

# Non-replication of association studies: “pseudo-failures” to replicate?

Prakash Gorroochurn, PhD<sup>1</sup>, Susan E. Hodge, DSc<sup>1,3</sup>, Gary A. Heiman, PhD<sup>2</sup>, Martina Durner, MD<sup>1</sup>, and David A. Greenberg, PhD<sup>1,3</sup>

Recently, serious doubts have been cast on the usefulness of association studies as a means to genetically dissect complex diseases because most initial findings fail to replicate in subsequent studies. The reasons usually invoked are population stratification, genetic heterogeneity, and inflated Type I errors. In this article, we argue that, even when these problems are addressed, the scientific community usually has unreasonably high expectations on replication success, based on initial low *P* values, a phenomenon known as the replication fallacy. We present a modified formula that gives the replication power of a second association study based on the *P* value of an initial study. When both studies have similar sample sizes, this formula shows that: (1) a *P* value only slightly lower than the nominal  $\alpha$  results in only approximately 50% replication power; (2) very low *P* values are required to achieve a replication power of at least 80% (e.g., at  $\alpha = 0.05$ , a *P* value of  $<0.005$  is required). Because many initially significant findings result in low replication power, replication failure should not be surprising or be interpreted as necessarily refuting the initial findings. We refer to replication failures for which the replication power is low as “pseudo-failures.” ***Genet Med* 2007;9(6):325–331.**

**Key Words:** Complex diseases, replication fallacy, replication probability, effect size

## INTRODUCTION

Despite much initial optimism, genetic association studies have been far from entirely successful. A decade ago, Risch and Merikangas<sup>1</sup> argued that association analysis could provide a more statistically powerful framework to dissect complex diseases than linkage analysis under certain circumstances. This fact, compounded by an avalanche of SNPs that later became steadily available, triggered great enthusiasm in association studies. Several association studies have been performed and published in journals spanning a wide range of interdisciplinary interests. However, the results have been less than impressive: most of these studies seem to have failed one of the benchmarks of scientific success, namely replicability.<sup>2</sup> A review article by Hirschhorn et al.<sup>3</sup> found that, of more than 600 associations previously reported, 166 were studied at least thrice; of these, only six were consistently replicated. Thus, the reality today seems to be one of skepticism, if not downright pessimism.<sup>4–9</sup> What has gone wrong? Should we eventually abandon association studies? More precisely, is a 3.6% replication rate really bad? In this article, we show that the concept of

replication has been largely misunderstood and that the genetics community has unintentionally put overly high replication expectations on initially significant findings (even when these are quite significant). In brief, we contend that a 3.6% replication rate, when looked at in the right statistical light, is really not surprising and should not give way to disillusionment.

## MAPPING BY POPULATION ASSOCIATION

A genetic association exists between a phenotype (e.g., disease) and a specific allele at a genetic marker if the allele occurs more often (than would be expected by chance) in a group of individuals with the disease (cases) compared with a group without the disease (controls). Whereas linkage analysis is concerned with the co-segregation of marker loci with disease within families, association analysis studies the dependence between marker alleles and disease at a population level (see, for example, Hodge<sup>10</sup>).

For association mapping by linkage disequilibrium to be successful, two conditions are desirable.<sup>11</sup> First, the disease allele must have arisen only once in the population so that there is complete linkage disequilibrium between the marker and disease alleles (hence small isolated populations are often preferred). Second, the marker and disease loci must be in very close physical proximity, so that the disequilibrium between the marker and disease alleles is the result of tight linkage (typically the recombination fraction  $\theta \approx 10^{-5}$ ). Problems in linkage disequilibrium mapping arise when the disequilibrium in question is instead the result of population history, natural selection, or population stratification. In all these cases, an as-

From the <sup>1</sup>Division of Statistical Genetics, Department of Biostatistics, and <sup>2</sup>Department of Epidemiology, Mailman School of Public Health, Columbia University; <sup>3</sup>Clinical-Genetic Epidemiology Unit, New York State Psychiatric Institute, New York, New York.

P. Gorroochurn, Columbia University, Department of Biostatistics, Rm 620, 722 W 168th Street, New York, NY 10032. E-mail: pg2113@columbia.edu

Disclosure: The authors declare no conflict of interest.

Submitted for publication February 23, 2007.

Accepted for publication March 30, 2007.

DOI: 10.1097/GIM.0b013e3180676d79

sociation between marker allele and disease can be observed, but the disease locus is not necessarily close to the marker locus (in the cases of population history or selection) or is even non-existent (in the case of population stratification).

Population stratification can be especially deleterious to association studies<sup>12,13</sup>, and much effort has recently been channeled to further understand and correct for this confounder.<sup>14–20</sup> If population stratification was to be corrected for and other design issues were to be improved in future association studies, should we expect a sudden leap in successful replication rates? We explain why such an expectation will likely remain unfulfilled in years to come. The main problem lies in a common misunderstanding of the meaning of a replication and in how likely it is to occur. Furthermore, we will show that many of the apparent failures to replicate have, in fact, been “pseudo-failures.”

