

Appendix

We note that all notation for this appendix is provided in Table 1 in the main manuscript.

Component probability mass function for Causal Variant (CV) count data for single locus

Before proceeding with the log-likelihoods and the test statistics, we compute the binomial probability mass function $Bin(x_{1,k} | v_{1,k}^t, i_{1,k}^t, j_1^t)$ for a single locus. Here, the subscript "1" indicates that it is the "first" locus. The notation becomes important later on, when we consider multiple loci. We comment that determination of the probability mass function is necessary for computation of the log-likelihoods.

Given CV counts, coverage, and phenotype data $(x_{1,k}, v_{1,k}^t, i_{1,k}^t)$, we prove that the conditional probability of the CV counts $x_{1,k}$ for individual k , given coverage $v_{1,k}^t$, phenotype $i_{1,k}^t$, and genotype j_1^t follows a binomial distribution, $Bin(x_{1,k}; v_{1,k}^t; p(i_{1,k}^t, j_1^t))$, where

$p(i_{1,k}^t, j_1^t) = \left(\frac{2-j_1^t}{2} \varepsilon_{i_{1,k}^t} + \frac{j_1^t}{2} (1 - \varepsilon_{i_{1,k}^t}) \right)$ is the probability of a "success", namely the probability

of observing the causal variant A .

To see this, consider the following: for genotype $j_1^t = 0$ and $j_1^t = 2$, it is straightforward. For instance, when $j_1^t = 0$, the individual's genotype is a/a . Thus every observed CV A is really an a that has been misread. The number of errors is $x_{1,k}$, so that $x_{1,k}$ follows a Binomial distribution with probability $\varepsilon_{i_{1,k}^t}$ given $v_{1,k}^t$ trials. Similarly, when $j_1^t = 2$, every read of CV A is now correct read from the genotype A/A . Thus, $x_{1,k}$ follows a Binomial distribution with probability $(1 - \varepsilon_{i_{1,k}^t})$ given $v_{1,k}^t$ trials. For the heterozygote genotype $j_1^t = 1$, the probability of CV A being read on a single trial is the sum of two probabilities:

$$\begin{aligned}
& \Pr(a)\Pr(\text{read} = A | \text{correct variant} = a) + \Pr(A)\Pr(\text{read} = A | \text{correct variant} = A), \\
&= \frac{1}{2} \varepsilon_{i_{1,k}}^t + \frac{1}{2} (1 - \varepsilon_{i_{1,k}}^t), \\
&= \frac{1}{2}.
\end{aligned}$$

Here, we specify a symmetric error model in which:

$$\Pr(\text{read} = A | \text{correct variant} = a) = \Pr(\text{read} = a | \text{correct variant} = A).$$

An asymmetric model, where:

$$\Pr(\text{read} = A | \text{correct variant} = a) \neq \Pr(\text{read} = a | \text{correct variant} = A),$$

may be considered. However, the estimation of $\Pr(\text{read} = a | \text{correct variant} = A)$ is quite unstable due to low frequency of causal variant (even for 1000 samples).

Log-likelihood of the observed data

The log-likelihood function of the multiple-locus CV count data involves 3^M genotype configurations. If we set $\mathbf{x}_k = (x_{1,k}, \dots, x_{M,k})$, $\mathbf{v}_k^t = (v_{1,k}^t, \dots, v_{M,k}^t)$, $\mathbf{G}^t = (j_1^t, \dots, j_M^t)$, we can write:

$$\ln(L_{H_d}) = \sum_{k=1}^N \ln[\Pr(\mathbf{x}_k, \mathbf{v}_k^t, i_k^t)],$$

$$\begin{aligned}
&= \sum_{k=1}^N \ln \left[\sum_{\mathbf{G}^t} (\Pr(\mathbf{x}_k, \mathbf{v}_k^t, i_k^t, \mathbf{G}^t)) \right], \\
&= \sum_{k=1}^N \ln \left[\sum_{\mathbf{G}^t} (\Pr(\mathbf{x}_k | \mathbf{v}_k^t, i_k^t, \mathbf{G}^t) \times \Pr(i_k^t | \mathbf{v}_k^t, \mathbf{G}^t) \times \Pr(\mathbf{v}_k^t | \mathbf{G}^t) \times \Pr(\mathbf{G}^t)) \right], \\
&= \sum_{k=1}^N \ln \left[\sum_{\mathbf{G}^t} (\Pr(\mathbf{x}_k | \mathbf{v}_k^t, i_k^t, \mathbf{G}^t) \times \Pr(i_k^t | \mathbf{G}^t) \times \Pr(\mathbf{v}_k^t | \mathbf{G}^t) \times \mathbf{g}_{(j_1^t, \dots, j_M^t)}) \right], \quad (A1) \\
&= \sum_{k=1}^N (1 - i_k^t) \times \ln \left[\sum_{\mathbf{G}^t} (\Pr(\mathbf{x}_k | \mathbf{v}_k^t, i_k^t = 0, \mathbf{G}^t) \times \Pr(i_k^t = 0 | \mathbf{G}^t) \times \mathbf{g}_{(j_1^t, \dots, j_M^t)}) \right] \\
&+ \sum_{k=1}^N (i_k^t) \times \ln \left[\sum_{\mathbf{G}^t} (\Pr(\mathbf{x}_k | \mathbf{v}_k^t, i_k^t = 1, \mathbf{G}^t) \times \Pr(i_k^t = 1 | \mathbf{G}^t) \times \mathbf{g}_{(j_1^t, \dots, j_M^t)}) \right] \\
&+\gamma,
\end{aligned}$$

where $\gamma = \sum_{k=1}^N \ln(\Pr(\mathbf{v}_k^t))$ and $\mathbf{g}_{(j_1^t, \dots, j_M^t)} = \Pr(\mathbf{G}^t) = \Pr((j_1^t, \dots, j_M^t))$. The last equality follows from the fact that the phenotypes and genotypes are independent of the coverage.

