# Cost-Effective Designs for Linkage Disequilibrium Mapping of Complex Traits

Susan K. Service,[1] Lodewijk A. Sandkuijl,[2,*] and Nelson B. Freimer[1]

[1]Center for Neurobehavioral Genetics, University of California at Los Angeles, Los Angeles; and [2]Department of Medical Statistics, Leiden University Medical Center, Leiden, The Netherlands

The current development of densely spaced collections of single nucleotide polymorphisms (SNPs) will lead to genomewide association studies for a wide range of diseases in many different populations. Determinations of the appropriate number of SNPs to genotype involve a balancing of power and cost. Several variables are important in these determinations. We show that there are different combinations of sample size and marker density that can be expected to achieve the same power. Within certain bounds, investigators can choose between designs with more subjects and fewer markers or those with more markers and fewer subjects. Which designs are more cost-effective depends on the cost of phenotyping versus the cost of genotyping. We show that, under the assumption of a set cost for genotyping, one can calculate a "threshold cost" for phenotyping; when phenotyping costs per subject are less than this threshold, designs with more subjects will be more cost-effective than designs with more markers. This framework for determining a cost-effective study will aid in the planning of studies, especially if there are choices to be made with respect to phenotyping methods or study populations.

There is much current interest in the use of SNPs to conduct genomewide association studies of common diseases. The density of SNPs that will be required for such studies remains the subject of intense debate. The determination of the optimal number of SNPs for a particular study involves a balance between obtaining adequate power to detect association and keeping the cost of the study at a realistic level. The issue of cost in genotyping is vastly more important for association studies than for linkage studies, because the number of markers used in association studies is so much greater. For example, if one wished to double the density of microsatellite markers for a genomewide linkage study (e.g., from 400 to 800), the increase in cost would be trivial compared with that incurred in doubling the density of SNPs for a genomewide association study (e.g., from 40,000 to 80,000).

The power of an association study depends on several variables, which fall into two categories. The first category consists of known variables determined by the study design, including the number of markers genotyped and the sample size. The age and pattern of growth of the population under study are central determinants

of the power of genomewide association studies; for some populations, such as isolates with detailed genealogic records, these may be considered "known" variables. The second category consists of variables that are unknown, such as the magnitude of the difference between cases and controls in the frequency of the SNP allele associated with disease. In determining the optimal strategy for balancing cost and power of an association study, it is useful to examine the impact of making different choices for the known variables, on the basis of particular assumptions about the value of the unknown variables. We present here a simple framework for conducting such evaluations.

The choice of study population affects the power and cost of association studies, in that populations in which extensive linkage disequilibrium (LD) surrounds disease-susceptibility variants will require genotyping fewer markers than samples with less LD. Indeed, the discovery that, even in the most outbred populations, the genome is organized in blocks of conserved haplotypes has important implications for the cost of genotyping, and it is a major rationale for the "Hap-Map" project, which is aimed at identifying each of these blocks. For example, the initial proposal for genomewide association studies by Risch and Merikangas (1996) envisioned genotyping about one million SNPs per study. Consider one of their scenarios (a genotype relative risk [*GRR*] of 2.0 and a disease variant with frequency of 0.10); this situation would require genotyping 695 trios for adequate power, using a transmission/disequilibrium test. Under the assumption of a cost of $0.10 per genotype (which is near

the low end of currently obtainable genotyping costs), this genotyping study would cost ~$210 million. If one simply extrapolates for the entire genome—from the published haplotype-block information on chromosome 21 (Patil et al. 2001)—one would require 422,500 SNPs for complete genome coverage, for a much-reduced cost of ~$88 million. Many fewer numbers of SNPs may be needed in isolated populations, particularly in recently founded isolates, in which LD may extend for considerable distances (Service et al. 2001; Hall et al. 2002).

We show here a simple framework for evaluating the costs of case-control association studies designed to achieve a stated power. This framework is based on modeling the decay of LD around a disease locus as a function of the genetic distance between the associated marker allele and the disease locus and of the number of generations since the mutation was introduced into the population. By specifying the depth of a study population (i.e., the number of generations since its founding), it is possible to estimate the expected value of the LD coefficient ($D'$) for various map distances. Using the estimate of $D'$—together with assumptions about the frequency of the disease allele, the prevalence of the disease, the *GRR* associated with various genotypes, and the frequency of the associated SNP allele in the control sample—one can calculate the power of a specific sample size, and one can calculate the sample size necessary to detect a signal at a specified power. In calculating the expected value of $D'$, the present analysis does not consider other factors, such as SNP mutations, that may affect the decay of LD with genetic distance. This simple analysis is not intended to serve as a primer for estimating the power of association studies per se but rather as a heuristic device to demonstrate how one may consider, for a particular study population, two solutions to the problem of balancing cost and power: increasing the density of SNPs or increasing the size of the study sample. In addition to using values of $D'$ based on theoretical expectations, we have applied this methodology to estimates of $D'$ that were empirically derived from an outbred population (Reich et al. 2001).