**REPLICATION FALLACY, REPLICATION PROBABILITY, AND P VALUES**

We use “replication” to refer to a situation in which a null hypothesis that has been rejected at Time 1 is rejected again, and with the same direction of outcome, on the basis of a new study at Time 2<sup>21</sup> (Table 1).

If a test of association is rejected with  $P = 0.01$ , what is the probability the study will be replicated? Oakes<sup>22</sup> reports that, in a survey of 70 academic psychologists each with at least 2 years’ research experience, 60% endorsed the statement that if an initial study is significant with  $P = 0.01$ , then “You have a reliable experimental finding in the sense that if, hypothetically, the experiment were repeated a great many times, you would obtain a significant result on 99% of occasions.” Nothing could be further from the truth. In point of fact, as we will soon show, when  $\alpha = 0.05$ , an initial study with  $P = 0.05$  implies a replication probability of only 0.5, one with  $P = 0.02$  implies only approximately 0.64 replication probability, and one with  $P = 0.01$  results in no more than approximately 0.73 probability of replication! These results: (1) hold irrespective of the sample size used as long as the sample sizes in the initial and subsequent studies are equal (the case of unequal sample sizes will be discussed later); (2) assume a  $z$  test (we discuss

deviations from this assumption below); and (3) are based on the concept of *power*. The incorrect belief that  $1 - p$  is the probability of replication of an initial study is known as the *replication fallacy*.<sup>23,24</sup> The replication fallacy is a reason why the scientific community generally has overly high expectations that an initial study that yields a relatively low  $P$  value should be replicated, and if that does not happen, then surely the initial finding must have been false. Nothing, again, could be further from the truth. The whole issue is succinctly described by Oakes<sup>22</sup>:

“We have seen that the power of a replication of an independent means  $t$ -test design when the first experiment has an associated probability of 0.02 is approximately 0.67 . . . Suppose then, psychologist B, suspecting that A’s results were an artefact . . . decided to perform an exact replication of A’s study. Suppose B’s results were in the same direction as A’s but were not significant. It would be folly surely for B to assert that A’s findings were indeed artefactual.”

It is precisely this kind of apparent replication failure that we describe as a pseudo-failure because the probability of a successful replication is a priori quite low, and a subsequent “failure” should come as neither a surprise nor a contradiction of the initial positive finding.

If the  $P$  value is not a direct measure of the replicability of an initial study, is there anything at all it can tell about the likelihood of replication? The short answer is yes. Although neither the exact  $P$  value nor its complement can be interpreted as the probability of a successful replication, it has been shown that, for a given test size  $\alpha$ , “the replicability of null hypothesis rejection is a continuous, increasing function of the complement of its  $P$  value”.<sup>25</sup> Although this fact is well known in the social and behavioral sciences<sup>22,24–27</sup>, it seems to have been overlooked by the genetics community, and its implications have not before been explored as we do in this article.

Let us now consider an initial association study based on a simple  $\chi^2$  test (Table 2), which yields a  $P$  value  $p_1$  for a test size  $\alpha$ , where  $p_1 \leq \alpha$ . It can be shown that the initial finding, assuming it is not a false positive, will be replicated with probability (see Appendix)

$$p_{REP} = \Phi \left\{ \frac{z_1 \sqrt{n_2/n_1} - z_{crit}}{1 - z_1^2/(2n_1)} \right\} \quad (1)$$

where  $n_1, n_2$  are the sample sizes (i.e., the number of cases or controls) in the first and second studies, respectively,

**Table 1**

The meaning of a replication

	Second study	
	$P_2 \leq \alpha^a$	$P_2 > \alpha$
First Study		
$P_1 \leq \alpha$	Successful replication	Failure to replicate
$P_1 > \alpha^b$	Failure to replicate	Failure to establish effect

Adapted from Rosenthal.<sup>21</sup>

$P_1$ ,  $P$  value of first study;  $P_2$ ,  $P$  value of second study;  $\alpha$ , test size for both studies.

<sup>a</sup>For a true replication, the effect in the second study must be in the same direction as the first.

<sup>b</sup>A second study might be contemplated even when the first study was not significant if (1)  $P_1$  was only slightly larger than  $\alpha$ ; (2) the nominal  $\alpha$  used was extremely small; or (3) there was sufficient biological justification.

