

Supplemental Material:

Statistical physics methods provide the exact solution to a long-standing problem of genetics

Areejit Samal^{1,2,3} and Olivier C. Martin^{4,*}

¹*Laboratoire de Physique Théorique et Modèles Statistiques (LPTMS),
CNRS and Univ Paris-Sud, UMR 8626, F-91405 Orsay, France*

²*Max Planck Institute for Mathematics in the Sciences, 04103 Leipzig, Germany*

³*The Abdus Salam International Centre for Theoretical Physics, Trieste 34151, Italy*

⁴*INRA, UMR 0320 / UMR 8120 Génétique Quantitative et Evolution – Le Moulon, F-91190 Gif-sur-Yvette, France*

I. PRODUCTION OF A RECOMBINANT INBRED LINE VIA SINGLE SEED DESCENT

Assume given two homozygous diploid parents P_a and P_A at the F0 generation. Without loss of generality, label the L loci or markers of interest by $1, 2, \dots, L$. We denote by a_1, a_2, \dots, a_L the alleles of P_a , and by A_1, A_2, \dots, A_L the alleles of P_A . The first generation of offspring or F1 individuals are produced by crossing P_a and P_A (Fig. 1 in Main Text and Fig. S1). Note that the F1 individuals are all identical and are heterozygous at each locus. Thus, the genotype of F1 individuals is $\{a_1/A_1, a_2/A_2, \dots, a_L/A_L\}$. The construction of the next generation (F2) depends on whether individuals can be selfed or not. Most plants are hermaphrodites, the same individual being capable of producing both male and female gametes. Such plants can be selfed to produce offspring for the next generation, a process referred to as *single seed descent* (SSD) and illustrated in Fig. 1 in Main Text. For animals, it is necessary to cross brothers and sisters to produce offspring, and this is referred to as *sib* mating. The present work concerns SSD, the sib case being significantly more complex.

Each individual arising during the successive generations (F1, F2, F3, ...) has a genomic content corresponding to the union of two gametes produced within its progenitor: one via female meiosis and the other via male meiosis. These gametes often involve crossovers that mix alleles within chromosomes. For example, the F2 genotype in Fig. 1 in Main Text is $\{a_1/a_1, a_2/A_2, A_3/a_3\}$ and so the bottom chromosome is recombined for both intervals (1,2) and (2,3) due to the occurrence of a crossover in each interval. Recombination occurs during a meiosis if there are an odd number of crossovers between the 2 loci under consideration and as a result the interval (1,3) of the example given is not recombinant. The probability that a recombination occurs between locus i and locus j is referred to as the (meiotic) recombination rate $r_{i,j}$ for that pair of loci. Crossovers form stochastically and their statistics has to be modeled. For pedagogical reasons, we follow standard practice and consider that female and male meioses are described by the same stochastic process and so that in particular female and male recombination rates are identical. Nevertheless, our framework is easily extended to the case of distinct female and male recombination rates (see Sections II and VII in Supplemental Material). Many models have been proposed to describe the statistics of crossover formation. In the simplest model, crossovers arise as independent events in each meiosis, a hypothesis due to Haldane [1]. Other models take the crossovers to exhibit *interference* with close-by crossovers being very rare. Our framework allows any kind of crossover formation model to be treated since model dependencies are restricted to the probabilities $P(g)$ and $P(g')$ in Eq. 3 in Main Text. However, it is only in the case of no interference that the analytical calculations (using the mappings to trees) can be pushed very far.

If the L loci are not physically linked, the calculation of the probabilities of genotypes at successive generations becomes trivial because the allelic content at each locus is passed on independently, corresponding to $r_{i,j} = 1/2$. The whole complexity of finding the probabilities of multi-locus genotypes stems from the linkage between loci, *i.e.*,

$$r_{i,j} \neq 1/2$$

Thus, without loss of generality, we assume that all L loci are on the same chromosome. After an F2 individual is produced, it is used to produce an F3 individual, which itself is used to produce an F4 individual, and so forth. If a locus becomes homozygous at one generation, it will remain *fixed* (neglecting mutations) in all future generations. If a locus is heterozygous at one generation, the probability that it will remain heterozygous at the next generation is $1/2$. Thus, with the increase in the number of generations, more loci become homozygous and fixed. After a large

*olivier.martin@moulon.inra.fr

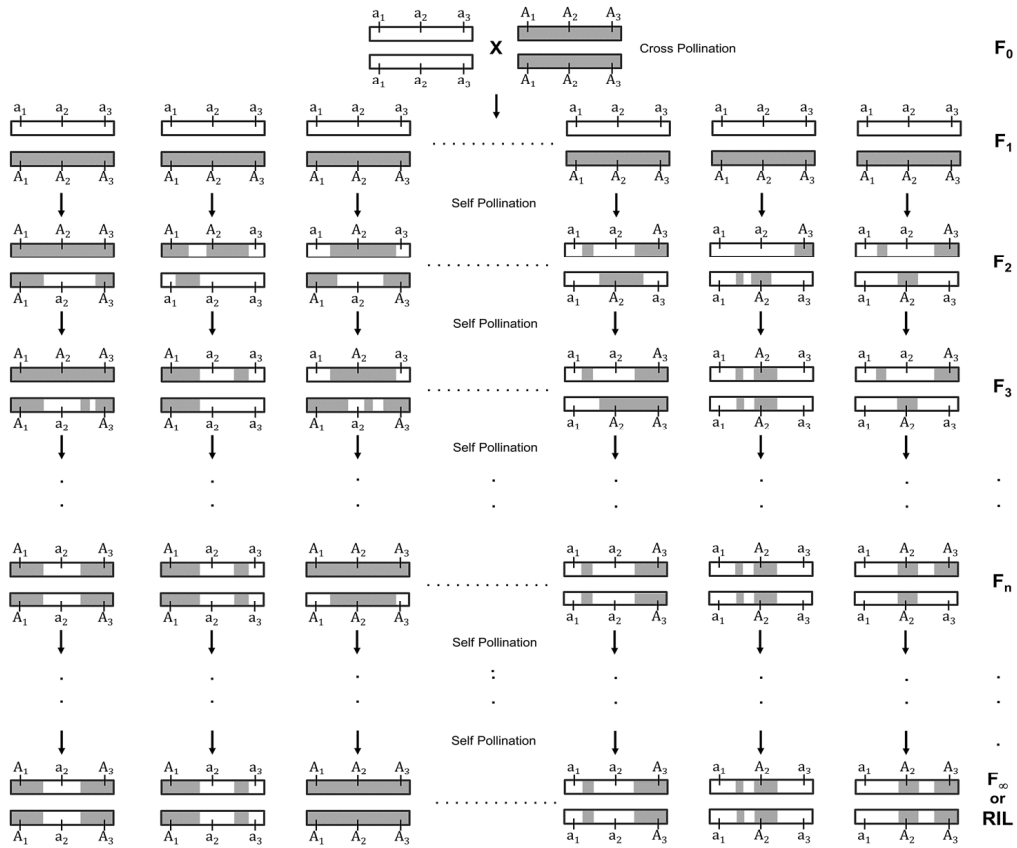


FIG. S1. **Many Recombinant Inbred Lines produced in parallel using Single Seed Descent.** Two homozygous parents are crossed to produce the F₁ generation of genetically identical individuals. Thereafter, at each generation, each plant produces one gamete via female meiosis and one gamete via male meiosis, and then these two gametes are fused to produce the genomic content of the individual of the next generation. Crossovers may arise during the meioses, leading to intra-chromosomal shuffling of allelic content. After enough generations, all loci become homozygous.

number of generations, all alleles will become fixed (Fig. 1 in Main Text). If this SSD process is performed in parallel for a number of *lines* as illustrated in Fig. S1, one obtains a population of recombinant inbred lines (RILs) where each genome is a homozygous mosaic of the two parental genomes. The different RIL individuals are inbred and any pair of loci may have recombined the alleles of the initial parents P_a and P_A , thus the term RILs.