## Methods

We modeled the decay of LD around a disease locus as a function of the genetic distance between the associated marker allele and the disease locus and of the number of generations since the mutation was introduced into the population. Given a number of genotyped SNP markers ($M$), the average genetic distance from a disease-susceptibility locus to the nearest marker (assuming the worst-case scenario of the disease locus being located midway between two markers) can be found by $3,673/(2 \cdot M)$, where 3,673 is the length of the genome in centimorgans (Kong et al. 2002). With the number of mark-

ers we are considering in the present article (a minimum of 10,000), when this distance is divided by 100 it also approximates the recombination fraction between disease locus and marker locus when no interference is assumed. If the disease locus is assumed to have alleles C and c and if the nearest marker is assumed to have alleles B and b (with the B allele associated with disease), the frequency of the B allele on chromosomes carrying the high-risk disease allele C can be found by
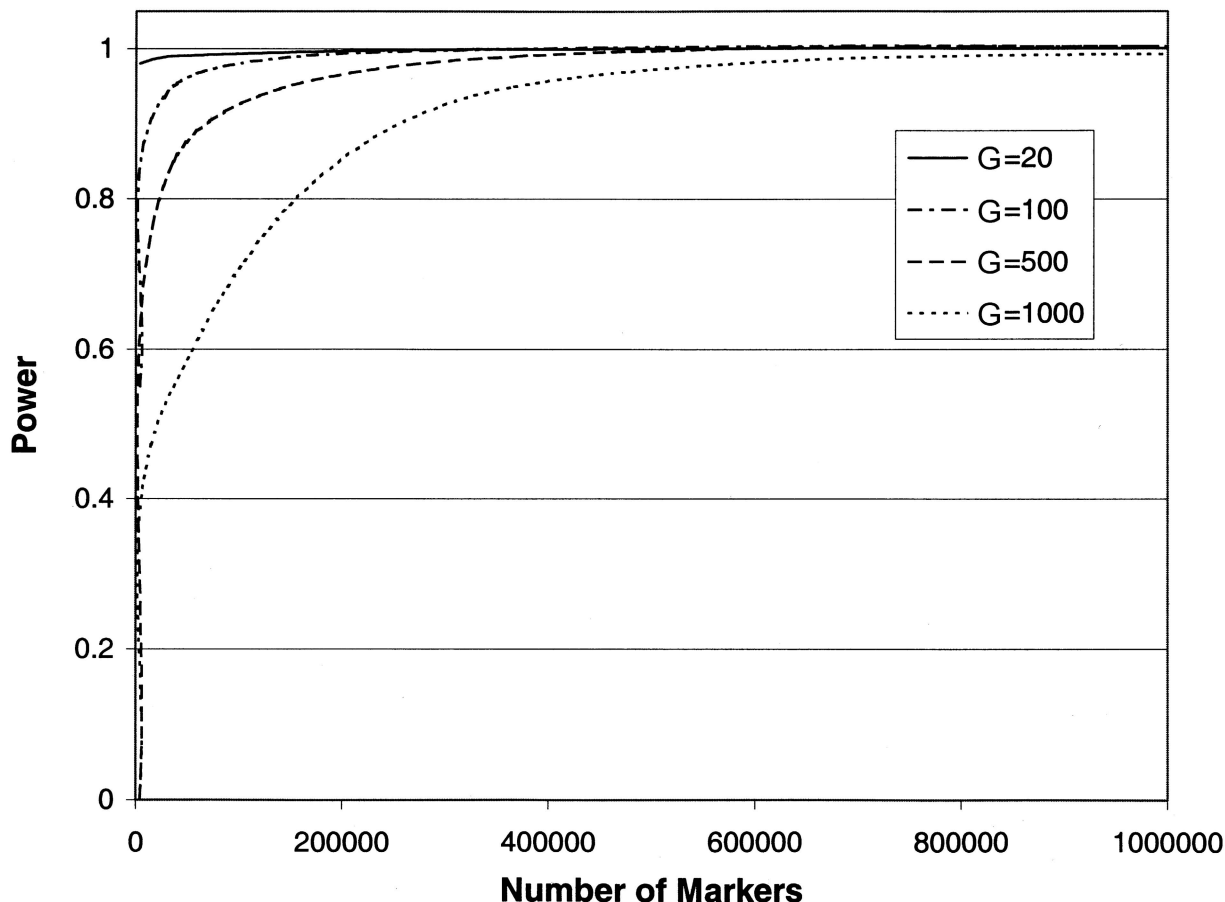
$$P(B \mid C) = (1 - \theta)^G + [1 - (1 - \theta)^G] \cdot P(B) ,$$

