Genetics.

12  1954;52(3):631-5.

13  10.     Levene H. On a matching problem arising in genetics. The annals of mathematical

14  statistics. 1949:91-4.

15  11.     Rohlfs RV, Weir BS. Distributions of Hardy–Weinberg equilibrium test statistics.

16  Genetics. 2008;180(3):1609-16.

17  12.     Graffelman J, Moreno V. The mid p-value in exact tests for Hardy-Weinberg

18  equilibrium. Statistical Applications in Genetics and Molecular Biology. 2013;12(4):433-48.

1    13.    Yong Zou G, Donner A. The merits of testing Hardy‐Weinberg equilibrium in the

2    analysis of unmatched case‐control data: A cautionary note. Annals of human genetics.

3    2006;70(6):923-33.

4    14.    Lohmueller KE, Sparsø T, Li Q, Andersson E, Korneliussen T, Albrechtsen A, et al.

5    Whole-exome sequencing of 2,000 Danish individuals and the role of rare coding variants

6    in type 2 diabetes. The American Journal of Human Genetics. 2013;93(6):1072-86.

7    15.    Psaty BM, O'Donnell CJ, Gudnason V, Lunetta KL, Folsom AR, Rotter JI, et al. Cohorts

8    for heart and aging research in genomic epidemiology (CHARGE) consortium design of

9    prospective meta-analyses of genome-wide association studies from 5 cohorts. Circulation:

10    Cardiovascular Genetics. 2009;2(1):73-80.

11    16.    Li H, Durbin R. Fast and accurate long-read alignment with Burrows–Wheeler

12    transform. Bioinformatics. 2010;26(5):589-95.

13    17.    Reid JG, Carroll A, Veeraraghavan N, Dahdouli M, Sundquist A, English A, et al.

14    Launching genomics into the cloud: deployment of Mercury, a next generation sequence

15    analysis pipeline. BMC bioinformatics. 2014;15(1):1.

16    18.    Huang Z, Rustagi N, Veeraraghavan N, Carroll A, Gibbs R, Boerwinkle E, et al. A

17    hybrid computational strategy to address WGS variant analysis in >5000 samples. BMC

18    Bioinformatics. 2016;17(1):1-12. doi: 10.1186/s12859-016-1211-6.

1   19.     McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The

2   Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA

3   sequencing data. Genome research. 2010;20(9):1297-303.

4   20.     Auwera GA, Carneiro MO, Hartl C, Poplin R, del Angel G, Levy‐Moonshine A, et al.

5   From FastQ data to high‐confidence variant calls: the genome analysis toolkit best

6   practices pipeline. Current protocols in bioinformatics. 2013:11-0.

7   21.     Wang Y, Lu J, Yu J, Gibbs RA, Yu F. An integrative variant analysis pipeline for

8   accurate genotype/haplotype inference in population NGS data. Genome research.

9   2013;23(5):833-42.

10  22.     Jun G, Wing MK, Abecasis GR, Kang HM. An efficient and scalable analysis framework

11  for variant extraction and refinement from population-scale DNA sequence data. Genome

12  Res. 2015;25. doi: 10.1101/gr.176552.114.

13

14

15  **Additional Files**

16  Additional hwe.submit.supp.docs

17