## Supplementary I: theoretical results

The theoretical mean of T can be decomposed as  $\mu = E(L) + E(M) + E(R)$ . Because M is an indicator random variable,  $E(M) = \Pr(M = 1)$  and  $\operatorname{Var}(M) = \Pr(M)[1 - \Pr(M)]$ . If we let  $X_k$  be the unordered genotype of SNP k, then it is natural to calculate E(L) as  $E[E(L | X_0)]$ . Since  $X_0$  takes only three possible values, the outer expectation in  $E[E(L | X_0)]$  is trivial to compute. The case  $X_0 = 1/2$  is easiest of all because L = 0 when  $X_0 = 1/2$  and M = 0. Similar comments apply to E(R). The most natural route to calculating the variance  $\sigma^2$  follows the formula

$$\operatorname{Var}(T) = \operatorname{Var}(L) + \operatorname{Var}(M) + \operatorname{Var}(R) + 2\operatorname{Cov}(L, M) + 2\operatorname{Cov}(L, R) + 2\operatorname{Cov}(M, R).$$

Again it is productive to condition on  $X_0$ . For instance,

$$\operatorname{Var}(L) = \operatorname{Var}[\operatorname{E}(L \mid X_0)] + \operatorname{E}[\operatorname{Var}(L \mid X_0)],$$
  
$$\operatorname{Var}(R) = \operatorname{Var}[\operatorname{E}(R \mid X_0)] + \operatorname{E}[\operatorname{Var}(R \mid X_0)],$$

and, assuming L and R are independent given  $X_0$ ,

$$Cov(L, R) = Cov[E(L \mid X_0), E(R \mid X_0)] + E[Cov(L, R \mid X_0)] = Cov[E(L \mid X_0), E(R \mid X_0)].$$

It is also worth pointing out that E(LM) = E(L) and E(RM) = E(R), since L and R equal 0 when M does, and when M = 1, LM equals L and RM equals R. Thus, one has

$$Cov(L, M) = E(LM) - E(L)E(M) = E(L)[1 - E(M)],$$
  

$$Cov(R, M) = E(RM) - E(R)E(M) = E(R)[1 - E(M)].$$

These considerations emphasize the importance of finding the distributions of L and R conditional on  $X_0 = 1/1$  and  $X_0 = 2/2$ . The next few sections tackle this issue.

To compute the conditional means and variances of L and R numerically, we recommend the right-tail sums

$$E(Y) = \sum_{j=1}^{\infty} \Pr(Y \ge j), \quad E(Y^2) = \sum_{j=1}^{\infty} (2j-1) \Pr(Y \ge j),$$
 (1)

valid for any nonnegative random variable Y with integer values; see Section XI.1 of (Feller, 1968). The sums defining E(Y) and  $E(Y^2)$  can be truncated as soon as they stabilize. If we define  $h_{i_0,\ldots,i_r}$  to be the population frequency of the haplotype  $(i_0,\ldots,i_r)$  extending from SNP 0 to SNP r, then the formula

$$\Pr(R \ge r \mid X_0 = i_0/i_0) = \frac{1}{p_{0,i_0}^2} \sum_{i_1=1}^2 \cdots \sum_{i_r=2}^2 h_{i_0,\dots,i_r}^2$$

delivers the required right-tail probabilities. Here  $p_{0,i_0}$  is the frequency of allele  $i_0$  of the core SNP 0. When all conceivable haplotypes are possible, there are  $2^r$  terms in the multiple sum, and the formula as it stands is cumbersome. On the other hand, if only a few haplotypes are possible, then the sum is straightforward to evaluate. The moment formulas (1) are still applicable. The haplotype frequencies  $h_{i_0,...,i_r}$  can be estimated from haplotype data.