where $G$ is the number of generations since the mutation was introduced, $\theta$ is the recombination fraction between marker and disease locus, and $P(B)$ is the frequency of the B allele in the population. The frequency of the B-C haplotype, $P(\text{B-C})$, can be found by $P(B|C)P(C)$, where $P(C)$ is the frequency of the C allele. The expected value of $D'$ can be found as the difference between $P(\text{B-C})$ and $P(B)P(C)$, standardized by the maximum this value could take, given the allele frequencies of C and B. For example, if $\theta = 0.0005$ (corresponding to a density of 1 marker per 0.1 cM, or 36,730 markers in a genome of 3,673 cM) and if $P(C) = 0.1$ and $P(B) = 0.25$, then the expected value of $D'$ for a population with $G = 500$ is 0.78. The power of a genetic association study can then be calculated using $D'$, $P(C)$, $P(B)$, the prevalence of the disease, and the *GRR* associated with the CC and Cc genotypes (see Genetic Power Calculator Web site) (Purcell et al. 2003). For all calculations, we used an $\alpha$ level of 0.0001, and, when calculating sample sizes, we assumed equal numbers of cases and controls and a power level of 80%. Unless otherwise specified, we set the following values: $P(B) = 0.25$, $P(C) = 0.10$, prevalence of disease = 0.15, and GRRs = 2.5 for both the CC and Cc genotypes. The assumption of these values for $P(B)$, $P(C)$, *GRR*, and disease prevalence results in genotype frequencies for the disease and marker loci shown in table 1.

**Table 1**

**Genotype Frequencies at the Disease and Marker Loci**

| | | FREQUENCY AT | | | |
| | | BB/Bb with $D' =$ | | | |
| GROUP | CC/Cc | .9 | .6 | .4 | .2 |
|---|---|---|---|---|---|
| Case | .02/.35 | .09/.46 | .08/.43 | .08/.41 | .07/.39 |
| Control | .01/.15 | .06/.37 | .06/.37 | .06/.37 | .06/.37 |

NOTE.—Genotype frequencies at the disease and marker locus when $P(C)$ (the frequency of the risk allele at disease locus) is .10, $P(B)$ (the frequency of the associated marker locus) is .25, disease prevalence is .15, and GRRs associated with the CC and Cc genotypes are both 2.5. The genotype frequencies at the marker locus in cases depend on the magnitude of LD between disease and marker; four sample values of $D'$ are presented.

**Figure 1**     Power to detect a false null hypothesis at the 0.0001 $\alpha$ level versus number of markers. A sample size of 1,000 cases and controls was used in all calculations. We assume the frequency of the disease mutation to be 0.10, the prevalence of the disease to be 0.15, the frequency of the associated marker allele to be 0.25, and the *GRR* associated with disease allele homozygotes and heterozygotes to be 2.5.