II. REDERIVING THE HALDANE-WADDINGTON FORMULA VIA THE NEW FRAMEWORK AND THE CASE OF SEX-SPECIFIC RECOMBINATION RATES

In 1918, Robbins [2] determined the probabilities of RIL genotypes produced using single seed descent (SSD) for the case of 2 loci. In 1931, Haldane and Waddington [3] reconsidered the question using a simpler method and went on to solve the problem when using sib mating. Furthermore, Haldane and Waddington [3] also showed that the 2-locus RIL probabilities completely determine the 3-locus RIL probabilities. Here we show how our framework can simplify the derivation of the Haldane-Waddington formula for 2 loci in the SSD case, thus illustrating, albeit on a very simple case, the logic of the Schwinger-Dyson (SD) approach, approach that generalizes to many more loci.

To rederive the (2-locus) Haldane-Waddington formula, we start with the $2^2 = 4$ possible RIL genotypes, $\{a_1/a_1, a_2/a_2\}$, $\{a_1/a_1, A_2/A_2\}$, $\{A_1/A_1, a_2/a_2\}$, and $\{A_1/A_1, A_2/A_2\}$ where the top (bottom) allele specified at a locus is that on the chromosome generated during female (male) meiosis. In our spin notation, these homozygous genotypes are denoted as $\vec{S} = \{1, 1\}$, $\{1, -1\}$, $\{-1, 1\}$ and $\{-1, -1\}$, respectively. The RIL recombination rate R is defined as the probability of having *recombinant* genotypes: $R = P(\{1, -1\}) + P(\{-1, 1\})$. R is related to the

expectation of the 2-spin product over all RIL genotypes with their respective probabilities via:

$$E[S_1S_2] = P(\{1, 1\}) + P(\{-1, -1\}) - P(\{1, -1\}) - P(\{-1, 1\}) = 1 - 2R \quad (5)$$

The difficulty in determining R comes from the fact that producing RILs involves in principle an infinite number of generations. The heart of our method consists in transforming such an infinite process into a finite one based on self-consistent equations as follows. The probability of a RIL genotype is associated with sums over all possible meioses across generations F1, F2, ... leading to that RIL genotype. Now think of classifying these *trajectories* according to the genotype produced at generation F2. In our framework, we must calculate the probability of each F2 genotype and the contribution of associated trajectories to $E[S_1S_2]$. There are 4^2 F2 genotypes and the probability of each is easy to compute.

Consider for instance the F2 genotype $\{a_1/A_1, A_2/A_2\}$ which occurs with probability $r(1-r)/4$ where $r = r_{1,2}$ is the recombination rate (for one meiosis) between the 2 loci. How much do trajectories passing through that F2 genotype contribute to $E[S_1S_2]$? Clearly, since the second locus is fixed to type ‘‘A’’, we have $S_2 = -1$ necessarily. Furthermore, the first locus will fix to either $S_1 = 1$ or $S_1 = -1$ with probability $1/2$ for each and summing these outcomes gives 0 for the expectation value of the 2-spin product. Thus trajectories passing through that F2 genotype contribute nothing to $E[S_1S_2]$. The same result will hold for all F2 genotypes that are heterozygous at one locus and homozygous at the other.

Consider then the F2 genotype $G = \{a_1/A_1, a_2/A_2\}$ arising with probability $P_G = (1-r)^2/4$. This genotype is identical to the F1 genotype, so its contribution is $P_G E[S_1S_2]$. The same result holds for the genotype $\{A_1/a_1, A_2/a_2\}$ because of the global invariance under exchange of all ‘‘a’’s for ‘‘A’’s and vice versa.

A bit more subtle is the case of the genotype $G' = \{a_1/A_1, A_2/a_2\}$ arising with probability $P_{G'} = r^2/4$. This case is *similar* to that of the F1 genotype except that the alleles at the second locus have been inverted. All trajectories produced from this genotype G' can be mapped to those produced from the F1 genotype if we perform the *substitution* of the alleles at the second locus, exchanging ‘‘a’’s and ‘‘A’’s. The probabilities of these substituted trajectories will be the same as before the substitution but when we consider the contribution of genotype G' in the RILs we have to also substitute $S_2 \rightarrow -S_2$. Thus the trajectories passing through the genotype G' contribute the amount $-P_{G'} E[S_1S_2]$. The same strategy applies to the genotype $\{A_1/a_1, a_2/A_2\}$ for which the required substitution is $S_1 \rightarrow -S_1$.

Lastly, there are F2 genotypes that are homozygous at both loci. Their contribution to $E[S_1S_2]$ is easily read off since each locus is fixed, and in fact the RIL fixation has been accomplished in just one generation.

Adding up the contributions associated with all 4^2 F2 genotypes gives the self-consistent equation:

$$E[S_1S_2] = \left[\frac{(1-r)^2}{2} - \frac{r^2}{2} \right] \times E[S_1S_2] + [2(1-r)r] \times 0 + \left[\frac{(1-r)^2}{2} - \frac{r^2}{2} \right] \times 1 \quad (6)$$

where we have ordered the terms according to F2 genotypes having 0, 1 and 2 fixed loci. This SD equation leads to $E[S_1S_2](1+2r) = 1-2r$ from which one obtains:

$$R = \frac{2r}{1+2r} \quad (7)$$

i.e., the Haldane-Waddington formula.

Robbins [2] and then Haldane and Waddington [3] also determined the probabilities of 2-locus SSD RIL genotypes in the case of sex-specific recombination rates, *i.e.*, where the female and male recombination rates are different. Interestingly, that generalization does not affect much our framework, the only modification arises in the probabilities of the F2 genotypes. Denoting the female and male recombination rates between the 2 loci by r^f and r^m , respectively, the generalization of Eq. 6 along with obvious simplifications gives:

$$E[S_1S_2] = \left[\frac{1-r^f-r^m}{2} \right] \times E[S_1S_2] + [(1-r^f)r^m + r^f(1-r^m)] \times 0 + \left[\frac{1-r^f-r^m}{2} \right] \times 1 \quad (8)$$

It is interesting to note that although the individual $P(G)$ s depend on both r^f and r^m , the above SD equation depends only on the mean of r^f and r^m because the middle factor is multiplied by 0. This property is not general, and in particular, we will show later that it does not hold for 4 loci (see Section VII in Supplemental Material). Solving for R leads to:

$$R = \frac{r^f + r^m}{1 + r^f + r^m} \quad (9)$$

III. PROBABILITIES FOR 2 AND 3 LOCI DO NOT DETERMINE THOSE FOR 4 LOCI

In their 1931 paper, Haldane and Waddington [3] derived the formula for 2-locus RIL probabilities using recursions from one generation to the next and then took the limit of an infinite number of generations. Furthermore, they provided a simple mathematical trick involving standard algebra to obtain 3-locus RIL probabilities from 2-locus RIL probabilities (Fig. 2a in Main Text). Thus, those authors showed that the 2-locus RIL probabilities completely determine the 3-locus RIL probabilities. However, the simple trick of Haldane and Waddington does not extend to the case of 4 loci. We now elucidate this difference between going from 2 to 3 loci versus from 3 to 4 loci, and show that 2- and 3-locus RIL probabilities do *not* determine the 4-locus probabilities (Fig. 2 in Main Text).

Let us first calculate the probabilities of all ($2^3 = 8$) 3-locus RIL genotypes. One can use symmetries such as $P(a_1/a_1, a_2/a_2, A_3/A_3) = P(A_1, A_1, A_2/A_2 a_3/a_3)$ to formulate the problem in terms of 4 unknowns $Q(0, 0)$, $Q(0, 1)$, $Q(1, 0)$ and $Q(1, 1)$ defined in Fig. 2a in Main Text. For these 4 quantities, the binary entry 0 (respectively 1) denotes absence (respectively presence) of a recombination event in the corresponding interval (1 or 2). To determine these 4 unknowns, one requires 4 independent equations. A first independent equation is that the sum of the 4 probabilities equals 1. Furthermore, three additional equations are obtained from the 2-locus RIL probabilities using $R_{1,2}$, $R_{1,3}$ and $R_{2,3}$. Having as many independent equations as unknowns (the rank of the matrix constraining the 4 unknowns is 4), one concludes that the 2-locus probabilities uniquely determine the 3-locus probabilities. The mathematics behind the 3-locus case is provided in Fig. 2a in Main Text.