Summing over all true genotype vectors \mathbf{G}^t is equivalent to summing each locus m 's genotype value from 0 to 2. In this work, we specify that, conditional on the underlying data (including the genotype vectors \mathbf{G}^t), the observed CV counts are independent. Written another way, we have:

$$\begin{aligned}
&\Pr(\mathbf{x}_k | i_k^t, \mathbf{v}_k^t = (v_{1,k}^t, \dots, v_{M,k}^t), \mathbf{G}^t = (j_1^t, \dots, j_M^t)) \\
&= \prod_{m=1}^M \text{Bin}(x_{m,k}; v_{m,k}^t; p(i_{1,k}^t, j_m^t)). \quad (A2)
\end{aligned}$$

It follows that we can rewrite the Equation (A1) under the null as:

$$\begin{aligned} \ln(L_{H_0}) &= \sum_{k=1}^N (1 - i_k) \ln \left[\sum_{j_1=0}^2 \cdots \sum_{j_M=0}^2 \left(\prod_{m=1}^M \text{Bin}(x_{m,k}; v_{m,k}^t; p(0, j_m^t)) \right) (f_{0,*} \times \mathbf{g}_{(j_1^t, \dots, j_M^t)}) \right] \\ &+ \sum_{k=1}^N (i_k) \ln \left[\sum_{j_1=0}^2 \cdots \sum_{j_M=0}^2 \left(\prod_{m=1}^M \text{Bin}(x_{m,k}; v_{m,k}^t; p(1, j_m^t)) \right) (f_{1,*} \times \mathbf{g}_{(j_1^t, \dots, j_M^t)}) \right] + \gamma. \end{aligned} \quad (A3)$$

Similarly, for the alternative hypothesis, we have:

$$\begin{aligned} \ln(L_{H_1}) &= \sum_{k=1}^N (1 - i_k) \ln \left[\sum_{j_1=0}^2 \cdots \sum_{j_M=0}^2 \left(\prod_{m=1}^M \text{Bin}(x_{m,k}; v_{m,k}^t; p(0, j_m^t)) \right) (f_{0,(j_1^t, \dots, j_M^t)} \times \mathbf{g}_{(j_1^t, \dots, j_M^t)}) \right] \\ &+ \sum_{k=1}^N (i_k) \ln \left[\sum_{j_1=0}^2 \cdots \sum_{j_M=0}^2 \left(\prod_{m=1}^M \text{Bin}(x_{m,k}; v_{m,k}^t; p(1, j_m^t)) \right) (f_{1,(j_1^t, \dots, j_M^t)} \times \mathbf{g}_{(j_1^t, \dots, j_M^t)}) \right] + \gamma. \end{aligned} \quad (A4)$$

EM Algorithm Estimates

We provide closed formula solutions of the $(r + 1)^{\text{st}}$ -step estimates of the parameters necessary to compute the log-likelihoods. The main difference is that estimates are determined as a function of a vector of genotypes rather than a single genotype. We determine that the posterior probability that individual k has genotype vector $\mathbf{G}^t = (j_1^t, \dots, j_M^t)$ is given by:

$$\begin{aligned} \tau_{\mathbf{G}^t, k}^{(r)} &= \tau_{(j_1^t, \dots, j_M^t), k}^{(r)} \\ &= \frac{f_{i_k^t, (j_1^t, \dots, j_M^t)} \Pr((j_1^t, \dots, j_M^t)) \left[\prod_{m=1}^M \text{Bin}(x_{m,k}; v_{m,k}^t; p(i_{1,k}^t, j_m^t)) \right]}{\sum_{s_1^t=0}^2 \cdots \sum_{s_M^t=0}^2 \left(f_{i_k^t, (s_1^t, \dots, s_M^t)} \Pr((s_1^t, \dots, s_M^t)) \left[\prod_{m=1}^M \text{Bin}(x_{m,k}; v_{m,k}^t; p(i_{1,k}^t, s_m^t)) \right] \right)}. \end{aligned} \quad (A5)$$

For the parameters α and β , we make use of the following notation:

$$\tau_{(j_1^t, \dots, j_M^t)_{+(1)}}^{(r)} = \sum_{k=1}^N (i_k) \tau_{(j_1^t, \dots, j_M^t), k}^{(r)}$$

$$\tau_{++(1)}^{(r)} = \sum_{k=1}^N \sum_{j_1=0}^2 \sum_{j_2=0}^2 \dots \sum_{j_M=0}^2 (i_k) \tau_{(j_1^t, \dots, j_M^t), k}^{(r)} = N_1,$$

$$\tau_{(j_1^t, \dots, j_M^t)_{+(0)}}^{(r)} = \sum_{k=1}^N (1 - i_k) \tau_{(j_1^t, \dots, j_M^t), k}^{(r)}$$

$$\tau_{++(0)}^{(r)} = \sum_{k=1}^N \sum_{j_1=0}^2 \sum_{j_2=0}^2 \dots \sum_{j_M=0}^2 (1 - i_k) \tau_{(j_1^t, \dots, j_M^t), k}^{(r)} = N_0,$$

$$\tau_{(j_1^t, \dots, j_M^t)_+}^{(r)} = \sum_{k=1}^N \tau_{(j_1^t, \dots, j_M^t), k}^{(r)}$$