We also applied our methodology to empirically derived estimates of $D'$. We used the data of Reich et al. (2001), who estimated pairwise values of $D'$ between SNP markers in 19 randomly selected genome regions in a sample of 44 persons from an outbred U.S. population of northern European descent. Reich et al. found a great deal of variability in estimates of $D'$, both by physical distance between markers and across genome regions. We used their reported average values of physical distance in relation to $D'$. In estimating the number of markers needed to cover the genome, at a density specified by a given $D'$ value, we assume the genome is 3,000 Mb.
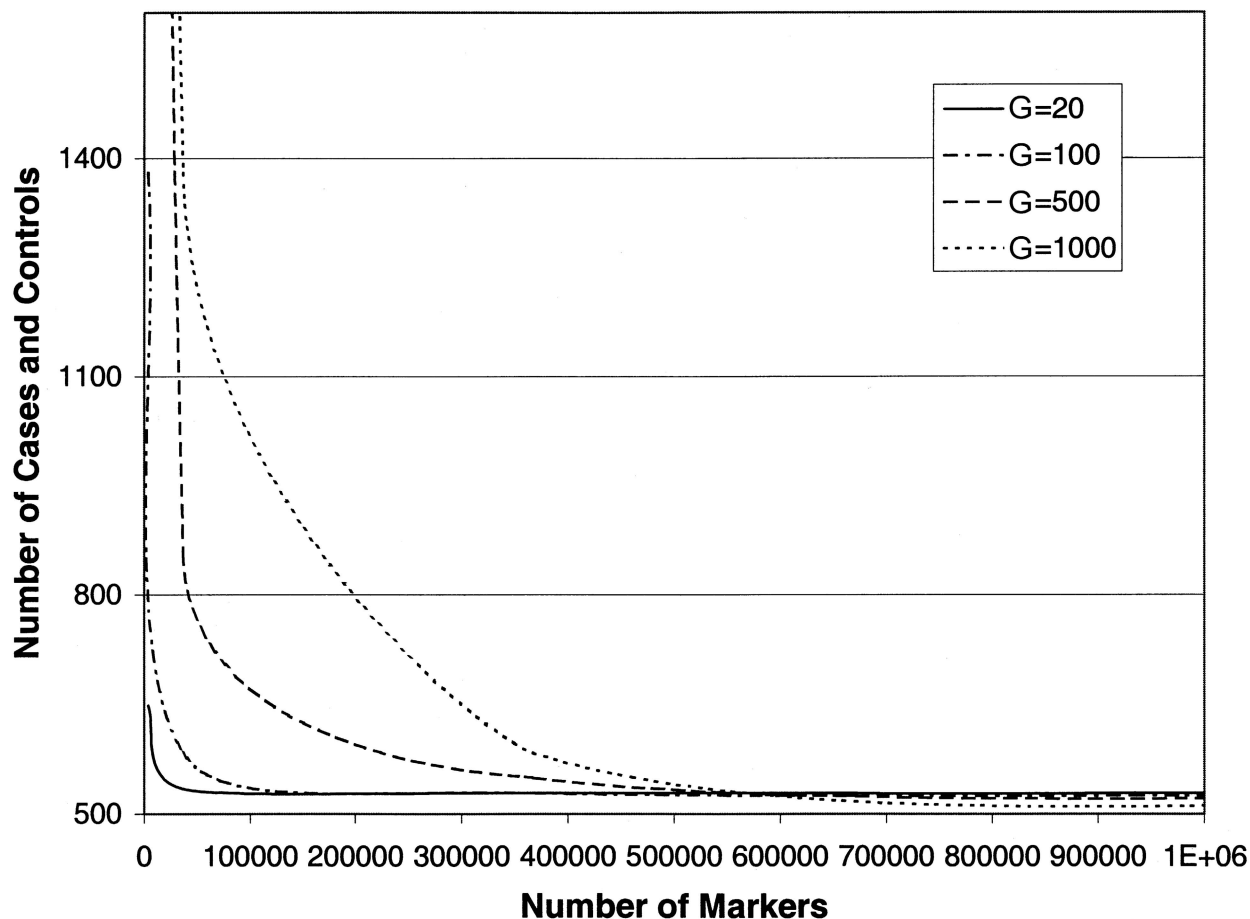
## Results

### Power in Relation to Population Depth

We first examined the relationship between marker density and power for populations of four different depths, equivalent to a recently founded isolate (20 gen-

erations), an older isolate (100 generations), and two outbred populations (500 and 1,000 generations). For all populations, sample sizes of 1,000 cases and controls were used. We show the results of this analysis for a single disease-variant frequency and associated SNP allele frequency (fig. 1). As expected, an increase in the marker density enhances the power to detect association. In the younger populations, far fewer markers are needed to obtain adequate power than in the older populations. For all population depths, the increase in power reaches a plateau as the numbers of markers increase; however, this plateau occurs at very different points in young populations (<5,000 markers) compared with old populations (~600,000 markers).

### Power in Relation to Sample Size

One can also increase power to detect association by increasing the sample size. As shown in figure 2, for any depth of population, one can identify alternative strategies to obtain a specified level of power (in this case,

**Figure 2** Each line represents combinations of numbers of markers and numbers of cases and controls necessary to achieve 80% power to reject a false null hypothesis at an $\alpha$ level of 0.0001. We assume the frequency of the disease mutation to be 0.10, the prevalence of the disease to be 0.15, the frequency of the associated marker allele to be 0.25, and the *GRR* associated with disease allele homozygotes and heterozygotes to be 2.5.

