1 MATHEMATICAL APPENDIX

1.1 Distribution of distinct haplotypes within infections

Denote the number of distinct haplotype in an infection by the random variable C. Assuming the statistical model outlined in the main text, the probability distribution of C can be derived.

Remember *H* is the number of possible haplotypes. First, for subset of possible haplotypes $A \subseteq \{1, \ldots, H\}$ we define the set of all MOI vectors *m* with exactly the haplotypes in *A* infecting by

$$M_A := \{ \boldsymbol{m} \in \mathbb{N}^H \, | \, m_h > 0 \Leftrightarrow h \in A \}$$
(A.1)

and the set of all MOI vectors m with not necessarily all haplotypes in A infecting by

$$\tilde{M}_A := \{ \boldsymbol{m} \in \mathbb{N}^H \mid m_h = 0 \text{ for } h \notin A \}.$$
(A.2)

Assuming c > 0, the probability of an infection with exactly c distinct haplotypes is calculated to be

$$P[C=c] = \sum_{m=c}^{\infty} \kappa_m \sum_{\substack{A \subseteq \{1,\dots,H\}: \ \boldsymbol{m} \in M_A \\ |A|=c}} \sum_{\boldsymbol{m} \in M_A} \binom{m}{\boldsymbol{m}} p^{\boldsymbol{m}}.$$
(A.3)

Using the inclusion-exclusion principle, followed by the binomial theorem, an interchange of the summation, and the definition of the modification of the probability generating function (10) (see also Table 1) this becomes

$$P[C = c] = \sum_{m=1}^{\infty} \kappa_m \sum_{\substack{A \subseteq \{1, \dots, H\}: B \subseteq A}} \sum_{B \subseteq A} (-1)^{|A| - |B|} \sum_{\boldsymbol{m} \in \tilde{M}_B} \binom{m}{\boldsymbol{m}} \boldsymbol{p}^{\boldsymbol{m}}$$

$$= \sum_{m=1}^{\infty} \kappa_m \sum_{\substack{A \subseteq \{1, \dots, H\}: B \subseteq A}} \sum_{B \subseteq A} (-1)^{|A| - |B|} \left(\sum_{h \in B} p_h\right)^m$$

$$= \sum_{\substack{A \subseteq \{1, \dots, H\}: B \subseteq A}} \sum_{B \subseteq A} (-1)^{|A| - |B|} G\left(\sum_{h \in B} p_h\right).$$
(A.4)

Note that the last step holds also in the case in which disease-free samples are considered ($\kappa_0 > 0$). Moreover, the disease free case, i.e., c = 0, occurs if and only if MOI = 0, which occurs with probability κ_0 , i.e., we obtain

$$\mathbf{P}[C=0] = \kappa_0. \tag{A.5}$$

1.2 Distribution of the maximum number of alleles across markers

In the absence of phased haplotype information, an observation x provides in general only ambiguous information about the haplotypes being actually present in the corresponding infection (compare 'absence/presence' with 'infecting haplotypes' in Figure 4). On the one hand, if a haplotype h carries an allele which is not observed in x, it cannot be present in the infection, i.e. it is incompatible with the observation. On the other hand, we say that a haplotype h is compatible if all the alleles of which it is comprised are observed, i.e., it cannot be ruled out that h is actually present in the infection.

To calculate the distribution of the maximum number of alleles observed across the L loci (K), let us define the set of all observations \boldsymbol{x} for which this number is exactly k. Namely,

$$U_k := \{ \boldsymbol{x} = (\boldsymbol{x}_1, \dots, \boldsymbol{x}_L) \mid \forall l : |\boldsymbol{x}_l| \le k \land \exists l : |\boldsymbol{x}_l| = k \}.$$
(A.6)

Then, the probability of such an observation is

$$P[K = k] = \sum_{\boldsymbol{x} \in U_k} P[\boldsymbol{x}].$$
(A.7)

This equation involves the probabilities to observe x, i.e., P[x], which we need to derive in a more explicit form.