We determine:

$$\alpha^{(r+1)} = \ln \left(\tau_{(0, \dots, 0)_{+(1)}}^{(r)} \right) - \ln \left(\tau_{(0, \dots, 0)_{+(0)}}^{(r)} \right), \quad (A6)$$

$$\beta^{(r+1)}$$

$$= \ln \left(N_1 - \tau_{(0, \dots, 0)_{+(1)}}^{(r)} \right) - \ln \left(N_0 - \tau_{(0, \dots, 0)_{+(0)}}^{(r)} \right) - \ln \left(\tau_{(0, \dots, 0)_{+(1)}}^{(r)} \right) + \ln \left(\tau_{(0, \dots, 0)_{+(0)}}^{(r)} \right). \quad (A7)$$

Under the null hypothesis, we have:

$$\alpha = \ln(N_1) - \ln(N_0).$$

For the $(r + 1)^{\text{st}}$ -step estimates of the genotype frequencies , we determine:

$$\mathbf{g}_{(j_1^t, \dots, j_M^t)}^{(r)} = \frac{\sum_{k=1}^N \tau_{(j_1^t, \dots, j_M^t), k}^{(r)}}{N}. \quad (\text{A8})$$

Finally, for the $(r + 1)^{\text{st}}$ -step estimates of the sequence error probabilities, we can show that:

$$\begin{aligned} \varepsilon_0^{(r+1)} &= \frac{\sum_{k=1}^N \left((1 - i_k) \left[\sum_{m=1}^M \sum_{(j_1^t, \dots, j_M^t), j_m^t=0} \left(\tau_{(j_1^t, \dots, j_M^t), k}^{(r)} x_{m,k} \right) + \sum_{m=1}^M \sum_{(j_1^t, \dots, j_M^t), j_m^t=2} \left(\tau_{(j_1^t, \dots, j_M^t), k}^{(r)} (v_{m,k}^t - x_{m,k}) \right) \right] \right)}{\sum_{k=1}^N \left((1 - i_k) \left(\sum_{m=1}^M \sum_{(j_1^t, \dots, j_M^t), j_m^t=0} \left(\tau_{(j_1^t, \dots, j_M^t), k}^{(r)} v_{m,k}^t \right) + \sum_{m=1}^M \sum_{(j_1^t, \dots, j_M^t), j_m^t=2} \left(\tau_{(j_1^t, \dots, j_M^t), k}^{(r)} v_{m,k}^t \right) \right) \right)}, \\ \varepsilon_1^{(r+1)} &= \frac{\sum_{k=1}^N \left((i_k) \left[\sum_{m=1}^M \sum_{(j_1^t, \dots, j_M^t), j_m^t=0} \left(\tau_{(j_1^t, \dots, j_M^t), k}^{(r)} x_{m,k} \right) + \sum_{m=1}^M \sum_{(j_1^t, \dots, j_M^t), j_m^t=2} \left(\tau_{(j_1^t, \dots, j_M^t), k}^{(r)} (v_{m,k}^t - x_{m,k}) \right) \right] \right)}{\sum_{k=1}^N \left((i_k) \left(\sum_{m=1}^M \sum_{(j_1^t, \dots, j_M^t), j_m^t=0} \left(\tau_{(j_1^t, \dots, j_M^t), k}^{(r)} v_{m,k}^t \right) + \sum_{m=1}^M \sum_{(j_1^t, \dots, j_M^t), j_m^t=2} \left(\tau_{(j_1^t, \dots, j_M^t), k}^{(r)} v_{m,k}^t \right) \right) \right)}. \end{aligned} \quad (\text{A9})$$

We comment that these probabilities are locus-independent, that is, the subscripts do not contain the individual locus number. Having said that, the formulas indicate that the error probabilities are computed as a composite of the individual locus data values ($v_{m,k}^t$ and $x_{m,k}$), suggesting an "average" over all loci.

Derivation of test statistic

We use the log-likelihoods for each hypothesis for a given iteration value r to ultimately determine the maximum log-likelihoods under each scenario (Null and Alternative), and then use these maxs to determine the value of the test statistic. We use the notation:

$$\ln(L_{H_d}) = \sum_{k=1}^N \ln[\Pr(\mathbf{x}_k, \mathbf{v}_k^t, i_k^t)] = \text{Log-likelihood Equation of the observed data.}$$

Note:

$$d = 0 \text{ for Null Hypothesis: } \quad (H_0: \beta = 0),$$

$$d = 1 \text{ for Alternative Hypothesis: } \quad (H_1: \beta \neq 0).$$

$\ln(\hat{L}_{H_d})$ = Maximum log-likelihood of the data for each hypothesis. This maximum is achieved by applying the EM algorithm in the following way:

1. Specify a certain number of starting points (i.e., randomly generated vectors $\vec{\psi}$ of parameter settings for α, β , etc.).
2. For each vector $\vec{\psi}$ in Item 1, update the log-likelihoods under each hypothesis until some stopping condition is satisfied, such as:

$$|(r + 1)^{\text{st}} \text{- step of } \ln(L_{H_d}) - (r)^{\text{th}} \text{- step of } \ln(L_{H_d})| < \delta, \quad (2)$$

for some tolerance δ . In this work, we use $\delta = 0.00001$. The maximum log-likelihood is then the $(r)^{\text{th}}$ - step of $\ln(L_{H_d})$. We denote this value by: $\ln(L_{H_d})_{r(\vec{\psi})}$.

a. NOTE: For an arbitrary vector $\vec{\psi}$ in Item 1, if the stopping condition (2) is not met after the maximum number of steps, we define the log-likelihood as:

$\ln(L_{H_d})_{r_{\max}(\vec{\psi})}$, where r_{\max} is the total number of steps specified for the EM algorithm. For example, in the Simulation section in the paper, $r_{\max}=1000$ (Item (xi) in *Simulations*).