80%) on the basis of varying combinations of marker density and sample size. For example, for the population in which $G = 1,000$, one can achieve 80% power with a sample of ~800 individuals (cases and controls) and 200,000 markers or with a sample of ~550 cases and controls and 400,000 markers.
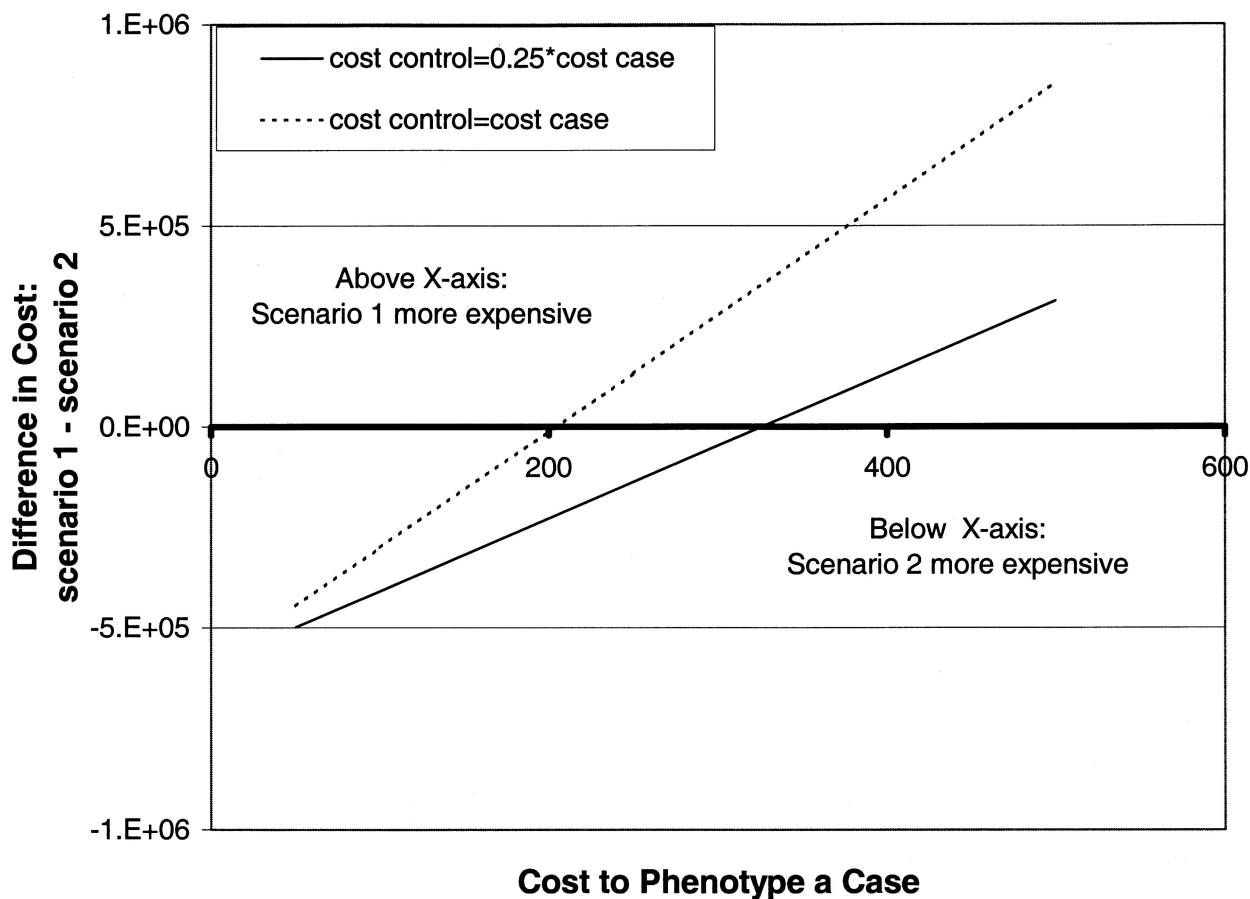
*Optimization of Strategy in Terms of Cost*

In determining whether the optimal study design involves increasing the marker density or increasing the sample size, one can estimate the cost to achieve a given power. That is, one can estimate how expensive it is to add one individual to the sample (case or control), compared with the cost of genotyping one additional marker. Although genotyping costs are essentially fixed for a given genotyping platform, the costs of sample collection are specific to particular projects, depending on such

variables as the cost of ascertainment in particular locales and the type of phenotyping that must be done. For most projects these costs can be varied by altering the study design, and it is possible to identify study designs that minimize cost. For a given genotyping cost, one can achieve 80% power through either of two strategies: larger sample and fewer markers or a smaller sample with more markers. Consider an example with $G$ of 500. To achieve 80% power with 12,243 markers (corresponding to $\theta = 0.0015$ and $D' = 0.47$, when $P[C] = 0.1$ and $P[B] = 0.25$) requires 2,299 cases and controls (scenario 1); with 36,370 ($D' = 0.78$) markers, this power level requires only 855 cases and controls (scenario 2). Which is the cost-effective strategy? The difference in cost of these two scenarios depends on the cost to phenotype a case. Assume that determinations of genotypes cost $0.10 each. We make two assumptions

about the cost to ascertain/phenotype a control. In one example, this cost is 25% of the cost to ascertain and phenotype a case; in the other example, this cost is the same as the cost to ascertain and phenotype a case. In the first example, the cost of a study can be calculated as $N \cdot (T + M \cdot 0.10) + N \cdot (T \cdot 0.25 + M \cdot 0.10)$, where $N$ is the number of cases and controls, $M$ is the number of markers, and $T$ is the cost to ascertain and phenotype a case. In the second example, the cost of a study can be calculated as $2N \cdot (T + M \cdot 0.10)$. We have calculated the difference in cost of these two scenarios for varying values of $T$ (fig. 3). The point at which the line crosses the $X$-axis in figure 3 determines the value at which the cost of the two scenarios is equal. The examples depicted in this figure show that, when controls are 25% as expensive to phenotype as cases, scenario 1 is more expensive than scenario 2 when the

cost of collecting each case is >\$327. We term the cost depicted by this point the "threshold phenotyping cost." When cases and controls are equally expensive to phenotype, the threshold phenotyping cost decreases to \$205.