First, let the set of all possible disease-positive observations be denoted

$$\mathscr{O} := \prod_{l=1}^{L} \left(\mathscr{P}(\{1, \dots, n_l\}) \setminus \emptyset \right).$$
(A.8)

In case also disease-negative samples are included, the observation $\mathbf{0} = (\emptyset, \dots, \emptyset)$ is possible, hence \mathscr{O} has to be replaced by $\mathscr{O} \cup \{\mathbf{0}\}$. Clearly,

$$\mathbf{P}[\mathbf{0}] = \kappa_0. \tag{A.9}$$

We can say an observation $\boldsymbol{y} = (\boldsymbol{y}_1, \dots, \boldsymbol{y}_L)$ is subsumed by an observation $\boldsymbol{x} = (\boldsymbol{x}_1, \dots, \boldsymbol{x}_L)$ if at each locus all alleles which occur in observation \boldsymbol{y} also occur in observation \boldsymbol{x} , we write

$$\boldsymbol{y} \preceq \boldsymbol{x} \quad :\Leftrightarrow \quad \boldsymbol{y}_l \subseteq \boldsymbol{x}_l \quad \text{for all } l.$$
 (A.10a)

Further define the set of all observations subsumed by \boldsymbol{x} as

$$\mathscr{A}_{\boldsymbol{x}} := \{ \boldsymbol{y} \in \mathscr{O} \mid \boldsymbol{y} \leq \boldsymbol{x} \}.$$
(A.10b)

Let the set of all haplotypes h, which are compatible with observation x, be denoted by

$$A_{\boldsymbol{x}} := \{ \boldsymbol{h} = (h_1, \dots, h_L) \in \mathscr{H} \mid h_l \in \boldsymbol{x}_l \text{ for all } l \}.$$
(A.10c)

The set of all MOI vectors m which are compatible with observation x are denoted by

$$M_{\boldsymbol{x}} = \{ \boldsymbol{m} \in \mathbb{N}^{H} \mid \boldsymbol{m} \to \boldsymbol{x} \}$$

= $\{ \boldsymbol{m} \in \mathbb{N}^{H} \mid \forall l : \boldsymbol{x}_{l} = \{ h_{l} \mid \forall \boldsymbol{h} = (h_{1}, \dots, h_{L}) \in \mathscr{H} \text{ with } m_{\boldsymbol{h}} > 0 \} \}.$ (A.10d)

Furthermore, the set of all MOI vectors m which are compatible with an observation subsumed by x, i.e., with an observation in \mathcal{A}_x , is denoted by

$$\tilde{M}_{\boldsymbol{x}} = \left\{ \boldsymbol{m} \in \mathbb{N}^H \, \middle| \, \forall \boldsymbol{h} \notin A_{\boldsymbol{x}} : \, m_{\boldsymbol{h}} = 0 \right\}.$$
(A.10e)

Using the definition of the set M_x , the probability of observing x can be rewritten as

$$P[\boldsymbol{x}] = \sum_{m=0}^{\infty} \kappa_m \sum_{\substack{|\boldsymbol{m}|=m\\ \boldsymbol{m} \to \boldsymbol{x}}} \binom{m}{\boldsymbol{m}} \boldsymbol{p}^{\boldsymbol{m}} = \sum_{\boldsymbol{m} \in M_{\boldsymbol{x}}} \kappa_{|\boldsymbol{m}|} \binom{|\boldsymbol{m}|}{\boldsymbol{m}} \boldsymbol{p}^{\boldsymbol{m}}.$$
(A.11)

By an inclusion-exclusion argument the above becomes

$$P[\boldsymbol{x}] = \sum_{\boldsymbol{y} \in \mathscr{A}_{\boldsymbol{x}}} (-1)^{\sum_{l=1}^{L} |\boldsymbol{x}_l| - |\boldsymbol{y}_l|} \sum_{\boldsymbol{m} \in \tilde{M}_{\boldsymbol{y}}} \kappa_{|\boldsymbol{m}|} \binom{|\boldsymbol{m}|}{\boldsymbol{m}} \boldsymbol{p}^{\boldsymbol{m}}$$

In the case that $\boldsymbol{y} = \boldsymbol{0}$, the above equation needs some adjustment, namely the power $\sum_{l=1}^{L} |\boldsymbol{x}_l| - |\boldsymbol{y}_l|$ has to be replaced by $L - 1 + \sum_{l=1}^{L} |\boldsymbol{x}_l| - |\boldsymbol{y}_l|$. Therefore, we define

$$d(\boldsymbol{x}, \boldsymbol{y}) := \begin{cases} L - 1 + \sum_{l=1}^{L} |\boldsymbol{x}_l| - |\boldsymbol{y}_l| & \text{for } \boldsymbol{x} \neq \boldsymbol{0} \text{ and } \boldsymbol{y} = \boldsymbol{0}, \\ \sum_{l=1}^{L} |\boldsymbol{x}_l| - |\boldsymbol{y}_l| & \text{else.} \end{cases}$$
(A.12)