We next ask if all 2- and 3-locus probabilities similarly determine the 4-locus probabilities. There are $2^4 = 16$ 4-locus RIL genotypes. Again one can use symmetries to formulate the problem in terms of 8 unknowns (see $Q(0, 0, 0)$ etc. defined in Fig. 2b in Main Text). To determine these 8 unknowns, one needs 8 independent equations. As before, one equation follows from the fact that the sum of the 8 probabilities equals 1. Furthermore, there are 6 equations associated with 2-locus constraints ($R_{1,2}$, $R_{1,3}$, $R_{1,4}$, $R_{2,3}$, $R_{2,4}$, $R_{3,4}$). We need one more independent equation to solve for all 8 unknowns. It is tempting to use one of the equations based on 3-locus constraints (Fig. 2b in Main Text). However all those equations are consequences of the 2-locus constraints: they are *automatically* satisfied and provide no further constraints on the unknowns. In the 4-locus case the rank of the matrix constraining the 8 unknowns is at most 7 (Fig. 2b in Main Text) regardless of which of the 3-locus constraints are added since these follow from the 2-locus constraints. In conclusion, to obtain all 4-locus RIL probabilities, *one additional piece of information is needed* that is not incorporated in 2- or 3-locus RIL probabilities. Finding and exploiting this missing information has prevented researchers from extending the Haldane-Waddington result for over 80 years.

IV. MAPPING ALL F2 GENOTYPES WITH 4 LOCI TO TREES

To derive the SD Eq. 4 in Main Text for $E[S_1 S_2 S_3 S_4]$, we classify the F2 genotypes according to whether they are homozygous (H) or heterozygous (h) at the different loci. There are 2^4 such classes and each class contains 2^4 genotypes because a homozygous (respectively, heterozygous) locus i can have the allelic content a_i/a_i or A_i/A_i (respectively, a_i/A_i or A_i/a_i). For illustration, consider all F2 genotypes belonging to the class $hhhh$. The first locus can be in the state a_1/A_1 or A_1/a_1 , the second in the state a_2/A_2 or A_2/a_2 , the third in the state a_3/A_3 or A_3/a_3 , and so on. This succession of possibilities can be represented by a binary tree whose root is associated with the state of the first locus. Thus, there are 2 trees for the $hhhh$ class: one rooted on the a_1/A_1 state (Fig. 3 in Main Text) and the other on the A_1/a_1 state. These two trees are related to each other: one goes from one tree to the other via a global exchange of “a”s into “A”s and vice versa. In the Main Text, we mentioned this *exchange invariance* at the level of RIL probabilities but in fact it also holds for probabilities of genotypes at any generation of the RIL construction. Each F2 genotype G can be identified with a path on its associated tree which goes from the root to a leaf of that tree. Furthermore, the probability of a genotype G , composed of its two gametes g and g' , is the product of the following terms if crossovers arise without interference: a factor $1/4$ for the root node and a factor $(1-r)^2$, $(1-r)r$, $r(1-r)$ or r^2 for each interval between adjacent loci depending on whether the interval is recombinant or not for g and for g' . Finally, each genotype G comes with a sign, denoted here by $\text{sign}(G)$, which arises from the substitution rules (see Main Text).

It is easy to enumerate all the classes to cover: $hhhh$, $hhhH$, $hhHh$, $hhHH$, and so on. Each class gives rise to two trees. However, just as in the example considered above, these two trees are related by the exchange invariance under swapping of “a”s and “A”s. Thus, it is enough to consider one tree per class and to multiply its contribution by 2 in the SD equation. Thus, without loss of generality, we show in all our figures only those trees which have at locus 1 of their female chromosome the a_1 allele. With this first simplification, the number of trees to be considered is 2^4 .

A second major simplification arises by noting that each class is associated with a different multi-spin product to average. For instance, if one considers the terms on the right-hand side of the SD Eq. 4 in Main Text, $hhhh$ is

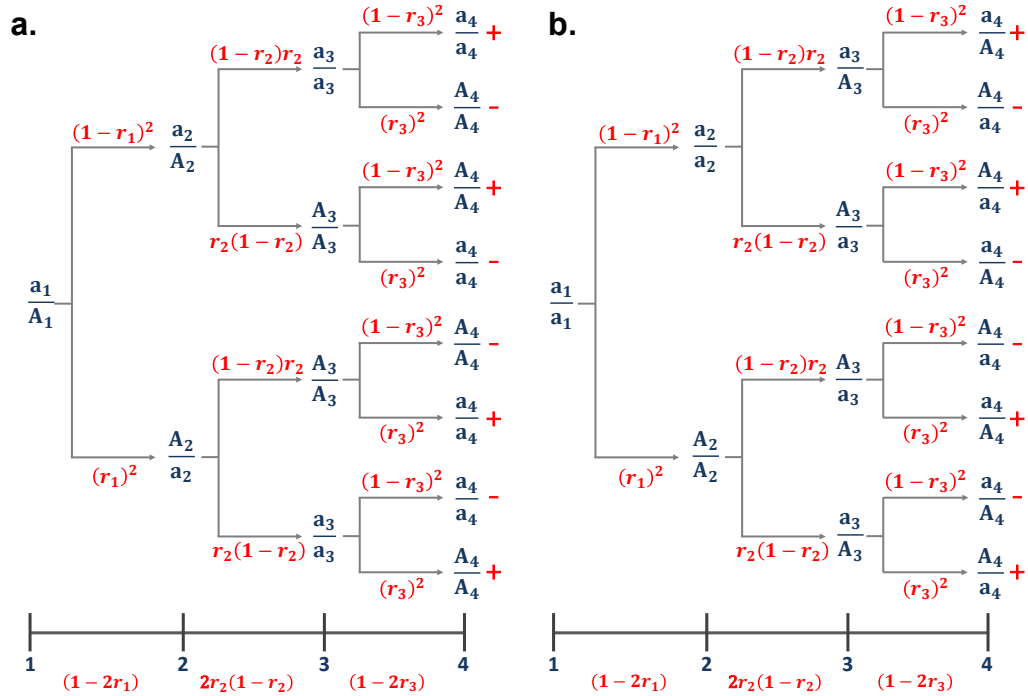


FIG. S2. **Four-locus trees for classes of F2 genotypes with 2 fixed loci leading to non-vanishing contributions in the Schwinger-Dyson equation.** (a) $hhHH$ and (b) $HhhH$.

associated with $E[S_1S_2S_3S_4]$, $hhhH$ with $E[S_1S_2S_3]$, $hhHh$ with $E[S_1S_2S_4]$, and so on. Because the probabilities of RIL genotypes are invariant under $\vec{S} \rightarrow -\vec{S}$, expectations are also invariant. Then if there is an odd number of spins in a spin-product, its expectation value vanishes. As a consequence, amongst the 2^4 classes, we only need to consider those with 0, 2 and 4 loci of type H .

Consider the tree for the class $hhhh$ rooted at a_1/A_1 (Fig. 3 in Main Text). The *sign* carried by a genotype is specified on the leaf of the path representing G on its tree. By summing $\text{sign}(G)P(G)$ over all genotypes G belonging to this tree, and multiplying by 2 to take into account the other tree for this class (*i.e.*, the tree rooted at A_1/a_1), we obtain the factor in the right-hand side of the SD equation associated with $E[S_1S_2S_3S_4]$. To derive the formula for this sum, we start with the right-most of the three intervals and collect the paths into pairs that differ only in this last interval. This pools together contributions of double recombinants and double non-recombinants with opposite signs, leading to the factor $(1-r_3)^2 - r_3^2 = 1 - 2r_3$ and a sign that depends on the pair. The factor for the third interval is given at the bottom of the tree in Fig. 3 in Main Text. An important point is that this factor $1 - 2r_3$ is common to all pairs of paths which differ only in the last interval, and so these pairs can be identified with shortened paths restricted to just the first two intervals. This property allows us to *iterate* the procedure. Thus we consider now all (shortened) paths covering just the first two intervals and pair these up if they differ only on the second interval. Again, the pairing requires pooling the contributions of double recombinants and double non-recombinants. Because the two paths to be added in the pair both have the same sign (which was not true for the third interval), the common factor for the second interval is $(1-r_2)^2 + r_2^2$ (Fig. 3 in Main Text). This pooling leaves us with just two shortened paths of one segment with opposite signs for the first interval (Fig. 3 in Main Text). Thus, adding the contributions of these two paths we obtain the factor for the first interval $(1-r_1)^2 - r_1^2 = 1 - 2r_1$ (Fig. 3 in Main Text). One final factor must be included: the probability of having the first locus in the given (a_1/A_1) state, *i.e.*, $1/4$. The resulting product is this tree's contribution to $E[S_1S_2S_3S_4]$, coming from its 8 F2 genotypes. There are also the other 8 F2 genotypes of the class $hhhh$ associated with a tree rooted on A_1/a_1 which leads to exactly the same result as can be seen either by direct calculation or by using the previously mentioned exchange invariance under global swaps of “a”s and “A”s. Putting all this together, the factor in front of $E[S_1S_2S_3S_4]$ on the right-hand side of Eq. 4 in Main Text is:

$$A_{1,1,1,1} = \frac{(1-2r_1) [(1-r_2)^2 + r_2^2] (1-2r_3)}{2} \quad (10)$$

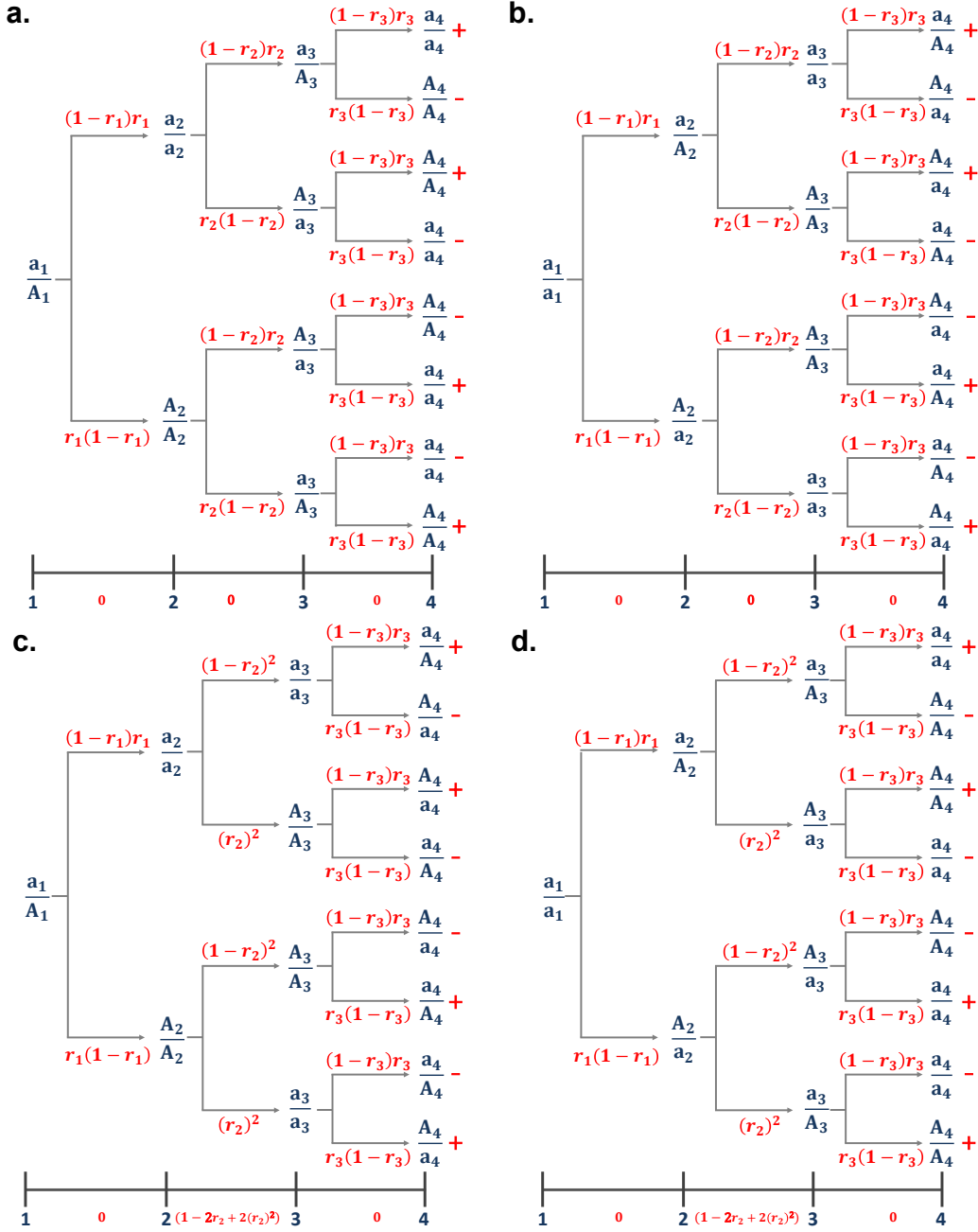


FIG. S3. Four-locus trees for classes of F2 genotypes with 2 fixed loci leading to vanishing contributions in the Schwinger-Dyson equation. (a) $hHhH$, (b) $HhHh$, (c) $hHHh$ and (d) $HhhH$.

where the indices of $A_{1,1,1,1}$ refer to the powers n_i arising in the associated spin product $E[S_1^{n_1} S_2^{n_2} S_3^{n_3} S_4^{n_4}]$.

Suppose one repeats the calculation that led to $A_{1,1,1,1}$ but replaces the male chromosome in each G by a modified one where all “a”s have been exchanged for “A”s and vice versa. This transformation takes one from a heterozygous G to a homozygous G' . Interestingly, this transformation affects neither the probabilities arising in each interval nor the signs ($\text{sign}(G) = \text{sign}(G')$). Thus, $A_{0,0,0,0} = A_{1,1,1,1}$. Furthermore, it is easy to see that this invariance applies to any of the F2 genotypes. As a result, for any choices of the n_i , ($n_i = 0$ or 1), $A_{n_1, n_2, n_3, n_4} = A_{1-n_1, 1-n_2, 1-n_3, 1-n_4}$, providing the third major simplification and reduction in the set of trees to be considered.

Finally we are left with the mixed cases where G has two homozygous loci and two heterozygous loci. Let us begin with the class $hhHH$ which determines the factor $A_{1,1,0,0}$ (Fig. S2a in Supplemental Material). The contributions

from the associated F2 genotypes can be combined just as in the calculation of $A_{1,1,1,1}$: the third interval again leads to the factor $(1 - 2r_3)$; the second interval has one recombinant gamete and one non-recombinant gamete, leading to the factor $2(1 - r_2)r_2$; and finally the first interval again leads to the factor $(1 - 2r_1)$. Using the fact that two trees contribute to the class $hhHH$ and the invariance result from the previous paragraph for deducing $A_{0,0,1,1}$, (Fig. S2b in Supplemental Material), we obtain:

$$A_{1,1,0,0} = A_{0,0,1,1} = \frac{(1 - 2r_1)(2(1 - r_2)r_2)(1 - 2r_3)}{2} \quad (11)$$

The other mixed cases lead to an even simpler result. Consider for instance the class $hHhH$ and the associated tree (Fig. S3a in Supplemental Material). When paths differing only on the last interval are paired, the interval factor for each genotype is $(1 - r_3)r_3$ but the signs are opposite and so the sum vanishes. The same is true for the remaining classes ($HhHh$, $hHHh$, and $HhhH$), and thus $A_{1,0,1,0} = A_{0,1,0,1} = A_{1,0,0,1} = A_{1,0,0,1} = 0$ (Fig. S3 in Supplemental Material).