We evaluated how the threshold phenotyping cost would vary in different populations, ranging from a depth of 20 generations to a depth of 1,000 generations (table 2). In each case, we compared a scenario with 10,000 markers and a scenario with 50,000 markers, calculating the sample size needed to achieve power of 80%. We calculated threshold phenotyping costs under the assumption that controls are 25% as expensive to phenotype as cases and also under the assumption that the cost of phenotyping the two groups is equal. In both examples, the threshold phenotyping cost decreases progressively in the deeper populations. The threshold of



**Figure 3**     Difference in cost for two scenarios versus the cost to phenotype a case. Scenario 1 uses 12,243 markers and 2,299 cases and controls, scenario 2 uses 36,370 markers and 855 cases and controls. Two lines are plotted: one in which the cost to phenotype a control is set at 25% of the cost to phenotype a case and the other in which the phenotyping costs for cases and controls are equal. When the cost per case is less than \$327 (and the cost for controls is 25% of that for cases) or \$205 (and the costs for controls and cases are equal), it is less expensive to genotype 12,243 markers in 2,299 persons; when the cost is >\$327 (\$205 when costs for cases and controls are equal), it is less expensive to genotype 36,730 markers in 855 cases and controls.

**Table 2**

**Threshold Phenotyping Costs for Various Population Depths**

| G | N1 | N2 | THRESHOLD PHENOTYPING COST WHEN | |
|---|---|---|---|---|
| | | | Control Cost = 25% of Case Cost | Control Cost = 100% of Case Cost |
| 20 | 572 | 534 | 88,337 | 55,210 |
| 100 | 758 | 572 | 18,082 | 11,301 |
| 500 | 3,157 | 758 | 422 | 264 |
| 1,000 | 19,350 | 1,085 | 0 | 0 |

NOTE.—$N1$ is the sample size to achieve 80% power in the scenario in which the number of markers = 10,000, and $N2$ is the sample size needed to achieve 80% power in an alternative scenario in which the number of markers = 50,000. As long as the cost to ascertain and phenotype a case is less than the threshold cost, $N1$ and $M1$ represent the more cost-effective strategy.

zero in the oldest population ($G = 1,000$) indicates that, given the assumptions we have made in this example, it is more cost-effective to genotype a larger number of markers than to add to the sample, regardless of the cost of phenotyping each case.

We also modeled the effect on threshold costs of different values for variables that are unknown at the outset of a study. Table 3 shows a range of threshold costs for two such variables. First, we varied the frequency of the SNP allele associated with disease susceptibility. Next we varied the *GRR* associated with the CC and Cc genotypes at the disease locus. In these scenarios, we used the same assumptions as in table 2 (marker sets of 10,000 and 50,000 and the sample size needed to obtain 80% power). In table 3 we show the results for a single population depth ($G = 500$), again with two examples for the cost of phenotyping controls relative to the cost of phenotyping cases. The threshold costs are lowest with increasing marker allele frequencies and with decreasing values for *GRR*. This result reflects the fact that there is generally greater power to detect association when the frequency of the SNP allele is close to the frequency of the disease allele (Garner and Slatkin 2003) and the fact that there is greater power in a higher *GRR*. However, the different values for these variables have a relatively small effect on the threshold cost compared to different values of population depth.

### Empirically Derived Estimates of D′

Reich et al. (2001) found the average value of $D′$ to be 0.65 for markers 0.04 Mb apart and 0.35 for markers 0.16 Mb apart. These physical distances correspond to 75,000 markers and 18,750 markers, respectively, for a genome of 3,000 Mb. To calculate the sample size of cases and controls needed to identify an association at an $\alpha$ level of 0.0001 with 80% power for these two marker densities, we first assumed the same estimates for $P(C)$, disease prevalence, $P(B)$, and *GRR* that we

specified in the "Methods" section. The sample size required for 75,000 markers was estimated to be 1,219 cases and controls, and for 18,750 markers it was estimated to be 4,107 cases and controls. When the cost of phenotyping controls is assumed to be 25% of the cost of phenotyping cases, the threshold phenotyping cost is ~$800; with equal costs to phenotype cases and controls, the threshold decreases to ~$500. When a slightly stronger effect size ($GRR = 3.0$) is used, these thresholds are increased slightly, to $826 and $526, respectively.