Using the theorem of total probability, then the binomial theorem and finally the definition of the generating functions (10) leads to

$$P[\boldsymbol{x}] = \sum_{\boldsymbol{y} \in \mathscr{A}_{\boldsymbol{x}}} (-1)^{d(\boldsymbol{x},\boldsymbol{y})} \sum_{m=1}^{\infty} \kappa_m \sum_{\substack{\boldsymbol{m} \in \tilde{M}_{\boldsymbol{y}} \\ |\boldsymbol{m}| = m}} {\binom{m}{\boldsymbol{m}} \boldsymbol{p}^{\boldsymbol{m}}}$$
$$= \sum_{\boldsymbol{y} \in \mathscr{A}_{\boldsymbol{x}}} (-1)^{d(\boldsymbol{x},\boldsymbol{y})} \sum_{m=1}^{\infty} \kappa_m \left(\sum_{\boldsymbol{h} \in A_{\boldsymbol{y}}} p_{\boldsymbol{h}}\right)^m$$
$$= \sum_{\boldsymbol{y} \in \mathscr{A}_{\boldsymbol{x}}} (-1)^{d(\boldsymbol{x},\boldsymbol{y})} G\left(\sum_{\boldsymbol{h} \in A_{\boldsymbol{y}}} p_{\boldsymbol{h}}\right).$$
(A.13)

Hence, the distribution of the maximum number of alleles across loci is

$$P[K = k] = \sum_{\boldsymbol{x} \in U_k} \sum_{\boldsymbol{y} \in \mathscr{A}_{\boldsymbol{x}}} (-1)^{d(\boldsymbol{x}, \boldsymbol{y})} G\bigg(\sum_{\boldsymbol{h} \in A_{\boldsymbol{y}}} p_{\boldsymbol{h}}\bigg).$$
(A.14)

1.3 Distribution of the average number of alleles across markers

Here, the distribution of the average number of alleles across markers is derived. For this purpose define the set of all observations in which the number of observed alleles across the loci sum up to k by

$$V_k := \left\{ \boldsymbol{x} = (\boldsymbol{x}_1, \dots, \boldsymbol{x}_L) \, \middle| \, \sum_{l=1}^L |\boldsymbol{x}_l| = k \right\}.$$
(A.15)

The distribution of the average number of alleles, \overline{K} , hence becomes

$$P\left[\overline{K} = \frac{k}{L}\right] = \sum_{\boldsymbol{x} \in V_k} P[\boldsymbol{x}].$$
(A.16)

Using the expression (A.13) for $P[\mathbf{x}]$, one obtains

$$P\left[\overline{K} = \frac{k}{L}\right] = \sum_{\boldsymbol{x} \in V_k} \sum_{\boldsymbol{y} \in \mathscr{A}_{\boldsymbol{x}}} (-1)^{d(\boldsymbol{x}, \boldsymbol{y})} G\left(\sum_{\boldsymbol{h} \in A_{\boldsymbol{y}}} p_{\boldsymbol{h}}\right).$$
(A.17)

1.4 Prevalence

The prevalence of haplotype h is the probability that it occurs in an infection. To derive this we use the index notation h rather than the vector notation h, which yields

$$P[M_h > 0] = \sum_{m=0}^{\infty} P[M = m] P[M_h > 0 | M = m] = \sum_{m=0}^{\infty} \kappa_m \sum_{\substack{\boldsymbol{m}:\\m_h > 0}} \binom{m}{\boldsymbol{m}} \boldsymbol{p}^{\boldsymbol{m}}$$
$$= \sum_{m=0}^{\infty} \kappa_m \left(\sum_{\boldsymbol{m}=1}^{H} p_g \right)^{\boldsymbol{m}} - \sum_{\substack{\boldsymbol{m}:\\m_h = 0}} \binom{m}{\boldsymbol{m}} \boldsymbol{p}^{\boldsymbol{m}} \right)$$
$$= \sum_{m=0}^{\infty} \kappa_m \left(\left(\sum_{g=1}^{H} p_g \right)^{\boldsymbol{m}} - \left(\sum_{\substack{g=1\\g \neq h}}^{H} p_g \right)^{\boldsymbol{m}} \right)$$