Thus, collecting the results from all classes of F2 genotypes leads to the 4-locus SD equation (Eq. 4 in Main Text) where the expectation of the 4-spin product arises on both sides of this equation. Using the formula for the expectation of the 2-spin product: $E[S_i S_j] = 1 - 2R_{i,j} = (1 - 2r_{i,j})/(1 + 2r_{i,j})$ in Eq. 4 of Main Text one obtains the expectation of the 4-spin product as:

$$E[S_1 S_2 S_3 S_4] = \frac{F + A}{1 - F} \quad (12)$$

where the terms F and A in this equation are given by:

$$F = \frac{(1 - 2r_1)((1 - r_2)^2 + r_2^2)(1 - 2r_3)}{2} \quad \text{and} \quad A = \frac{(1 - 2r_1)(2(1 - r_2)r_2)(1 - 2r_3)(1 - 4r_1r_3)}{(1 + 2r_1)(1 + 2r_3)}. \quad (13)$$

The probabilities of RIL genotypes can be obtained by substituting the expectations of the 2-spin products given by $E[S_i S_j] = 1 - 2R_{i,j} = (1 - 2r_{i,j})/(1 + 2r_{i,j})$ and the expectation of the 4-spin product from Eq. 12 into Eq. 2 in Main Text. For instance, the probability of the four-locus genotype $\{a_1/a_1, a_2/a_2, a_3/a_3, a_4/a_4\}$ is

$$P(\{1, 1, 1, 1\}) = \frac{1 + \sum_{i < j} \frac{1 - 2r_{i,j}}{1 + 2r_{i,j}} + \frac{F + A}{1 - F}}{16} \quad (14)$$

V. USEFUL PROPERTIES FOR SIMPLIFYING THE DERIVATION OF THE SCHWINGER-DYSON EQUATIONS

In the 4-locus case we made use of a number of identities to reduce the number of F2 genotypes that had to be considered in the SD equation. Here we make such properties explicit for the general case of any number of loci and also introduce one additional invariance.

Rule 1: For each class of F2 genotypes (denoted by a succession of L letters in $\{H, h\}$), there are two associated trees: one with allele a_1 and the other with allele A_1 at locus 1 for the female chromosome. In fact the two trees lead to the same contribution to the SD equation. So, in practice one can force the allele at locus 1 for the female chromosome to be a_1 , reducing by a factor 2 the number of trees to be considered.

Rule 2: For a given class of F2 genotypes, the spin product $E[S_1^{n_1} S_2^{n_2} \dots S_L^{n_L}]$ generated in the SD equation has $n_i = 1$ if the locus i is of type h and $n_i = 0$ if the locus i is of type H . The number of spins in the spin product is then equal to the number of heterozygous loci in the class. Given the invariance of expectations values under the change of sign of all spins, the expectation value of a k -spin product vanishes when k is odd. Thus a second simplification consists in keeping only the classes of genotypes having an even number of h 's, again reducing by a factor of 2 the number of trees to be considered.

Rule 3: A further useful property is *chromosome choice invariance*. Consider exchanging “a”s and “A”s on just *one* of the chromosomes of an F2 genotype. In terms of meiosis, this corresponds to exchanging the two (F1) parental chromosomes when producing that gamete. In terms of the classes of F2 genotypes, it leads to the global swap of H s and h s, taking one class to a transformed one. A tree of the first class is transformed to a tree of the second class but the probabilities and signs are left invariant. However at the level of spin products, the transformation changes $n_i = 1$ into $n_i = 0$ and vice versa. As a result, factors in the SD equation come in equal pairs, for example $A_{1,0,0,1,0,0} = A_{0,1,1,0,1,1}$, reducing again by a factor 2 the number of trees to be considered. Note that if one applies chromosome choice invariance successively to both the male and the female chromosomes, *all* “a”s and “A”s are exchanged; then the class considered (list of H s and h s) is invariant but the allele at the first locus changes from a_1

to A_1 , leading to Rule 1 which is thus a special case of Rule 3. Collecting these results, there are always four trees that produce exactly the same factors (albeit multiplying different expectation values in the SD equation), these trees being rooted on a_1/A_1 , A_1/a_1 , a_1/a_1 and A_1/A_1 .

Rule 4: For a class of F2 genotypes to lead to a non-zero contribution in a SD equation, both the h loci and the H loci must come in adjacent pairs. To see why this is the case, consider a class of F2 genotypes in which there is a block of adjacent H loci, delimited by h loci, and let G be one genotype in this class. In the left interval bounding this block, one of the chromosomes of G is recombinant, the other not. The same property holds for the right interval bounding this block. Consider now the F2 genotype G' identical to G in terms of crossover locations except that for these two intervals we exchange which is the chromosome (female or male) that is recombinant. This transformation does not affect the probability of the genotype, but $\text{sign}(G') = -\text{sign}(G)$ if and only if the size of the H block is odd. The contribution of G' thus cancels exactly that of G in such a situation. This still holds if the block of H s has only one interval bounding it (*i.e.*, it goes to an end of the chromosome). And by symmetry, the whole argument can be repeated when considering blocks of h s instead of blocks of H s. This fourth rule was used while studying the 4-locus case, but it is completely general and greatly reduces the number of classes to consider when there are many loci.

VI. THE SCHWINGER-DYSON EQUATIONS FOR 6 LOCI AND BEYOND

Using the SD framework along with the simplification rules listed in Section V in Supplemental Material, the following four classes must be considered in the case of 6 loci: $hhhhhh$ that gives the factor for the $E[S_1S_2S_3S_4S_5S_6]$ and 1 terms; $hhhhHH$ that gives the factor for the $E[S_1S_2S_3S_4]$ and $E[S_5S_6]$ terms; $HHhhhh$ that gives the factor for the $E[S_3S_4S_5S_6]$ and $E[S_1S_2]$ terms; $hhHHhh$ that gives the factor for the $E[S_1S_2S_5S_6]$ and $E[S_3S_4]$ terms.

Consider the class $hhhhhh$ with the tree rooted at a_1/A_1 (Fig. S4 in Supplemental Material). Just as in the 4-locus case, an F2 genotype associated with this tree corresponds to a path from the root of the tree to one of the leaves of the tree. To determine the sum of $\text{sign}(G)P(G)$ over the F2 genotypes, we again start at the right-most (fifth) interval and collect into pairs the paths that differ only for that last interval. The calculation is identical to that performed for the 4-locus case and one obtains the factor $(1 - r_5)^2 - r_5^2$. Similarly (and not surprisingly in view of how the calculation proceeded in the 4-locus case), the fourth interval leads to the factor $(1 - r_4)^2 + r_4^2$. After having treated those two intervals, we see that the remaining paths correspond to a 4-locus tree that is identical with the one for the $hhhh$ class on loci 1 to 4 (Fig. 3 in Main Text). Thus the 6-locus tree for the $hhhhhh$ class gives a factor that is the product of $(1 - 2r_5)$, of $[(1 - r_4)^2 + r_4^2]$, and of the previously derived factor for the tree for the $hhhh$ class on loci 1 to 4, so that

$$A_{1,1,1,1,1,1} = \frac{(1 - 2r_1) [(1 - r_2)^2 + r_2^2] (1 - 2r_3) [(1 - r_4)^2 + r_4^2] (1 - 2r_5)}{2} \quad (15)$$

Consider next the class $hhhhHH$ and its tree rooted at a_1/A_1 (Fig. S5 in Supplemental Material). Proceeding as before, we pool together the paths that differ only in the last interval, leading to the common factor $(1 - r_5)^2 - r_5^2$. Moving on to the fourth interval, we see that it takes one from a locus of type h to a locus of type H , leading to the factor $2(1 - r_4)r_4$. After this, the remaining tree is identical with the one for the $hhhh$ class on loci 1 to 4 (Fig. 3 in Main Text), just as in the previous paragraph. From this we conclude that the 6-locus tree for the $hhhhHH$ class gives a factor that is the product of $(1 - 2r_5)$, of $[2(1 - r_4)r_4]$, and of the previously derived factor for the tree for the $hhhh$ class on loci 1 to 4, and thus

$$A_{1,1,1,1,0,0} = \frac{(1 - 2r_1) [(1 - r_2)^2 + r_2^2] (1 - 2r_3) [2(1 - r_4)r_4] (1 - 2r_5)}{2} \quad (16)$$

Moving on to the class $HHhhhh$, the pooling over the last two intervals of the tree (Fig. S6 in Supplemental Material) leads to the same factors as those obtained for the $hhhhhh$ class. After this, the remaining tree is identical with the one for the $HHhh$ class on loci 1 to 4 (Fig. S2b in Supplemental Material). Of course, the same result could also have been obtained from the formula for the class $hhhhHH$ by taking the convention that loci are ordered from right to left rather than from left to right. This gives

$$A_{0,0,1,1,1,1} = \frac{(1 - 2r_1) [2(1 - r_2)r_2] (1 - 2r_3) [(1 - r_4)^2 + r_4^2] (1 - 2r_5)}{2} \quad (17)$$