3. We define the log-likelihood of the observed data, denoted $\ln(\hat{L}_{H_d})$, as:

$$\ln(\hat{L}_{H_d}) = \max_{\vec{\psi}} \left(\ln(L_{H_d})_{r(\vec{\psi})} \right).$$

$LTT_{ae,NGS} = 2[\ln(\hat{L}_{H_1}) - \ln(\hat{L}_{H_0})]$. As noted above in Item (3), the carat symbol $\hat{\cdot}$ indicates that we have obtained the maximum log-likelihood of the data under the particular hypothesis. Note that $LTT_{ae,NGS}$ is asymptotically a chi-square distribution with 1 degree of freedom. We consider two versions of the $LTT_{ae,NGS}$ statistic. In this work, we allow for differential misclassification in the computation of the $LTT_{ae,NGS}$ statistic. Specifically, the two error model parameters are unconstrained (that is, it may be that $\varepsilon_0 \neq \varepsilon_1$).

Example log-likelihood calculations of observed data

In this appendix, we provide calculations of observed-data log-likelihood calculations, demonstrate how the EM algorithm is implemented, and compute the $LTT_{ae,NGS}$ statistic. We perform calculations for a scenario where we have 6 cases and 4 controls sequenced at a single locus.

Example Data Set

Consider the following set of data:

Table A1. Single locus NGS example data set.

Individual (k)	Phenotype Code (i_k^t)	CV counts ($x_{1,k}$)	Coverage($v_{1,k}^t$)
1	1	0	4
2	1	1	4
3	1	2	4
4	1	2	3
5	1	0	3
6	1	1	3
7	0	0	4
8	0	0	4
9	0	3	3
10	0	2	3

Note that there are 6 cases and 4 controls listed in Table A1, and the coverage values are either 3 or 4.

To compute the log-likelihood under the null and alternative hypotheses, we must compute the probabilities $\Pr(\mathbf{x}_k, \mathbf{v}_k^t, i_k^t)$ (Equation (A1)). Recall that:

$$\Pr(\mathbf{x}_k, \mathbf{v}_k^t, i_k^t) = \sum_{\mathbf{G}^t} \Pr(\mathbf{x}_k, \mathbf{v}_k^t, i_k^t, \mathbf{G}^t). \quad (A10)$$

Since there is only one locus, we may rewrite the summation over \mathbf{G}^t as the sum over j_1^t from 0 to 2 (the number of copies of the CV in the true genotype). Thus, equation (A10) may be rewritten as:

$$\Pr(\mathbf{x}_k, \mathbf{v}_k^t, i_k^t) = \sum_{j_1^t=0}^2 \Pr(\mathbf{x}_k, \mathbf{v}_k^t, i_k^t, j_1^t). \quad (A11)$$

We documented above that:

$$\Pr(\mathbf{x}_k, \mathbf{v}_k^t, i_k^t, j_1^t) = \Pr(\mathbf{x}_k | \mathbf{v}_k^t, i_k^t, j_1^t) \times \Pr(\mathbf{v}_k^t) \times \Pr(i_k^t | j_1^t) \times \mathbf{g}_{j_1^t}. \quad (A12)$$

The terms $\Pr(\mathbf{v}_k^t)$ are fixed for each individual over all iterations of the EM algorithm and so we treat it as a constant from this point forward. We shall not include it in our calculations moving forward. From the work done above (Equation (A3)), including Notation from Table 1, we may rewrite Equation (A12) as:

$$\Pr(\mathbf{x}_k, \mathbf{v}_k^t, i_k^t, j_1^t) = \text{Bin}\left(x_{1,k}; v_{1,k}^t p(i_{1,k}^t, j_m^t)\right) \times f_{i_k^t, j_1^t} \times \mathbf{g}_{j_1^t}. \quad (A13)$$

Equation (A13) is the critical equation for our log-likelihood computations. To compute this probability, we must determine the values $\varepsilon_{i_k^t}$, $f_{i_k^t, j_1^t}$, and $\mathbf{g}_{j_1^t}$. All values are unknown and computed via iteration by the EM algorithm. We use the superscript notation (r) to indicate the r^{th} -iteration step estimate of the parameters. The value $r = 0$ indicates the starting value

(randomly generated) for each parameter. Note that $f_{i_k^t, j_1^t} = \frac{\left(e^{\alpha + \beta w_{j_1^t}}\right)^{(i_k^t)}}{1 + e^{\alpha + \beta w_{j_1^t}}}$, so estimation of this parameter reduces to estimation of the parameters α and β . Under the null hypothesis, $\beta = 0$, so

that $f_{i_k^t, j_1^t} = \frac{(e^\alpha)^{(i_k^t)}}{1 + e^\alpha}$. We computed above that, under the null hypothesis, $\alpha^{(r)} = \ln(N_1) -$

$\ln(N_0)$; that is, $\alpha^{(r)}$ remains constant over all EM iterations r . Under the null, therefore, we can reduce the penetrances $f_{i_k^t, j_1^t}$ to:

$$f_{i_k^t, j_1^t} = f_{i_k^t, *} = \frac{(e^\alpha)^{(i_k^t)}}{1 + e^\alpha}. \quad (A14)$$

In other words, under the null hypothesis that $\beta = 0$, the penetrance values depend only on the phenotype code of the individual.