**Table 2**

Contingency table and  $\chi^2$  test statistics for first and second studies

	$i^{\text{th}}$ study ( $i = 1,2$ )		
	Genotype with at least one marker allele	Genotype with no marker allele	Row total
Disease	$a_i$	$b_i$	$n_i$
Non-disease	$c_i$	$d_i$	$n_i$

$$\chi^2_i = \frac{2(a_i d_i - b_i c_i)^2}{n_i(a_i + c_i)(b_i + d_i)}$$

$z_{crit} = \Phi^{-1}(1-\alpha/2)$ ,  $z_1 = \Phi^{-1}(1-p_1/2)$ , and  $\Phi$ ,  $\Phi^{-1}$  are the standard normal distribution function and inverse distribution function, respectively. Note that, in accordance with our definition of a replication, Equation 1 will be used only when: (1) the initial test is significant (i.e.,  $p_1 \leq \alpha$ ); and (2) the outcomes of both tests are in the same direction (i.e.,  $a_1 - c_1$  and  $a_2 - c_2$  in Table 2 have the same signs). Moreover, although Equation 1 assumes an equal number of cases and controls in the first study, and similarly equal numbers in the second study, situations with different numbers of cases and controls can be dealt with through the use of *harmonic means*. More specifically, any study containing  $r$  cases and  $s$  controls (where  $r \neq s$ ) can be treated as one with the same number of cases and controls where  $t = 2rs/(r + s)$ .<sup>28</sup> Equation 1 implicitly assumes the same test size  $\alpha$  for both the first and subsequent studies: for unequal test sizes, say  $\alpha^*$  and  $\alpha$ , respectively, the same equation can be used as long as  $p_1 \leq \alpha^*$ . Finally, we also note that the whole issue of replicability can also be approached from a Bayesian perspective.<sup>29</sup>

The critical assumption made in Equation 1 is that the expected effect size in the second study equals the observed effect size in the initial study (Greenwald et al.<sup>25</sup>). Therefore, the replication probability calculated is, in effect, the conditional power of the second study to be statistically significant given this assumption about effect size (hence, “replication probability” and “replication power” are used interchangeably from now on). Regarding this assumption, Greenwald et al.<sup>25</sup> state that it “involves a step of inductive reasoning that is (a) well recognized to lack rigorous logical foundation but is (b) nevertheless essential to ordinary scientific activity.” Of more pertinence is the fact that, when calculating replication probabilities, the literature has usually assumed the second study is an exact repetition of the first study<sup>25,27</sup>, in the sense that both studies: (1) have the same sample size; and (2) are performed using the same population. Clearly, this is the norm neither in association studies nor elsewhere, especially the second condition. Regarding the first condition, all replication calculations are based on the assumption of constancy of effect size, as we explained at the start of the paragraph. Now, because effect size is unaffected by sample size<sup>24</sup>, there is enough justification to use the effect size in a study with a different sample size. However, it is very desirable for the first study to have a reasonably large sample size so that the standard error of its observed effect size is relatively small. For examples on the use of different sample sizes in replication calculations, based on initial effect sizes, see Tversky and Kahneman<sup>30</sup> and Heils.<sup>31</sup> Regarding the second condition, Tan et al.<sup>32</sup> point out, “Estimates of effect size will tend to regress to the true effect size in subsequent [association] studies, which is usually less extreme.” The same feeling has been echoed elsewhere,<sup>7,8,33</sup> and Göring et al.<sup>5</sup> have provided a theoretical justification. Thus, the expected effect size of the second study will, on average, be smaller than the observed effect size of the first study. What this means is that, in general, even when a different population is used for the subsequent study, Equation 1 can still be used, with the understanding that the replication probabilities calculated from that

equation actually represent *upper bounds* for the true replication probabilities.<sup>25</sup>

Let us now examine three of the more unexpected consequences of Equation 1, when both the first and second studies have similar sample sizes (i.e.,  $n_2 \approx n_1$ ):

**Consequence 1**

A  $P$  value only slightly less than the nominal  $\alpha$  in the first study (e.g.,  $P = 0.04$  at  $\alpha = 0.05$ ) results in a replication power of only approximately 50% for the second study (see Fig. 1).

**Consequence 2**

Reasonably low  $P$  values in the first study do not necessarily result in high replication power of the second study (e.g.,  $P = 0.02$  at  $\alpha = 0.05$  implies a replication power of no more than 64%) (see Fig. 1).

**Consequence 3**

To achieve a replication power of 80% for the second study at  $\alpha = 0.05$ , a  $P$  value of at most 0.005 must be obtained in the first study (see Fig. 1).