Finally, consider the tree associated with the last class $hhHHhh$ (Fig. S7 in Supplemental Material). The pooling over the last two intervals of the tree leads to the factors $(1 - 2r_5)$ and $2(1 - r_4)r_4$. After this, the remaining tree is identical to the one for the $hhHH$ class on loci 1 to 4 (Fig. S2a in Supplemental Material). As a result,

$$A_{1,1,0,0,1,1} = \frac{(1 - 2r_1) [2(1 - r_2)r_2] (1 - 2r_3) [2(1 - r_4)r_4] (1 - 2r_5)}{2} \quad (18)$$

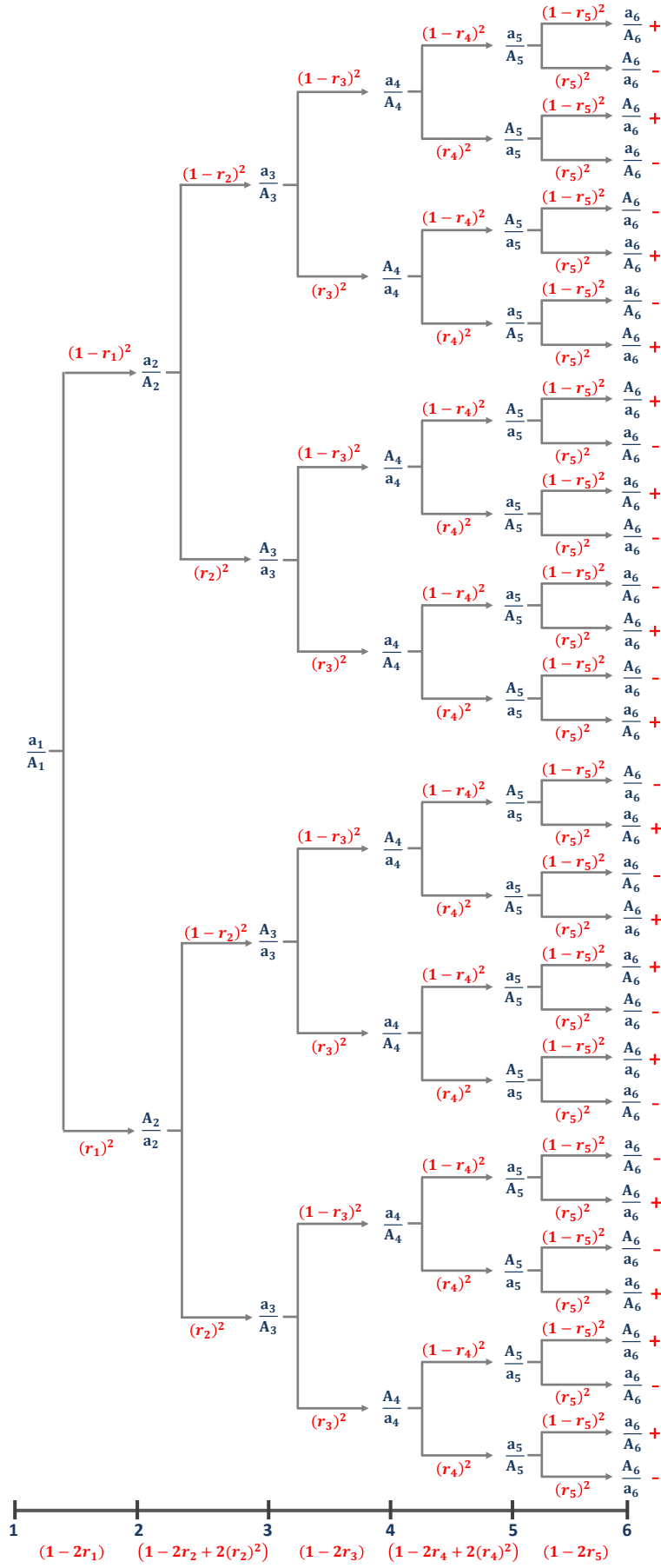
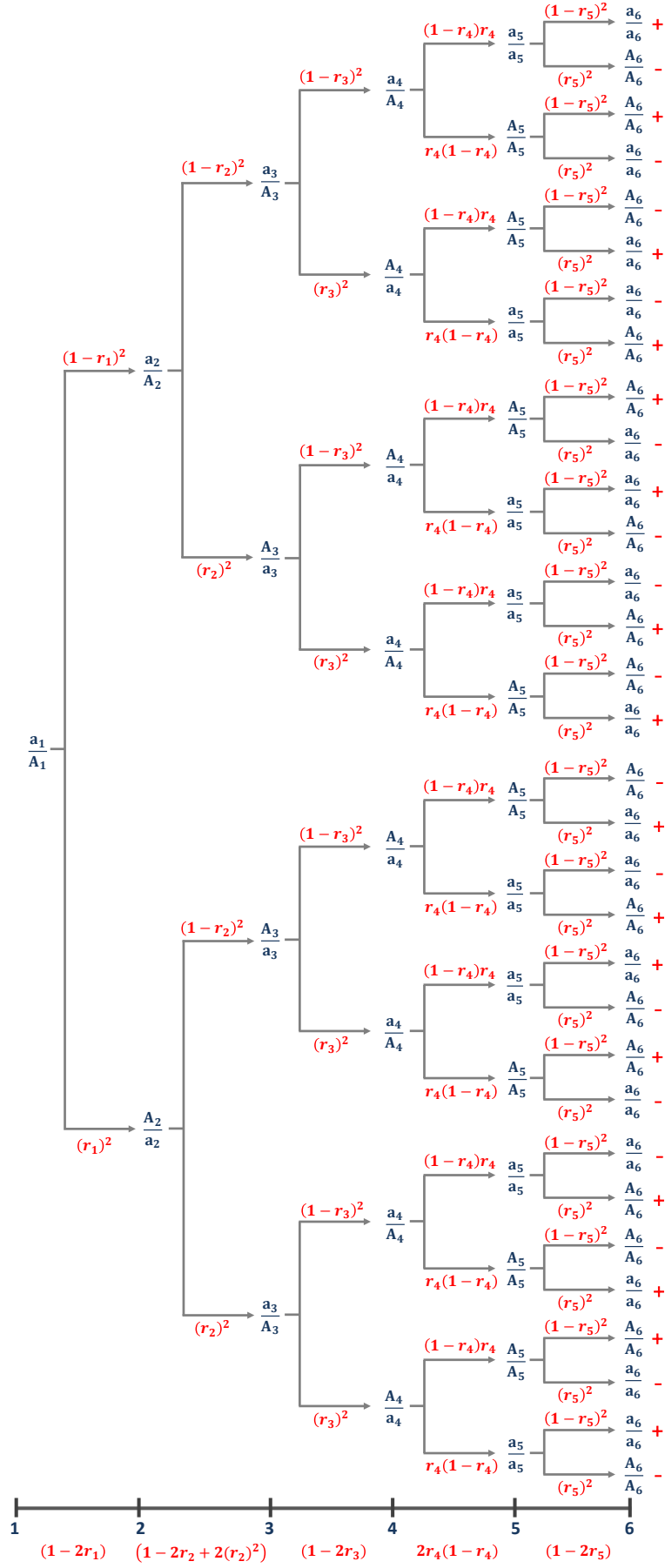


FIG. S4. Six-locus tree for the class $hhhhhh$ of F2 genotypes.

FIG. S5. Six-locus tree for the class $hhhhHH$ of F2 genotypes.