Log-likelihood of observed data under null hypothesis

Given the information above, we may now perform the log-likelihood calculations of the observed data. We start by randomly generating starting values for the aforementioned parameters. Let:

$$\varepsilon_{i_k^t=0}^{(0)} = 0.073;$$

$$\varepsilon_{i_k^t=1}^{(0)} = 0.005;$$

$$g_{j_1^t=0}^{(0)} = 0.316;$$

$$g_{j_1^t=1}^{(0)} = 0.279;$$

$$g_{j_1^t=2}^{(0)} = 0.405.$$

$$f_{i_k^t=0, j_1^t}^{(0)} = f_{i_k^t=0, *}^{(0)} = \frac{(e^{\alpha t})^{(0)}}{1 + e^{\alpha t}} = \frac{1}{1 + e^{\ln(N_1) - \ln(N_0)}} = \frac{1}{1 + \frac{N_1}{N_0}} = \frac{N_0}{N} = \frac{4}{10} = 0.40.$$

$$f_{i_k^t=1, j_1^t}^{(0)} = f_{i_k^t=1, *}^{(0)} = \frac{(e^{\alpha t})^{(1)}}{1 + e^{\alpha t}} = \frac{\frac{N_1}{N_0}}{1 + \frac{N_1}{N_0}} = \frac{N_1}{N} = \frac{6}{10} = 0.60.$$

(Recall that the superscripts here are the Iteration Steps (r). For Individual 1, using Equation (A3) and the values above, we have, at EM Iteration step $r = 0$, the log-likelihood of the observed data is:

$$\begin{aligned} \ln[\Pr(x_{1,1}, v_{1,1}^t, i_1^t)] &= \ln \left[\sum_{j_1^t=0}^2 \Pr(\mathbf{x}_1, \mathbf{v}_1^t, i_1^t, j_1^t) \right] \\ &= \ln \left[\sum_{j_1^t=0}^2 \text{Bin} \left(x_{1,1}; v_{1,1}^t, \left(\frac{2-j_1^t}{2} \varepsilon_{i_1^t}^{(0)} + \frac{j_1^t}{2} (1 - \varepsilon_{i_1^t}^{(0)}) \right) \right) \times f_{i_1^t, *}^{(0)} g_{j_1^t}^{(0)} \right]. \end{aligned} \quad (A15)$$

From Table 1, we know $i_1^t = 1$, $x_{1,1} = 0$, $v_{1,1}^t = 4$. It follows that the log-likelihood (A15) is equal to:

$$\begin{aligned} &= \ln \left[\text{Bin}(0; 4, \varepsilon_1^{(0)}) \times f_{1,*}^{(0)} g_0^{(0)} + \text{Bin} \left(0; 4, \frac{1}{2} \right) \times f_{1,*}^{(0)} g_1^{(0)} + \text{Bin} \left(0; 4, (1 - \varepsilon_1^{(0)}) \right) \times f_{1,*}^{(0)} g_2^{(0)} \right], \\ &= \ln[\text{Bin}(0; 4, 0.005) \times (0.60)(0.316) + \text{Bin}(0; 4, 0.50) \times (0.60)(0.279) + \text{Bin}(0; 4, 0.995) \\ &\quad \times (0.60)(0.405)], \\ &= \ln[(0.980)(0.60)(0.316) + (0.063)(0.60)(0.279) + (0)(0.60)(0.405)], \\ &= \ln[0.186 + 0.010], \\ &= -1.628. \end{aligned}$$

The log-likelihoods of every other individual are computed in the same fashion. To provide an example of someone with different data than the first individual, we choose Individual 10. This individual is unaffected, has coverage equal to 3, and has a non-zero CV count (2).

From Table 1, we have $i_{10}^t = 0$, $x_{10,1} = 2$, $v_{10,1}^t = 3$. It follows that the log-likelihood (A15) is equal to:

$$\begin{aligned} &= \ln \left[\text{Bin}(2; 3, \varepsilon_0^{(0)}) \times f_{0,*}^{(0)} g_0^{(0)} + \text{Bin}(2; 3, 0.5) \times f_{0,*}^{(0)} g_1^{(0)} + \text{Bin} \left(2; 3, (1 - \varepsilon_0^{(0)}) \right) \right. \\ &\quad \left. \times f_{0,*}^{(0)} g_2^{(0)} \right], \end{aligned}$$

$$\begin{aligned}
&= \ln[\text{Bin}(2; 3, 0.073) \times (0.40)(0.316) + \text{Bin}(2; 3, 0.50) \times (0.40)(0.279) + \text{Bin}(2; 3, 0.927) \\
&\quad \times (0.40)(0.405)], \\
&= \ln[(0.015)(0.40)(0.316) + (0.375)(0.40)(0.279) + (0.188)(0.40)(0.405)], \\
&= \ln[0.002 + 0.042 + 0.030], \\
&= -2.601.
\end{aligned}$$

Log-likelihoods for the other individuals may be computed similarly. We present log-likelihood values for all individuals in Table A2.

Table A2. Log-likelihood values for all individuals for Iteration Step 0.

Individual (k)	$\ln[\text{Pr}(x_{1,k}, v_{1,k}^t, i_k^t)]$
1	-1.628
2	-3.088
3	-2.767
4	-2.712
5	-1.572
6	-2.724
7	-2.299
8	-2.299
9	-1.945
10	-2.601
Total	-23.635

As noted in the table, we compute a total log-likelihood of -23.635.