When  $n_1$  and  $n_2$  are allowed to be different, two more consequences are

**Consequence 4**

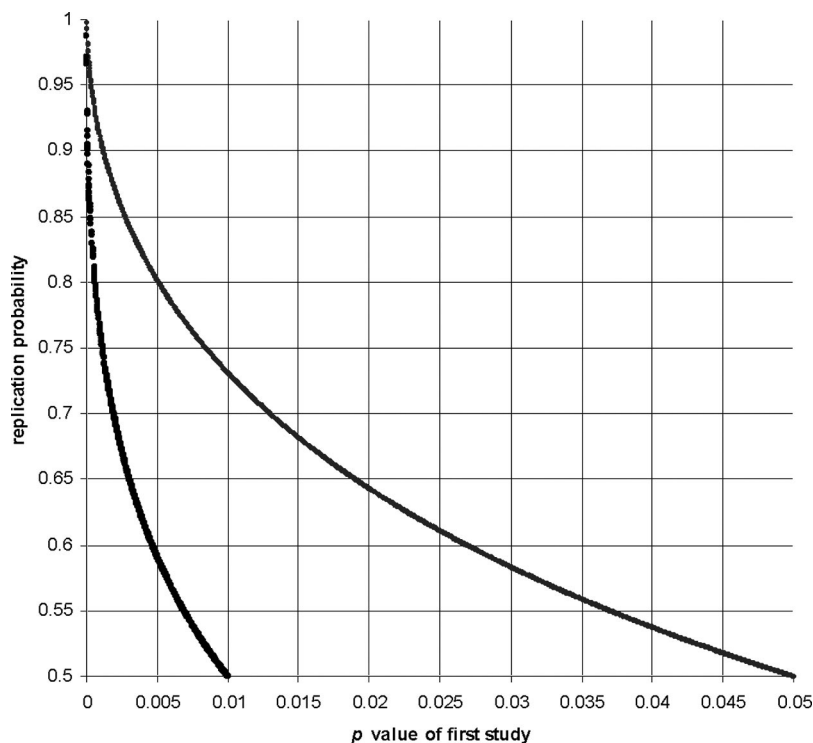
For reasonably large sample sizes, the replication power depends only on the *ratio* of the initial and subsequent sample sizes, not on their absolute values (e.g., if  $P = 0.02$  at  $\alpha = 0.05$ , then  $p_{REP} = .727$  when  $n_1 = 100$  and  $n_2 = 120$ , and  $p_{REP} = .723$  when  $n_1 = 500$  and  $n_2 = 600$ ).

**Consequence 5**

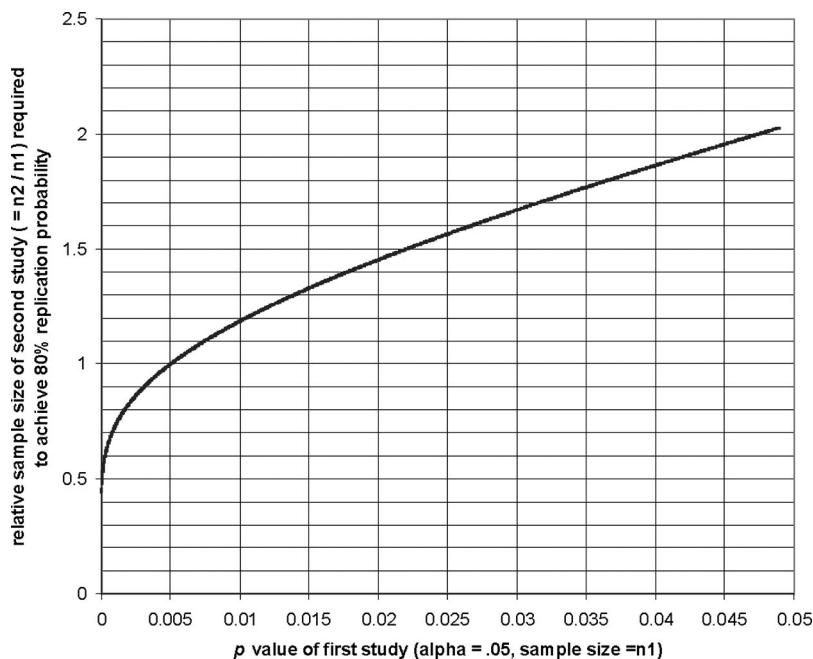
For a given sample size of the first study, if the initial  $P$  value is only slightly less than the nominal  $\alpha$ , the sample size required for the second study must be much larger to achieve a replication power of 80%. Suppose an initial association study with  $n_1$  cases and  $n_1$  controls yields a  $P$  value  $p_1$ . Then, the number of cases (or controls) required for the second study to achieve a replication power of  $p_{REP}$  follows directly from Equation 1

$$n_2 \approx n_1 \left\{ \frac{z_{crit} + \Phi^{-1}(p_{REP})}{z_1} \right\}^2 \tag{2}$$

where, as before,  $z_{crit} = \Phi^{-1}(1-\alpha/2)$ ,  $z_1 = \Phi^{-1}(1-p_1/2)$ ,  $\Phi^{-1}$  is the inverse standard normal distribution function, and we have assumed  $z_1^2 \ll 2n_1$ . Figure 2 illustrates the variation of  $n_2/n_1$  with  $p_1$ . For example, if the initial study has  $P = 0.04$  at  $\alpha = .05$ , then the sample size of the second study must be approximately 1.86 times that of the first study for a replication power of 80%. (Because the expected effect size of the second study will, on average, be smaller than the observed effect size of the first study, as we explained earlier, the actual required relative sample size will be *greater* than 1.86.)



**Fig. 1.** Variation of replication probability as a function of the  $P$  value of the  $\chi^2$  test in the first study. The first and second studies are assumed to have the same sample size. For the top curve,  $\alpha = 0.05$ ; for the bottom curve,  $\alpha = 0.01$ . For other values of  $\alpha$ , use Equation 1.



**Fig. 2.** Variation of the relative sample size (number of cases or controls) required for the second study as a function of the  $P$  value of the  $\chi^2$  test in the first study to achieve a replication probability of 80%.