## Discussion

The results presented here show how several variables interact to determine the costs of genomewide association studies. These costs reflect different study design strategies for achieving a specified power to detect association. The examples that we used represent only a few of the possible combinations of these variables. For example, we recognize that there may be a wide range of conceivable ratios for the cost of sampling and phenotyping cases compared with controls, depending on the disease and population. The values that we used in this analysis for the costs of collecting controls (25% of the cost of collecting a case and 100% of the cost of collecting a case) may be unrealistic in particular situations, and investigators may wish to consider other values for this variable. It is straightforward to employ alternative values for any of the variables that we are considering. The approach that we have employed therefore provides a framework for investigators to utilize information regarding these variables that is specific to their projects, in making cost-effective designs. For example, a given investigator may have little choice regarding the study population but may have several different methods available for phenotyping. With a limited budget in such a situation, phenotyping methods with

**Table 3**

**Impact of Varying Marker Allele Frequencies and *GRR* Associated with Homozygous and Heterozygous Genotypes at the Disease Locus**

| | | | | THRESHOLD PHENOTYPING COST WHEN | |
| ALLELE FREQUENCY | *GRR* | *N1* | *N2* | Control Cost = 25% of Case Cost | Control Cost = 100% of Case Cost |
| --- | --- | --- | --- | --- | --- |
| .05 | 2.5 | 2,463 | 649 | 689.75 | 431.09 |
| .10 | 2.5 | 1,156 | 302 | 663.23 | 414.52 |
| .20 | 2.5 | 2,407 | 587 | 464.18 | 290.19 |
| .30 | 2.5 | 4,014 | 953 | 392.55 | 245.34 |
| .25 | 1.5 | 20,230 | 4,763 | 370.85 | 231.78 |
| .25 | 2.0 | 6,037 | 1,437 | 399.30 | 249.57 |
| .25 | 2.5 | 3,157 | 758 | 422.17 | 263.86 |
| .25 | 3.0 | 2,064 | 498 | 435.25 | 272.03 |

NOTE.—*N1* (*N2*) is the sample size of cases and controls necessary to achieve 80% power with 10,000 (50,000) markers at an $\alpha$ level of .0001. This example is for $G = 500$.

different costs result in the ability to collect different sample sizes, and, in turn, these different sample sizes will require varying numbers of markers to achieve desired power levels. By using our approach, investigators can evaluate how the costs of these different phenotyping methods may affect the design and ultimately success of the study.

Of the variables that we have considered, the depth of the study population produces perhaps the most dramatic effect on both cost and power. All of the examples suggest that, for recently founded isolates, it may be possible to conduct association studies with fewer markers, smaller sample sizes, and lower costs than in outbred populations. These observations are concordant with the results of theoretical and empirical studies regarding the extent of LD in different populations (Reich et al. 2001; Service et al. 2001; Angius et al. 2002; Hall et al. 2002; Varilo et al. 2003). In addition, genomewide association studies have already been performed in some young isolates, using sets of ~1,000 microsatellite markers (whose information content is probably equivalent to ~3,000–5,000 standard SNPs, based on average heterozygosity) (Ophoff et al. 2002; Vaessen et al. 2002).

The approach that we have presented represents a first step in evaluation of the relationship between power and cost in genetic association studies. The framework we used requires assumptions about such variables as the heterogeneity of the disease, the degree of LD between disease and marker locus, or the frequency of the marker allele associated with disease. In particular, we made the simplifying assumption that the extent of LD in a given population is uniform throughout the genome and therefore that evenly spaced sets of markers will be employed to achieve a set level of power. We do not yet have an adequate theoretical framework or empirical data set to incorporate into our analysis such factors as variability between populations in the size of haplotype blocks in particular regions of the

genome. When more data become available, it may be useful to modify this framework to consider genomewide genotyping strategies that employ different spacing of markers in blocks of low LD compared with blocks of high LD. It will also likely be valuable to employ methods for evaluating power that use either empirically estimated or theoretically derived distributions for key parameters in power calculations, to avoid making assumptions about their values (Schork 2002). Regardless of the method used to calculate power, the procedure we describe can be applied to identify cost-effective study designs.