Critically important parameters for the EM-parameter-estimate updates (and consequently, the updated log-likelihoods) are the posterior probabilities $\tau_{j_1^t, k}^{(r)}$. These probabilities are used to determine the updates. To compute the terms, $\tau_{j_1^t, k}^{(r)}$, we use the formula provided above:

$$\tau_{j_1^t, k}^{(r)} = \frac{f_{i_k^t, * }^{(0)} g_{j_1^t}^{(0)} \text{Bin}\left(x_{1,1}; v_{1,1}^t; \left(\frac{2-j_1^t}{2} \varepsilon_{i_1^t}^{(0)} + \frac{j_1^t}{2} (1 - \varepsilon_{i_1^t}^{(0)})\right)\right)}{\sum_{s_1^t=0}^2 f_{i_k^t, * }^{(0)} g_{s_1^t}^{(0)} \text{Bin}\left(x_{1,1}; v_{1,1}^t; \left(\frac{2-s_1^t}{2} \varepsilon_{i_1^t}^{(0)} + \frac{s_1^t}{2} (1 - \varepsilon_{i_1^t}^{(0)})\right)\right)} = \frac{g_{j_1^t}^{(0)} \text{Bin}\left(x_{1,1}; v_{1,1}^t; \left(\frac{2-j_1^t}{2} \varepsilon_{i_1^t}^{(0)} + \frac{j_1^t}{2} (1 - \varepsilon_{i_1^t}^{(0)})\right)\right)}{\sum_{s_1^t=0}^2 g_{s_1^t}^{(0)} \text{Bin}\left(x_{1,1}; v_{1,1}^t; \left(\frac{2-s_1^t}{2} \varepsilon_{i_1^t}^{(0)} + \frac{s_1^t}{2} (1 - \varepsilon_{i_1^t}^{(0)})\right)\right)}.$$

Let us perform this calculation for Individuals 1 and 10. We compute the singleton terms

$$f_{i_k^t, * }^{(0)} g_{s_1^t}^{(0)} \text{Bin}\left(x_{1,1}; v_{1,1}^t; \left(\frac{2-s_1^t}{2} \varepsilon_{i_1^t}^{(0)} + \frac{s_1^t}{2} (1 - \varepsilon_{i_1^t}^{(0)})\right)\right) \text{ for } s_1^t = 0, 1, 2 \text{ first. Once we have these}$$

values, it is straightforward to compute each of the $\tau_{j_1^t, k}^{(r)}$. For Individual 1, we have:

$$\begin{aligned} f_{1, * }^{(0)} g_0^{(0)} \text{Bin}\left(x_{1,1}; v_{1,1}^t; \left(\frac{2-0}{2} \varepsilon_1^{(0)} + \frac{0}{2} (1 - \varepsilon_1^{(0)})\right)\right) & \quad (s_1^t = 0) \\ = (0.60)(0.316) \times \text{Bin}(0; 4, 0.005) & \\ = 0.186. & \end{aligned}$$

$$\begin{aligned} f_{1, * }^{(0)} g_1^{(0)} \text{Bin}\left(x_{1,1}; v_{1,1}^t; \left(\frac{2-1}{2} \varepsilon_1^{(0)} + \frac{1}{2} (1 - \varepsilon_1^{(0)})\right)\right) & \quad (s_1^t = 1) \\ = (0.60)(0.279) \times \text{Bin}(0; 4, 0.5) & \\ = 0.010. & \end{aligned}$$

$$\begin{aligned} f_{1, * }^{(0)} g_2^{(0)} \text{Bin}\left(x_{1,1}; v_{1,1}^t; \left(\frac{2-2}{2} \varepsilon_1^{(0)} + \frac{2}{2} (1 - \varepsilon_1^{(0)})\right)\right) & \quad (s_1^t = 2) \\ = (0.60)(0.405) \times \text{Bin}(0; 4, 0.995) & \\ = 1.52\text{E} - 10 \approx 0. & \end{aligned}$$

The sum of these terms is $0.186 + 0.010 = 0.196$.

Thus,

$$\tau_{0,1}^{(0)} = \frac{0.186}{0.196} = 0.947,$$

$$\tau_{1,1}^{(0)} = \frac{0.010}{0.196} = 0.053,$$

$$\tau_{2,1}^{(0)} = \frac{0.000}{0.196} = 0.000.$$

Based on these calculations, we conclude that the most likely genotype code for individual 1 is 0 (from just one iteration!).

For Individual 10, we have:

$$\begin{aligned} f_{0,*}^{(0)} g_0^{(0)} \text{Bin} \left(x_{1,10}; v_{1,10}^t, \left(\frac{2-0}{2} \varepsilon_0^{(0)} + \frac{0}{2} (1 - \varepsilon_0^{(0)}) \right) \right) & \quad (s_1^t = 0) \\ = (0.40)(0.316) \times \text{Bin}(2; 3, 0.073), & \\ = 0.002. & \end{aligned}$$

$$\begin{aligned} f_{0,*}^{(0)} g_1^{(0)} \text{Bin} \left(x_{1,10}; v_{1,10}^t, \left(\frac{2-1}{2} \varepsilon_0^{(0)} + \frac{1}{2} (1 - \varepsilon_0^{(0)}) \right) \right) & \quad (s_1^t = 1) \\ = (0.40)(0.279) \times \text{Bin}(2; 3, 0.50), & \\ = 0.042. & \end{aligned}$$

$$\begin{aligned} f_{0,*}^{(0)} g_2^{(0)} \text{Bin} \left(x_{1,10}; v_{1,10}^t, \left(\frac{2-2}{2} \varepsilon_0^{(0)} + \frac{2}{2} (1 - \varepsilon_0^{(0)}) \right) \right) & \quad (s_1^t = 2) \\ = (0.40)(0.405) \times \text{Bin}(2; 3, 0.927), & \\ = 0.030. & \end{aligned}$$

The sum of these terms is 0.074.