### IMPLICATION FOR ASSOCIATION STUDIES AND DISCUSSION

Initial association studies that achieve borderline significance, and even those that result in relatively low  $P$  values, result in low replication probability (or power) for a subsequent study. Therefore, failure to replicate such an initial find-

ing in a subsequent study should neither come as a surprise nor be deemed “troubling;” nor is this failure to replicate necessarily an outright refutation of the initial finding.

What should one do with an initial study that yields a  $P$  value of, say, 0.03? We leave it to the reader to decide, but consider this: if roughly the same sample size is used for both the initial

and subsequent studies, the replication power is <58%, and no clinical trial with a power of 58% would ever pass a review board. If one insists on replicating, one should do so with the understanding that there are big chances of not succeeding. There is no denying that consistent replication and subsequent biological confirmation should be the gold standard. However, replication of some initial findings just might not be within the realm of high probability.

We have, however, seen that an association study that yields a  $P$  value of 0.005 at  $\alpha = 0.05$  implies that a subsequent study will have almost 80% probability of replication (assuming the expected effect size of the second study is the observed effect size of the first). Does this mean that *any time* one obtains an initial  $P$  value of 0.005, one should expect 80% of future studies to indeed result in a confirmation of the initial finding? To answer this question, one must remember that the replication power in Equation 1 is calculated on the assumption that the effect size observed in the initial study is the population effect size in the subsequent study. If this is indeed the case, one would obtain 80% replication rate in future findings. (Because the original effect size will more likely be overestimated, the actual success rate will be slightly lower than 80%, as explained previously). However, if the initial finding is a false positive, which will occur at a rate of 5%, the replication rate will be much lower than 80%.<sup>29</sup>

This point leads us to a vital consideration before Equation 1 can be applied. It is extremely important for the association study to be well designed so that sources of bias either are minimal or have been removed. Although this will not completely eliminate false positives, it will keep them as low as possible. Otherwise, it will be less likely for the expected effect size in the subsequent study to be even close to the observed effect size of the first study. For example, suppose an investigator reports a  $P$  value of 0.005, but the design used suffers from considerable population stratification. Then, this significant finding will more likely be a false positive, the observed effect size will very likely be a gross overestimate of the actual effect size, and application of Equation 1 will lead to exaggerated confidence in the replication power of a subsequent study (i.e., the true replication power will be much lower than 80%).

A critical issue for genome screen association studies concerns corrections for multiple testing. The question then arises as to how to compute the replication power after observing one or more positive results in a genome screen. Which test size  $\alpha$  should one use, and which  $P$  value? The answer may vary depending on the situation. The simplest case is one in which a single association peak from a genome scan is chosen to study for replication, i.e., at a locus at which a few apparently highly associated SNPs are tested for replication. It would then be misleading to use the test with the smallest  $P$  value, while still assuming a test size  $\alpha$ , to calculate the replication power. Because a genome-wide association study could potentially consist of an extremely large number of tests, an FDR correction procedure<sup>34</sup> is appealing. The appropriate  $P$  values then to use in the computation of replication power are those that are significant with the FDR procedure. In the special case when

only a few tests are performed and a Bonferroni adjustment of the test size  $\alpha$  is chosen, a  $P$  value that is significant at  $\alpha/C$  could be used to compute the replication power for the subsequent test.

In this article, we focused on the replication power for a second study based on the  $P$  value of only one first study. This is mainly because, in many cases, only one or two replications have been attempted, as can be seen, for example, from the survey conducted by Hirschhorn.<sup>3</sup> However, it is also legitimate to ask how Equation 1 could be used if we wished to compute the replication power based on the results of several initial studies, i.e., based on a meta-analytic approach.<sup>21,35</sup> In this case, we propose computing the average of the effect sizes ( $\bar{\phi}_1$ ) of the initial studies (see Appendix) and estimating the initial effective sample size ( $n_h$ ) from the harmonic mean of the initial sample sizes. The value of  $z_1$  can then be calculated by using  $z_1 = \bar{\phi}_1 \sqrt{2n_h}$ , and Equation 1 can be used to calculate the replication power of the second study, with  $n_1 = n_h$ . However, there are three points to note with such an approach: (1) all initial studies (whether significant or not) must be used in the calculation of  $z_1$ ; (2) the value of  $z_1$  must be at least as large as  $z_{crit}$ ; and (3) sources of bias, e.g., population stratification, in the initial studies must be corrected before  $\bar{\phi}_1$  is calculated; otherwise the latter will be overestimated.

How robust is Equation 1 to distributional and sample size assumptions? It is important for the sample sizes of the second and especially the first studies to be reasonably large for at least two reasons: (1) so that  $\chi^2$  tests can be legitimately used; and (2) so that the effect size in the first experiment can be consistently estimated. When testing for two proportions using small sample sizes, replication power can still be calculated from a noncentral hypergeometric distribution, but it will not be accurate. When comparing two means using the  $t$ -test, formulas for the replication power have been given by Greenwald et al.<sup>25</sup> and Posavac.<sup>27</sup> Finally, Consequence 1 (i.e.,  $P_{REP} \approx .5$  when an initial study has a  $P$  value only slightly less than  $\alpha$ ) always holds irrespective of the underlying distribution, as long as the latter is symmetric.