Using the definition of the generation function (10) this becomes

$$P[M_h > 0] = 1 - G(1 - p_h).$$

This formula, translated into the vector notation of haplotypes becomes

$$P[M_{h} > 0] = 1 - G(1 - p_{h}).$$

1.5 Distribution of MOI conditional on a particular observation

Using the parameter estimates $\hat{\theta}$ as a plug-in for the true parameters, and denoting the generating function based on the estimates by \hat{G} , the estimated probability of MOI = m given observation \boldsymbol{x} is calculated to be

$$P[MOI = m | \boldsymbol{x}] = \frac{P[\boldsymbol{x}, m]}{P[\boldsymbol{x}]} = \frac{\frac{\hat{\kappa}_m \sum_{\boldsymbol{m}: \ |\boldsymbol{m}| = m}}{\sum_{\boldsymbol{y} \in \mathscr{A}_{\boldsymbol{x}}} (-1)^{d(\boldsymbol{x}, \boldsymbol{y})} \hat{G}\left(\sum_{\boldsymbol{h} \in A_{\boldsymbol{y}}} \hat{p}_{\boldsymbol{h}}\right)}{\sum_{\boldsymbol{y} \in \mathscr{A}_{\boldsymbol{x}}} (-1)^{d(\boldsymbol{x}, \boldsymbol{y})} \hat{G}\left(\sum_{\boldsymbol{h} \in A_{\boldsymbol{y}}} \hat{p}_{\boldsymbol{h}}\right)}.$$
(A.18)

By using an inclusion-exclusion argument in the numerator, the above becomes

$$P[MOI = m | \boldsymbol{x}] = \frac{P[\boldsymbol{x}, m]}{P[\boldsymbol{x}]} = \frac{\hat{\kappa}_m \sum_{\boldsymbol{y} \in \mathscr{A}_{\boldsymbol{x}}} (-1)^{d(\boldsymbol{x}, \boldsymbol{y})} \left(\sum_{\boldsymbol{h} \in A_{\boldsymbol{y}}} \hat{p}_{\boldsymbol{h}}\right)^m}{\sum_{\boldsymbol{y} \in \mathscr{A}_{\boldsymbol{x}}} (-1)^{d(\boldsymbol{x}, \boldsymbol{y})} \hat{G}\left(\sum_{\boldsymbol{h} \in A_{\boldsymbol{y}}} \hat{p}_{\boldsymbol{h}}\right)}.$$
(A.19)

The true probability of MOI = m given observation \boldsymbol{x} is obtained from (A.19) by substituting the true parameters $\boldsymbol{\theta}$ (and the generating function G) for the plug-in parameters $\hat{\boldsymbol{\theta}}$.

1.6 Distribution of the number of haplotypes given a particular observation

Let the set of all possible ensembles (sets) of exactly c haplotypes which together can create observation \boldsymbol{x} by

$$B_{\boldsymbol{x}}^{(c)} := \{\{\boldsymbol{h}_1, \dots, \boldsymbol{h}_c\} \subseteq \mathscr{H} \mid |\{\boldsymbol{h}_1, \dots, \boldsymbol{h}_c\}| = c, \ x_l = \{(\boldsymbol{h}_k)_l \mid k = 1, \dots, c\} \text{ for } l = 1, \dots, L\}.$$
(A.20)

For simplicity of notation we denote elements of the set $B_{\boldsymbol{x}}^{(c)}$, i.e., ensembles of c haplotypes leading to observation \boldsymbol{x} by A, without reference to c and \boldsymbol{x} .