Thus,

$$\tau_{0,10}^{(0)} = \frac{0.002}{0.074} = 0.025,$$

$$\tau_{1,10}^{(0)} = \frac{0.042}{0.074} = 0.564,$$

$$\tau_{2,10}^{(0)} = \frac{0.030}{0.074} = 0.411.$$

Based on these calculations, we can only reasonably rule out genotype code 0 for Individual 10. The other genotype probabilities are relatively large. The posterior probabilities may be determined in a similar fashion for the other individuals. We report all posterior probabilities in Table A3.

Table A3. Posterior probabilities of a given genotype for all individuals in Table A1 for Iteration 0.

Individual	$\tau_{0,k}^{(0)}$	$\tau_{1,k}^{(0)}$	$\tau_{2,k}^{(0)}$
1	0.947	0.053	8E-10
2	0.082	0.918	3E-06
3	0.000	0.999	0.001
4	0.000	0.946	0.054
5	0.899	0.101	1E-07
6	0.043	0.957	0.000
7	0.930	0.070	5E-05
8	0.930	0.070	5E-05
9	0.000	0.098	0.902
10	0.025	0.564	0.411

Using the formula above, we compute the 1st Iteration Step estimates for the genotype

frequencies as: $g_{j_1^t}^{(1)} = \frac{\sum_{k=1}^N \tau_{j_1^t,k}^{(0)}}{N}$. To determine, $g_0^{(1)}$, we sum up the first column of Table A3 and

divide by $N = 10$. We get: $g_0^{(1)} = \frac{3.858}{10} = 0.386$. Similarly, $g_1^{(1)} = \frac{4.774}{10} = 0.477$, $g_2^{(1)} = \frac{1.368}{10} =$

0.137. As a check, we note that the frequencies sum to 1.0.

To compute the error parameter updates, we use the formula (A9). In our special case of a single locus,

$$\begin{aligned}
\varepsilon_0^{(r+1)} &= \frac{\sum_{k=1}^N \left\{ (1 - i_k) \left(\tau_{0,k}^{(r)} x_{1,k} + \tau_{2,k}^{(r)} (v_{1,k}^t - x_{1,k}) \right) \right\}}{\sum_{k=1}^N \left\{ (1 - i_k) \left(\tau_{0,k}^{(r)} + \tau_{2,k}^{(r)} \right) v_{1,k}^t \right\}}, \\
\varepsilon_1^{(r+1)} &= \frac{\sum_{k=1}^N \left\{ i_k \left(\tau_{0,k}^{(r)} x_{1,k} + \tau_{2,k}^{(r)} (v_{1,k}^t - x_{1,k}) \right) \right\}}{\sum_{k=1}^N \left\{ i_k \left(\tau_{0,k}^{(r)} + \tau_{2,k}^{(r)} \right) v_{1,k}^t \right\}}.
\end{aligned} \tag{A16}$$

In Table A4, we compute the numerator and denominators for each of the error-parameter estimates. The interested reader can verify these values. Using the columns total, we compute that the updated error-parameter estimates for Iteration Step $r = 1$ are:

$$\varepsilon_0^{(1)} = \frac{0.4627}{11.4594} = 0.040,$$

$$\varepsilon_1^{(1)} = \frac{0.1822}{7.1097} = 0.026.$$

Table A4. Numerator and denominator terms from formula (A16).

Individual	Numerator-Term ($\varepsilon_0^{(1)}$)	Denominator-Term ($\varepsilon_0^{(1)}$)	Numerator-Term ($\varepsilon_1^{(1)}$)	Denominator-Term ($\varepsilon_1^{(1)}$)
1	0	0	3.1E-09	3.787
2	0	0	0.0820	0.3278
3	0	0	0.0020	0.0041
4	0	0	0.0548	0.1637
5	0	0	4.4E-07	2.6978
6	0	0	0.0435	0.1296
7	0.0002	3.7219	0	0
8	0.0002	3.7219	0	0
9	0.0010	2.7074	0	0
10	0.4613	1.3082	0	0
Total	0.4627	11.4594	0.1822	7.1097

Using these values and following the same steps listed above (with the updated parameter estimates), we obtain the following log-likelihoods (including the total log-likelihood) for Iteration Step 1.

Table A5. Log-likelihood values for all individuals for Iteration Step 1 (null hypothesis).

Individual (k)	$\ln[\Pr(x_{1,k}, v_{1,k}^t, i_k^t)]$
1	-1.485
2	-2.369
3	-2.220
4	-2.173
5	-1.387
6	-2.084
7	-1.946
8	-1.946
9	-2.628
10	-2.546
Total	-20.783

As noted in the table, we compute a total log-likelihood of -20.783.

Using the inequality (2) listed above and the information from the two tables, we have:

$$|(1)^{\text{st}}\text{- step of } \ln(L_{H_0}) - (0)^{\text{th}}\text{- step of } \ln(L_{H_0})| = |-20.783 - (-23.635)| = 2.852.$$

The interested reader can check that it takes up to the $r = 25^{\text{th}}$ Iteration Step before the tolerance inequality is achieved. The $(25)^{\text{st}}$ - step of $\ln(L_{H_0})$ is -19.48327932, and $(24)^{\text{th}}$ - step of $\ln(L_{H_0})$ is -19.48328661, so the difference is approximately 7.30E-06, less than the tolerance of 1E-05.