The rationale for assuming that an expected value of the effect size in a second study is equal to the observed value of the statistic in a first study, which is at the basis of Equation 1, can be legitimately questioned. Referring more specifically to the odds ratio, Fleiss et al.<sup>36</sup> point out that, whereas differences in proportions would vary between studies, measures of effect size could remain constant from study to study. Murphy and Myers<sup>37</sup> further state that the effect size “provides a simple metric that allows for comparison of treatment effects from different studies, areas or research.” In fact, as we mentioned previously, it is more justifiable to assume that the expected effect size in the second study is, on average, slight less than the observed effect size in the initial study. Consequently, Equation 1 really gives the replication power for the best-case scenario, so that the actual replication power of the second study will be slightly less than that given by Equation 1.

Despite the various deficiencies of the  $P$  value, which have been discussed at great length elsewhere,<sup>26,38–40</sup> we believe that, in addition to measures of effect size and confidence intervals,

researchers should continue to report  $P$  values “because tests of statistical significance provide information that effect sizes and confidence intervals do not”.<sup>27</sup> We advocate, as does Greenland et al.,<sup>25</sup> the reporting of exact  $P$  values, rather than expressions such as  $P < 0.01$  or  $P > 0.05$ . Whenever a subsequent study is planned, researchers should compute replication power based on the initial  $P$  value and on the sample sizes of the initial and subsequent studies.

The literature on genetic association studies is rife with admonitions and possible explanations for their non-replications. The reasons are usually one or more of the following<sup>6,8,9</sup>: (1) population stratification; (2) genetic heterogeneity; (3) inflation in Type I error; and (4) gene-environment interaction. We do not deny that these are important problems, and attempts should be made to correct for them. But even if these problems were to be remedied, trying to replicate many initial findings, even if they are quite significant, may be predisposed to failure and should not be interpreted as necessarily contradicting the initial association. To our knowledge, only one publication in the genetic-epidemiology literature<sup>5</sup> has acknowledged this fact. Moreover, only one publication<sup>31</sup> on association studies has actually reported calculations for the power of a replication (based on the initial observed effect size), but even that power calculation seems to be slightly inflated.

When we examined five of the six studies surveyed by Hirschhorn et al.,<sup>3</sup> for which  $P$  values could be obtained and which had been consistently replicated, we found that all of them had  $P < 10^{-3}$ . Thus, subsequent studies of these with similar sample sizes had approximately 99.8% replication power! We also randomly sampled 50 of the 166 initial studies that had reported exact  $P$  values. We found that: (1) 78% had  $P$  values  $> 0.005$ , implying a subsequent study of similar sample size would be underpowered (a replication power of  $< 80\%$ ); (2) 38% had  $P$  values greater or equal to 0.02, implying a subsequent study of similar sample size would be seriously underpowered (a replication power of  $< 64\%$ ). Although it is true that many of the replication failures reported by Hirschhorn et al.<sup>3</sup> could be the result of the several reasons mentioned by these authors (e.g., population stratification, heterogeneity), it is equally true that, even for the remaining cases in which these sources of error were actually minimal, replication failures were bound to occur, because of their low a priori replication power. Furthermore, even if sources of bias had actually been minimal in most cases, replication success would still be low. On a second look, therefore, Hirschhorn et al.’s<sup>3</sup> review should perhaps not look so depressing, because it contains many potential pseudo-failures, that is, replications that were just not meant to be. We must stress, however, that our article in no way attempts to challenge Hirschhorn et al.’s arguments. It acknowledges these, but it also adds a whole new perspective to the issue of replication: old problems still need to be grappled with, because they reduce the odds of false positives, *but even if they are addressed, low replication power will continue to deny high replication rates in future association studies.*

In conclusion, if the  $P$  value in one’s initial study is very small (e.g.,  $P = 0.005$  when  $\alpha = 0.05$ ), then one can indeed anticipate a high replication probability. However, for more commonly observed  $P$  values (e.g.,  $P = 0.02$  and even  $P = 0.01$ ,

when  $\alpha = 0.05$ ), replication probabilities are notably lower than one might hope. In these situations, one should not be surprised by a failure to replicate.