For c haplotypes $A = {\mathbf{h}_1, \dots, \mathbf{h}_c}$ leading to observation \mathbf{x} , we define the set of all possible superinfections with exactly the haplotypes in A as

$$M_A = \left\{ \boldsymbol{m} \in \mathbb{N}^H \, \middle| \, \boldsymbol{m}_{\boldsymbol{h}} > 0 \, \Leftrightarrow \, \boldsymbol{h} \in A \right\}.$$
(A.21)

The probability of an infection with an ensemble A of c haplotypes leading to observation \boldsymbol{x} , i.e., for $A \in B_{\boldsymbol{x}}^{(c)}$ we have

$$P[\boldsymbol{x}, A] = \sum_{m=0}^{\infty} \hat{\kappa}_m \sum_{\substack{\boldsymbol{m} \in M_A:\\ |\boldsymbol{m}|=m}} {\binom{m}{\boldsymbol{p}}} \hat{\boldsymbol{p}}^{\boldsymbol{m}}.$$
(A.22)

Assume $A \in B_{\boldsymbol{x}}^{(c)}$, i.e., A contains exactly c haplotypes and leads to observation \boldsymbol{x} . Then a proper subset $B \subsetneq A$ contains less than c haplotypes and does not necessarily lead to observation \boldsymbol{x} . For a set of arbitrary many haplotypes B, we denote the set of all possible super-infections with (not necessarily all) haplotypes in B by

$$\tilde{M}_B = \left\{ \boldsymbol{m} \in \mathbb{N}^H \, \middle| \, \boldsymbol{m}_{\boldsymbol{h}} = 0 \text{ if } \boldsymbol{h} \notin B \right\}.$$
(A.23)

By an inclusion-exclusion argument, (A.20) can be rewritten as

$$P[\boldsymbol{x}, A] = \sum_{B \subseteq A} (-1)^{|A| - |B|} \sum_{m=0}^{\infty} \hat{\kappa}_m \sum_{\substack{\boldsymbol{m} \in \tilde{M}_B:\\|\boldsymbol{m}| = m}} {\binom{m}{\boldsymbol{p}}} \hat{\boldsymbol{p}}^{\boldsymbol{m}}.$$
(A.24)

By the binomial theorem, the fact that |A| = c, and the definition of the generating function we obtain

$$P[\boldsymbol{x}, A] = \sum_{m=0}^{\infty} \hat{\kappa}_m \sum_{B \subseteq A} (-1)^{c-|B|} \left(\sum_{\boldsymbol{h} \in B} \hat{p}_{\boldsymbol{h}}\right)^m$$

$$= \sum_{B \subseteq A} (-1)^{c-|B|} \hat{G}\left(\sum_{\boldsymbol{h} \in B} \hat{p}_{\boldsymbol{h}}\right).$$
(A.25)

Clearly, a super-infection with exactly the haplotypes in a set $A \in B_{\boldsymbol{x}}^{(c)}$ cannot be a super-infection with exactly the haplotypes in a different set $\tilde{A} \in B_{\boldsymbol{x}}^{(c)}$. In other words, super-infections with exactly the haplotypes in a set $A \in B_{\boldsymbol{x}}^{(c)}$ and in a set $A' \in B_{\boldsymbol{x}}^{(c)}$ are disjoint events, or if $A \neq A'$, then $M_A \cap M_{A'} = \emptyset$. Consequently, the probability of an infection with exactly c haplotypes leading to observation \boldsymbol{x} is

$$P[\boldsymbol{x}, C = c] = \sum_{A \in B_{\boldsymbol{x}}^{(c)}} P[\boldsymbol{x}, A]$$

$$= \sum_{A \in B_{\boldsymbol{x}}^{(c)}} \sum_{B \subseteq A} (-1)^{c-|B|} \hat{G}\left(\sum_{\boldsymbol{h} \in B} \hat{p}_{\boldsymbol{h}}\right).$$
(A.26)

Therefore, given the observation \boldsymbol{x} , the probability that exactly C = c different haplotypes were infecting becomes

$$P[C = c | \boldsymbol{x}] = \frac{P[\boldsymbol{x}, C = c]}{P[\boldsymbol{x}]} = \frac{\sum_{A \in B_{\boldsymbol{x}}^{(c)}} \sum_{B \subseteq A} (-1)^{c - |B|} \hat{G}\left(\sum_{\boldsymbol{h} \in B} \hat{p}_{\boldsymbol{h}}\right)}{\sum_{\boldsymbol{y} \in \mathscr{A}_{\boldsymbol{x}}} (-1)^{d(\boldsymbol{x}, \boldsymbol{y})} \hat{G}\left(\sum_{\boldsymbol{h} \in A_{\boldsymbol{y}}} \hat{p}_{\boldsymbol{h}}\right)}.$$
(A.27)