The variability in the extent of LD across the genome raises an additional issue for consideration in modifying our approach, namely the uncertainty regarding how to specify the level of statistical significance that is appropriate in estimating the power of a genomewide LD study. In particular, the threshold for significance for genomewide association studies must be adjusted for the multiple independent comparisons that are involved. Determining the precise number of such comparisons, however, is complicated by the fact that the markers used in these studies will demonstrate varying degrees of LD with each other and are therefore neither fully dependent nor fully independent. New methods for correcting for multiple comparisons in genomewide association studies may facilitate such determinations (Sabatti et al., in press). The degree of independence of the markers will vary between populations, and therefore it is probably not appropriate to suggest a single significance threshold that should be used for cost-power analyses, regardless of the population being considered. It may therefore be prudent for investigators to consider a range of significance thresholds.

The calculations in the present article are based on a very simple statistical test (a difference of two proportions). Many other tests of association exist that are more powerful. We chose this test as an example, simply

because it is easy to calculate either the power for a stated sample size or the sample size needed for a given power. More powerful tests have more complicated distributions, and power is usually found by simulation. For any given statistical test, there are likely to be several combinations of sample size and marker density that can produce similar powers; the more cost-effective strategy can be determined for a given situation, and this information can aid in study planning.

## Acknowledgments

## Electronic-Database Information

The URL for data presented herein is as follows:

Genetic Power Calculator, http://statgen.iop.kcl.ac.uk/gpc/

## References

Angius A, Bebbere D, Petretto E, Falchi M, Forabosco P, Maestrale B, Casu G, Perisco I, Melis PM, Pirastu M (2002) Not all isolates are equal: linkage disequilibrium analysis on Xq13.3 reveals different patterns in Sardinian sub-populations. Hum Genet 111:9–15

Garner C, Slatkin M (2003) On selecting markers for association studies: patterns of linkage disequilibrium between two and three diallelic loci. Genet Epidemiol 24:57–67

Hall D, Wijsman EM, Roos JF, Gogos JA, Karaviorgou M (2002) Extended intermarker linkage disequilibrium in the Afrikaners. Genome Res 12:956–961

Kong A, Gudbjartsson DF, Sainz J, Jonsdottir GM, Gudjonsson SA, Richardsson B, Sigurdardottir S, Barnard J, Hallbeck B, Masson G, Shlien A, Palsson ST, Frigge ML, Thorgeirsson TE, Gulcher JR, Stefansson K (2002) A high-resolution recombination map of the human genome. Nat Genet 31:241–247

Ophoff R, Escamilla MA, Service SK, Spesny M, Meshi DB, Poon W, Molina J, Fournier E, Gallegos A, Mathews C, Neylan T, Batki SL, Roche E, Ramirez M, Silva S, DeMille MC, Dong P, Leon PE, Reus VI, Sandkuijl LA, Freimer NB (2002) Genomewide linkage disequilibrium mapping of severe bipolar disorder in a population isolate. Am J Hum Genet 71:565–574

Patil N, Berno AJ, Hinds DA, Barrett WA, Doshi JM, Hacker CR, Kautzer CR, Lee DH, Marjoribanks C, McDonough DP, Nguyen BT, Norris MC, Sheehan JB, Shen N, Stern D, Stokowski RP, Thomas DJ, Trulson MO, Vyas KR, Frazer KA, Fodor SP, Cox DR (2001) Blocks of limited haplotype diversity revealed by high-resolution scanning of human chromosome 21. Science 294:1719–1723

Purcell S, Cherny S, Sham P (2003) Genetic power calculator: design of linkage and association genetic mapping studies of complex traits. Bioinformatics 19:149–150

Reich DE, Cargill M, Bolk S, Ireland J, Sabeti PC, Richter DJ, Lavery T, Kouyoumjian R, Farhadian SF, Ward R, Lander ES (2001) Linkage disequilibrium in the human genome. Nature 411:199–204

Risch N, Merikangas K (1996) The future of genetic studies of complex human diseases. Science 273:1516–1517

Sabatti C, Service SK, Freimer NB. False discovery rate in linkage and association genome screens for complex disorders. Genetics (in press)

Schork NJ (2002) Power calculations for genetic association studies using estimated probability distributions. Am J Hum Genet 70:1480–1489

Service SK, Ophoff RA, Freimer NB (2001) The genomewide distribution of background linkage disequilibrium in a population isolate. Hum Mol Genet 10:545–551

Vaessen N, Heutink P, Houwing-Duistermaat JJ, Snijders PJ, Rademaker T, Testers L, Batstra MR, Sandkuijl LA, van Duijn CM, Oostra BA (2002) A genomewide search for linkage-disequilibrium with type 1 diabetes in a recent genetically isolated population from the Netherlands. Diabetes 51:856–859

Varilo T, Paunio T, Parker A, Perola M, Meyer J, Terwilliger JD, Peltonen L (2003) The interval of linkage disequilibrium (LD) detected with microsatellite and SNP markers in chromosomes of Finnish populations with different histories. Hum Mol Genet 12:51–59