### APPENDIX: PROOF OF EQUATION 1

Let  $p_1, p_2$  be the two population proportions being compared. Let  $P_1^{(i)}, P_2^{(i)}$  be the corresponding sample proportions (based on  $n_i$  cases and  $n_i$  controls) in the  $i^{\text{th}}$  study ( $i = 1, 2$ ). The  $\chi^2$  test statistic,  $X_i^2$ , for the  $i^{\text{th}}$  study (see Table 2) can be converted into a  $z$ -test,  $Z_i$ , where  $Z_i = \text{sign}(P_1^{(i)} - P_2^{(i)}) \times \sqrt{X_i^2}$ . The observed effect size for the  $i^{\text{th}}$  study is  $\hat{\phi}_i = \sqrt{X_i^2 / (2n_i)} = Z_i / \sqrt{2n_i}$ .<sup>41</sup> Let the population effect size for the  $i^{\text{th}}$  study be  $\phi_i$ . If the first test has size  $\alpha$ , then its critical value is  $z_{crit} = \Phi^{-1}(1 - \alpha/2)$ , where  $\Phi$  is the standard normal cumulative distribution function. Also, let the observed value of  $Z_1$  be  $z_1$ . Under  $H_0: p_1 = p_2$ , for the first study,  $Z_1 \sim N(0, 1)$ . Conditional on  $\phi_2 = \hat{\phi}_1$  and on  $n_1, n_2, Z_2 \sim N(\mu_2, \sigma_2^2)$  but  $\mu_2 \neq 0$  and  $\sigma_2^2 \neq 1$ . Indeed, using  $\hat{\phi}_i = Z_i / \sqrt{2n_i}$ , we have

$$\begin{aligned} \mu_2 &= E(Z_2 | \phi_2 = \hat{\phi}_1, n_1, n_2) \\ &= \phi_2 \sqrt{2n_2} \\ &= \hat{\phi}_1 \sqrt{2n_2} \\ &= z_1 \sqrt{\frac{n_2}{n_1}}, \end{aligned}$$

and, using  $\text{var} \hat{\phi}_i = (1 - \phi_i^2)^2 / (2n_i - 2)$  (see Rosenthal<sup>42</sup>) and assuming  $n_2$  is reasonably large, we have

$$\begin{aligned} \sigma_2^2 &= \text{var}(Z_2 | \phi_2 = \hat{\phi}_1, n_1, n_2) \\ &= 2n_2 \text{var} \hat{\phi}_2 \\ &= 2n_2 \frac{(1 - \phi_2^2)^2}{2n_2 - 2} \\ &\approx (1 - \phi_2^2)^2 \\ &= (1 - \hat{\phi}_1^2)^2 \\ &= \left(1 - \frac{z_1^2}{2n_1}\right)^2. \end{aligned}$$

Therefore, the replication probability is

$$\begin{aligned} p_{REP} &= \Pr \left\{ Z_2 \geq z_{crit} \mid Z_2 \sim N \left[ z_1 \sqrt{\frac{n_2}{n_1}}, \left(1 - \frac{z_1^2}{2n_1}\right)^2 \right] \right\} \\ &= \Pr \left\{ Z \geq \frac{z_{crit} - z_1 \sqrt{n_2/n_1}}{1 - z_1^2/(2n_1)} \right\}, \text{ where } Z \sim N(0, 1), \\ &= \Phi \left\{ \frac{z_1 \sqrt{n_2/n_1} - z_{crit}}{1 - z_1^2/(2n_1)} \right\}. \end{aligned}$$

### ACKNOWLEDGMENTS

This work was supported in part by NIH grants R01 AA013654, DK-31813, NS27941, DK3177. We wish to thank Dr Bruce Levin for helpful discussions and Dana Politis for her invaluable time.