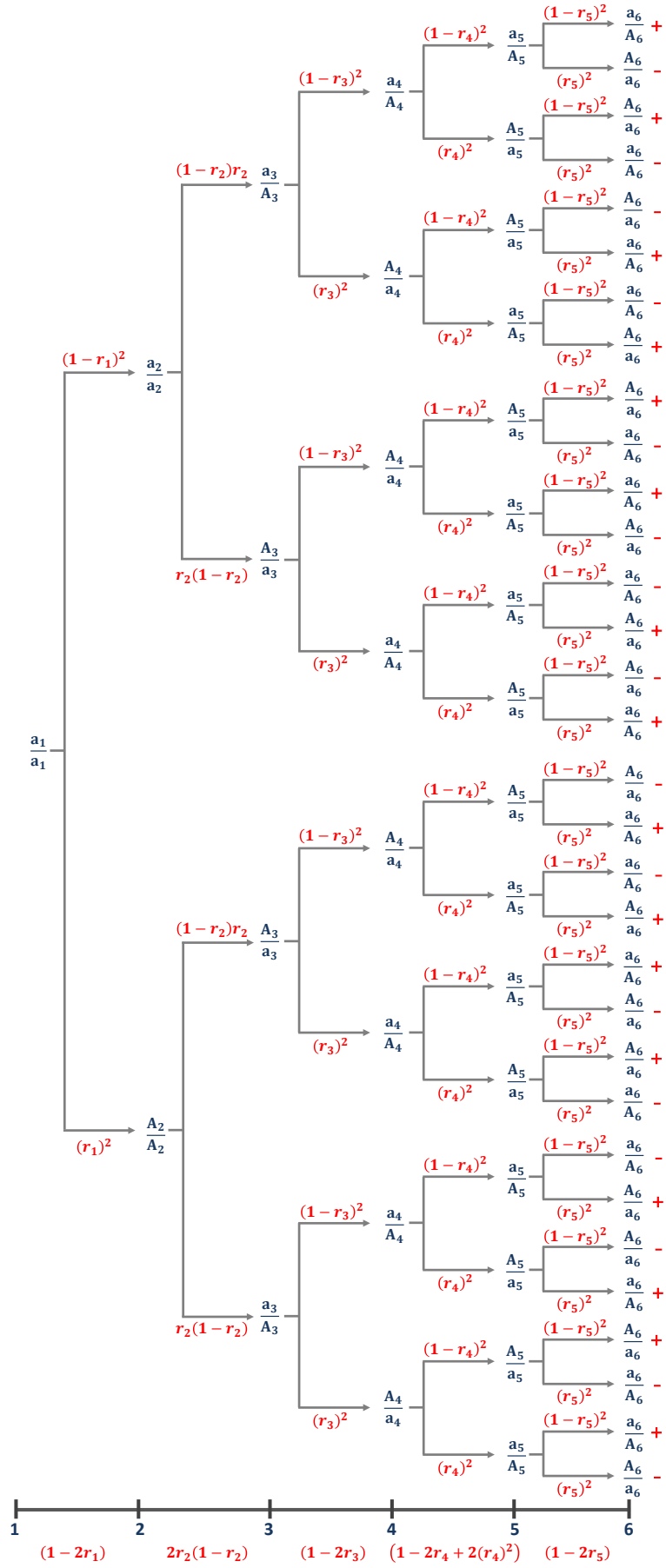
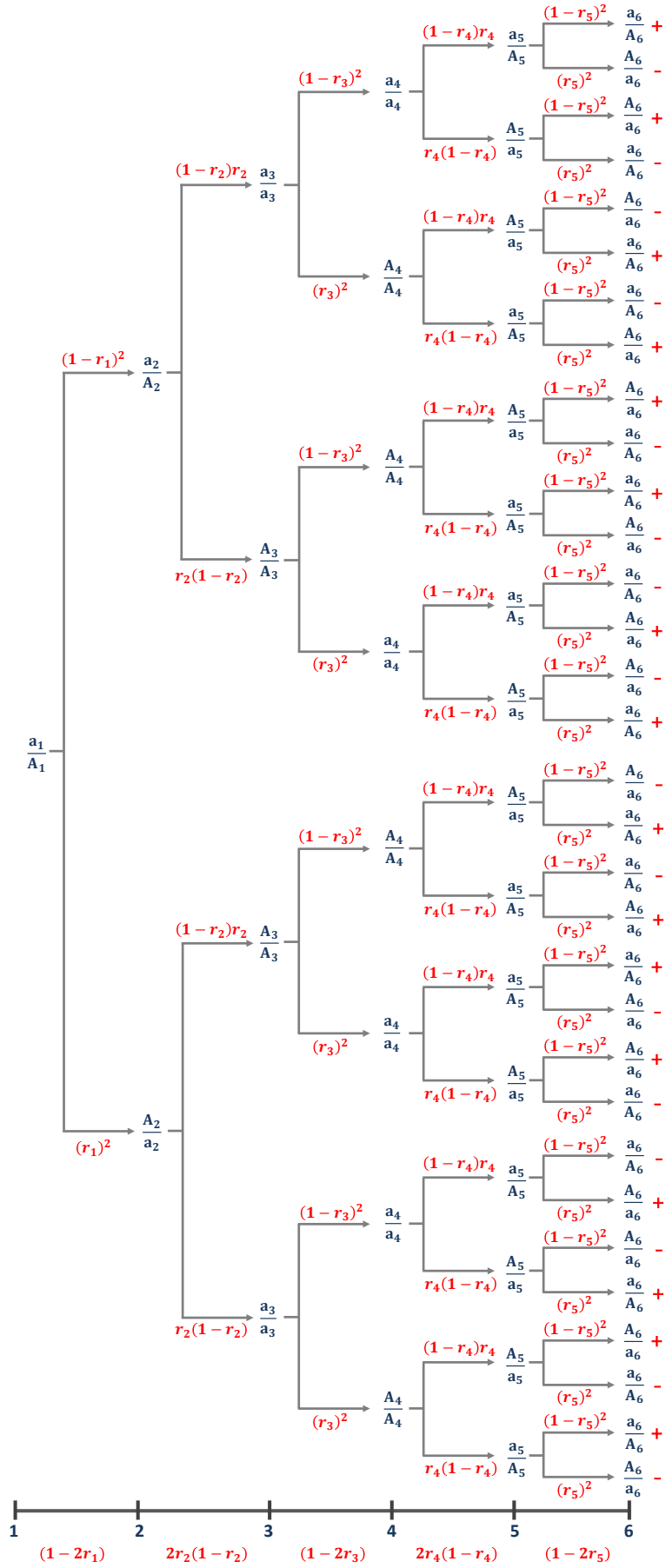


FIG. S6. Six-locus tree for the class $HHhhhh$ of F_2 genotypes.

FIG. S7. Six-locus tree for the class $hhHHhh$ of F_2 genotypes.

Collecting all terms, one obtains the 6-locus SD equation:

$$E[S_1S_2S_3S_4S_5S_6] = A_{1,1,1,1,1,1}(E[S_1S_2S_3S_4S_5S_6] + 1) + A_{1,1,1,1,0,0}(E[S_1S_2S_3S_4] + E[S_5S_6]) \\ + A_{0,0,1,1,1,1}(E[S_3S_4S_5S_6] + E[S_1S_2]) + A_{1,1,0,0,1,1}(E[S_1S_2S_5S_6] + E[S_3S_4]) \quad (19)$$

Note that these equations can be generalized to the case where recombination rates differ between male and female meiosis (see Section VII and Fig. S8 in Supplemental Material).

The patterns found in all these equations are easily extended to any number of loci. The SD equation for $E[S_1S_2 \dots S_L]$ can be written in terms of factors A_{n_1, n_2, \dots, n_L} and associated expectations of multi-spin products where the indices of these factors must satisfy the constraint of Rule 4 in Section V in Supplemental Material: 0s and 1s must come in adjacent pairs. For each such A_{n_1, n_2, \dots, n_L} , there is a global factor of $1/2$, a factor for each block of a given type (block of 0s or block of 1s), and one factor for each interval *connecting* blocks. The factor for connecting two blocks is $[2(1-r)r]$, r being the recombination rate in that connecting interval. The factor *within* a block is a product over all of its intervals, alternating between $(1-2r)$ terms and $[(1-r)^2 + r^2]$ terms and ending with a $(1-2r)$ term because the number of intervals is odd. These results show that the SD equations can be written down *automatically* for any number of loci.

VII. GENERALIZING THE FORMULAS TO SEX-SPECIFIC RECOMBINATION RATES

We saw how to generalize the standard Haldane-Waddington 2-locus formula for R to situations where the female and male meiotic recombination rates r^f and r^m differ (see Section II in Supplemental Material). Interestingly, it is also possible to generalize *all* our L -locus formulas to such a situation as follows.