## References

1. Risch N, Merikangas K. The future of genetic studies of complex human diseases. *Science* 1996;273:1516–1517.
2. Runyon RP, Haber A, Pittenger DJ, Coleman KA. Fundamentals of behavioral statistics, 8<sup>th</sup> edition. New York: McGraw-Hill, 1995.
3. Hirschhorn JN, Lohmueller K, Byrne E, Hirschhorn K. A comprehensive review of genetic association studies. *Genet Med* 2002;4:45–61.
4. Vieland VJ. The replication requirement. *Nat Genet* 2001;29:244–245.
5. Göring HH, Terwilliger JD, Blangero J. Large upward bias in estimation of locus-specific effects from genomewide scans. *Am J Hum Genet* 2001;69:1357–1369.
6. Hirschhorn JN, Altshuler D. Once and again—issues surrounding replication in genetic association studies. *J Clin Endocrinol Metab* 2002;87:4438–4441.
7. Dahlman I, Eaves IA, Kosoy R, Morrison VA, Heward J, Gough SC, et al. Parameters for reliable results in genetic association studies in common disease. *Nat Genet* 2002;30:149–150.
8. Colhoun HM, McKeigue PM, Davey SG. Problems of reporting genetic associations with complex outcomes. *Lancet* 2003;361:865–872.
9. Ott J. Association of genetic loci: replication or not, that is the question. *Neurology* 2004;63:955–958.
10. Hodge SE. Linkage analysis versus association analysis: distinguishing between two models that explain disease-marker associations. *Am J Hum Genet* 1993;53:367–384.
11. Halliburton R. Introduction to population genetics. New Jersey: Prentice Hall, 2003.
12. Schork NJ, Fallin D, Thiel B, Xu X, Broeckel U, Jacob HJ, et al. The future of genetic case-control studies. In: Rao DC, Province MA, editors. Genetic dissection of complex traits. San Diego, CA: Academic Press, 2001:191–212.
13. Rosenberg NA, Nordborg M. A general population-genetic model for the production by population structure of spurious genotype-phenotype associations in discrete, admixed or spatially distributed populations. *Genetics* 2006;173:1665–1678.
14. Devlin B, Roeder K. Genomic control for association studies. *Biometrics* 1999;55:997–1004.
15. Devlin B, Roeder K, Wasserman L. Genomic control, a new approach to genetic-based association studies. *Theor Popul Biol* 2001;60:155–166.
16. Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus genotype data. *Genetics* 2000;155:945–959.
17. Pritchard JK, Stephens M, Rosenberg NA, Donnelly P. Association mapping in structured populations. *Am J Hum Genet* 2000;67:170–181.
18. Heiman GA, Hodge SE, Gorroochurn P, Zhang J, Greenberg DA. Effect of population stratification on case-control association studies. I. Elevation in false positive rates and comparison to confounding risk ratios (a simulation study). *Hum Hered* 2004;58:30–39.
19. Gorroochurn P, Hodge SE, Heiman G, Greenberg DA. Effect of population stratification on case-control association studies. II. False-positive rates and their limiting behavior as number of subpopulations increases. *Hum Hered* 2004;58:40–48.
20. Gorroochurn P, Heiman GA, Hodge SE, Greenberg DA. Centralizing the non-central chi-square: a new method to correct for population stratification in genetic case-control association studies. *Genet Epidemiol* 2006;30:277–289.
21. Rosenthal R. Cumulating evidence. In: Keren G, Lewis C, editors. A handbook for data analysis in the behavioral sciences: methodological issues. Hillsdale, NJ: Lawrence Erlbaum, 1993:519–559.
22. Oakes M. Statistical inference: a commentary for the social and behavioral sciences. New York: Wiley; 1986.
23. Gigerenzer G. The superego, the ego, and the id in statistical reasoning. In: Keren G, Lewis C, editors. A handbook for data analysis in the behavioral sciences: methodological issues. Hillsdale, NJ: Lawrence Erlbaum, 1993:311–339.
24. Maxwell SE, Delaney HD. Designing experiments and analyzing data: a model comparison perspective, 2nd edition. Mahwah, NJ: Lawrence Erlbaum, 2004.
25. Greenwald AG, Gonzalez R, Harris RJ, Guthrie D. Effect sizes and p values: what should be reported and what should be replicated? *Psychophysiology* 1996;33:175–183.
26. Nickerson RS. Null hypothesis significance testing: a review of an old and continuing controversy. *Psychol Methods* 2000;5:241–301.
27. Posavac EJ. Using p values to estimate the probability of a statistically significant replication. *Understanding Stat* 2002;1:101–112.
28. Daly LE, Bourke GJ. Interpretation and uses of medical statistics, 5th edition. Oxford: Blackwell Science, 2000.
29. Sohn D. Statistical significance and replicability: why the former does not presage the latter. *Theory Psychol* 1998;8:291–311.
30. Tversky A, Kahneman D. Belief in the law of small numbers. In: Keren G, Lewis C, editors. A handbook for data analysis in the behavioral sciences: methodological issues. Hillsdale, NJ: Lawrence Erlbaum Associates, 1993:341–349.
31. Heils A, Haug K, Kunz WS, Fernandez G, Horvath S, Rebstock J, et al. Interleukin-1beta gene polymorphism and susceptibility to temporal lobe epilepsy with hippocampal sclerosis. *Ann Neurol* 2000;48:948–950.
32. Tan NC, Mulley JC, Berkovic SF. Genetic association studies in epilepsy: “the truth is out there.” *Epilepsia* 2004;45:1429–1442.
33. Ioannidis JP, Ntzani EE, Trikalinos TA, Contopoulos-Ioannidis DG. Replication validity of genetic association studies. *Nat Genet* 2001;29:306–309.
34. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Stat Soc B* 1995;57:289–300.
35. Wolf FM. Meta-Analysis. Thousand Oaks, CA: Sage Publications, 1986.
36. Fleiss JL, Levin B, Paik MC. Statistical methods for rates and proportions, 3rd edition. New York: Wiley, 2003.
37. Murphy KR, Myers B. Statistical power analysis. Mahwah, NJ: L. Erlbaum Associates, 2004.
38. Cohen J. The earth is round ( $p < .05$ ). *Am Psychol* 1994;49:997–1003.
39. Royall RM. Statistical evidence: a likelihood paradigm. London: Chapman and Hall, 1997.
40. Goodman SN. Toward evidence-based medical statistics. I. The P value fallacy. *Ann Intern Med* 1999;130:995–1004.
41. Tatsuoaka M. Effect size. In: Keren G, Lewis C, editors. A handbook for data analysis in the behavioral sciences: methodological issues. Hillsdale, NJ: Lawrence Erlbaum Associates, 1993:461–479.
42. Rosenthal R. Parametric measures of effect size. In: Cooper H, Hedrick PW, editors. The handbook of research synthesis. New York: Russell Sage Foundation; 1994:231–244.