First, the Glauber formula (Eq. 1 in Main Text) that gives the probabilities of the RIL genotypes in terms of expectation values of spin products is unchanged because it does not involve recombination rates and even less sex-specificity. Second, moving on to the SD equations, sex-dependence arises only at the level of the probabilities of gametes, *i.e.*, through the probabilities $P(g)$ and $P(g')$ (Eq. 3 in Main Text). The probabilities $P(g)$ and $P(g')$ must be modified but otherwise the logic is the same as in the sex-independent case. Specifically, one considers classes of F2 genotypes according to whether the successive loci are homozygous (H) or heterozygous (h). One maps these genotypes to binary trees as before to obtain a factor that multiplies an expectation value in the SD equation. That factor is a product of terms, one for each interval between adjacent loci. If, in the sex-independent case, an interval contributed the factor $(1-r)^2 - r^2$ (which simplifies to $1-2r$), it will now contribute $(1-r^f)(1-r^m) - r^f r^m$ (which simplifies to $(1-r^f - r^m)$). If an interval contributed $(1-r)^2 + r^2$ in the sex-independent case, it will now contribute $(1-r^f)(1-r^m) + r^f r^m$. If an interval contributed $2(1-r)r$ in the sex-independent case, it will now contribute $(1-r^f)r^m + r^f(1-r^m)$. However, this is not the end of the story: in the sex-independent case, a large number of trees were discarded because one of the intervals led to the factor 0 (Fig. S3 in Supplemental Material). For instance, for the class $hHhH$ in the sex-independent case, when one does the pooling of pairs in the right-most interval, one is led to $(1-r)r - r(1-r)$ which shows that the tree can be ignored (Fig. S3a in Supplemental Material). However, in the sex-specific case, that factor becomes $(1-r^f)r^m - r^f(1-r^m) = r^m - r^f$ which has no reason to vanish. Going back to the rules listed in Section V in Supplemental Material, it transpires that Rule 4 requires exchanging female and male meiosis. Thus, for sex-specific rates, this last rule and its associated simplifications have to be abandoned.

To illustrate the changes required for sex-specific recombination rates, consider the SD equation for $E[S_1S_2S_3S_4]$. The right-hand side of that SD equation contains one factor multiplying $E[S_1S_2S_3S_4]$ (the self term), factors for all of the $E[S_iS_j]$ terms (for *any* pair (i, j) of distinct loci, not just the (1,2) and (3,4) pairs found in the sex-independent case), and finally a factor multiplying 1 (no associated expectation). The fact that no other expectation values contribute is due to Rule 2 in Section V in Supplemental Material. Rule 3 in Section V in Supplemental Material implies that the factor in front of $E[S_1S_2S_3S_4]$ is the same as in front of 1, and also that the factor in front of $E[S_iS_j]$ is the same as in front of $E[S_kS_l]$ where $i, j, k,$ and l are all distinct. As an example, to obtain the factor multiplying $E[S_2S_4]$, we need to consider the F2 genotypes that are heterozygous at loci 2 and 4 and homozygous at loci 1 and 3. This class of genotypes gives two trees, one rooted at a_1/a_1 and other at A_1/A_1 . By Rule 1 in Section V in Supplemental Material, these contribute equally to the SD equation so it is sufficient to consider the first tree. For the sex-independent case, this tree (associated with class $HhHh$) satisfies all the rules listed in Section V in Supplemental Material except the last rule, so it vanishes when $r^f = r^m$ (Fig. S3b in Supplemental Material). The calculation of factors of this tree for the *sex-specific* case ($r^f \neq r^m$) are given in Fig. S8 in Supplemental Material to which must be taken into account the factor $1/4$ for the root of this tree. Putting these factors together, we conclude that the term multiplying $E[S_2S_4]$ in the SD equation is:

$$A_{0,1,0,1} = \frac{(r_1^m - r_1^f)[r_2^m + r_2^f - 2r_2^m r_2^f](r_3^m - r_3^f)}{2} \quad (20)$$

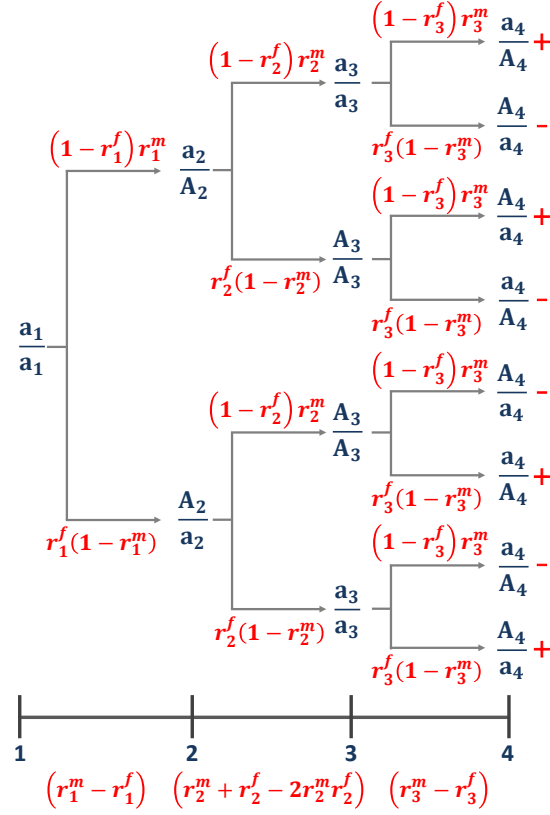


FIG. S8. Four-locus tree for the class $HhHh$ of F2 genotypes in the case of sex-specific recombination rates.

in the sex-specific case.

VIII. COMPUTER PROGRAMS FOR COMPUTING PROBABILITIES OF RIL GENOTYPES

Because the A_{n_1, n_2, \dots, n_L} coefficients of the SD equations follow such stereotyped patterns, it is possible to produce a computer program which determines them automatically. We have done so numerically within a C-language code that furthermore uses them to calculate all averages of k -spin products recursively for increasing k . Once all these averages have been tabulated, the program uses Glauber's equation to compute the probabilities of all 2^L RIL genotypes. Note that this last step naively takes on the order of 4^L operations, but in fact it is possible to use a multi-dimensional transform that requires only on the order of $L2^L$ operations [4]. The resulting computer program is available online as a Supplementary file. Its computation time grows by about a factor 10 when $L \rightarrow L + 2$. For illustration, the treatment of the case with $L = 14$ loci can be done in less than a second using a standard desktop computer.

To the extent that a purely numerical estimate of RIL probabilities is sufficient, other approaches are also possible. The most straightforward one consists in simulating the steps of production of a RIL, implementing the successive generations until the genotype produced is homozygous. If one repeats this process many times, one can get a large sample of RIL genotypes from which genotype probabilities can be estimated. However the number of different genotypes grows as 2^L ; if one wants to have reliable estimates of all genotypes, the sample must have several hundred realizations of each genotype, no matter how rare each genotype might be. Consequently, this approach is not very useful when the number of loci is greater than 10 and furthermore it is extremely inefficient if one needs precise estimates of the RIL probabilities. To overcome the statistical limitations of stochastic simulation, one may use instead the master equation. The procedure consists in following recursively (from one generation to the next) the probability of all 4^L possible genotypes. This recursion can be written as a $4^L \times 4^L$ matrix operating on a vector, and is in fact the approach provided by Haldane and Waddington. The limit of a large number of generations is associated with one of the leading eigenvectors. For L not too large, the appropriate eigenvector can be computed by

diagonalizing the matrix, but for large matrices (already at $L = 10$ the matrix has more than one million rows), one is forced to rely on the power method. Implementing this method requires one to apply the matrix to the vector a large number of times, large enough to see convergence of the recursion to a fixed point. We have coded this procedure in a C-language program that is also provided online as a Supplementary file. Its computation time grows by about a factor 100 when $L \rightarrow L + 2$. The treatment of $L = 10$ can be done in a few hours on a desktop computer depending on the number of iterations used to get close to the fixed point. It is thus about a million times slower at $L = 14$ than the program exploiting the Schwinger-Dyson equations and Glauber's formula.

-
- [1] J.B.S. Haldane, *Journal of Genetics* **8**, 299-309 (1919).
 - [2] R.B. Robbins, *Genetics*, **3**, 375-389 (1918).
 - [3] J.B.S. Haldane and C.H. Waddington, *Genetics* **16**, 357-374 (1931).
 - [4] F. Zanini and R.A. Neher, *Bioinformatics* **28**, 3332-3333 (2012).