It is interesting to note that the final error-parameters for this set of starting points are: $\varepsilon_0^{(25)} = 0.0, \varepsilon_1^{(25)} = 0.224$. This result suggests differential misclassification. Of course, these estimates must be taken “with a grain of salt” since the sample size is very small. Also, the genotype frequency estimates at this step are: $g_0^{(25)} = 0.587, g_1^{(25)} = 0.309, g_2^{(25)} = 0.104$. We

use these parameter estimates as starting values for the log-likelihood calculations under the alternative hypothesis.

Log-likelihood of observed data under alternative hypothesis

The calculations for log-likelihoods under the alternative hypothesis are nearly the same as under the null hypothesis, with the exception that the α and β parameters (and hence the penetrances) are updated with each iteration step.

For the 0th Iteration Step under the alternative, we use starting parameters determined by the last iteration step of the null ($r = 25$). As noted above, we have:

$$\varepsilon_0^{(0)} = 0.000;$$

$$\varepsilon_1^{(0)} = 0.224;$$

$$g_0^{(0)} = 0.587;$$

$$g_1^{(0)} = 0.309;$$

$$g_2^{(0)} = 0.104.$$

$$f_{0,j_1^t}^{(0)} = 0.40, 0 \leq j_1^t \leq 2.$$

$$f_{1,j_1^t}^{(0)} = 0.60, 0 \leq j_1^t \leq 2.$$

For Individual 1, since the 0th Step penetrances are the same as under the null, we may use Equation (A3) to compute the log-likelihood. We get:

$$\ln \left[Bin(0; 4, \varepsilon_1^{(0)}) \times f_{1,*}^{(0)} g_0^{(0)} + Bin\left(0; 4, \frac{1}{2}\right) \times f_{1,*}^{(0)} g_1^{(0)} + Bin\left(0; 4, (1 - \varepsilon_1^{(0)})\right) \times f_{1,*}^{(0)} g_2^{(0)} \right],$$

$$\begin{aligned}
&= \ln[Bin(0; 4, 0.224) \times (0.60)(0.587) + Bin(0; 4, 0.50) \times (0.60)(0.309) + Bin(0; 4, 0.776) \\
&\quad \times (0.60)(0.104)], \\
&= \ln[(0.363) \times (0.60)(0.587) + (0.0625) \times (0.60)(0.309) + (0.003) \times (0.60)(0.104)], \\
&= \ln[0.140], \\
&= -1.969.
\end{aligned}$$

As above, we also consider Individual 10.

The log-likelihood is equal to:

$$\begin{aligned}
&= \ln \left[Bin(2; 3, \varepsilon_0^{(0)}) \times f_{0,*}^{(0)} g_0^{(0)} + Bin(2; 3, 0.5) \times f_{0,*}^{(0)} g_1^{(0)} + Bin(2; 3, (1 - \varepsilon_0^{(0)})) \right. \\
&\quad \left. \times f_{0,*}^{(0)} g_2^{(0)} \right], \\
&= \ln[Bin(2; 3, 0.000) \times (0.40)(0.587) + Bin(2; 3, 0.50) \times (0.40)(0.309) + Bin(2; 3, 1.00) \\
&\quad \times (0.40)(0.104)], \\
&= \ln[(0.375) \times (0.40)(0.309)], \\
&= \ln[0.046], \\
&= -3.073.
\end{aligned}$$

Log-likelihoods for the other individuals may be computed similarly. We present log-likelihood values for all individuals in Table A6.

Table A6. Log-likelihood values for all individuals for Iteration Step 0 under alternative hypothesis.

Individual (k)	$\ln[\Pr(x_{1,k}, v_{1,k}^t, i_k^t)]$
1	-1.970
2	-1.630
3	-1.934
4	-1.996
5	-1.669
6	-1.518
7	-1.417
8	-1.417
9	-2.862
10	-3.073
Total	-19.483

As noted in the table, we compute a total log-likelihood of -19.483.

To compute the terms, $\tau_{j_1^t, k}^{(r)}$, as with the null, we use the formula:

$$\tau_{j_1^t, k}^{(r)} = \frac{f_{i_k^t, j_1^t}^{(0)} g_{j_1^t}^{(0)} \text{Bin}\left(x_{1,1}; v_{1,1}^t, \left(\frac{2-j_1^t}{2} \varepsilon_{i_1^t}^{(0)} + \frac{j_1^t}{2} (1 - \varepsilon_{i_1^t}^{(0)})\right)\right)}{\sum_{s_1^t=0}^2 \left(f_{i_k^t, s_1^t}^{(0)} g_{s_1^t}^{(0)} \text{Bin}\left(x_{1,1}; v_{1,1}^t, \left(\frac{2-s_1^t}{2} \varepsilon_{i_1^t}^{(0)} + \frac{s_1^t}{2} (1 - \varepsilon_{i_1^t}^{(0)})\right)\right) \right)}$$

The main difference between this calculation and the one for the null is that the penetrance value

$f_{i_k^t, j_1^t}^{(0)}$ now depends upon the genotype.

Let us perform this calculation for Individuals 1 and 10. For Individual 1, we have:

$$\begin{aligned} & f_{1,0}^{(0)} g_0^{(0)} \text{Bin}\left(x_{1,1}; v_{1,1}^t, \left(\frac{2-0}{2} \varepsilon_1^{(0)} + \frac{0}{2} (1 - \varepsilon_1^{(0)})\right)\right) \\ &= (0.60)(0.587) \times \text{Bin}(0; 4, 0.224) \quad (s_1^t = 0) \\ &= 0.128. \end{aligned}$$

$$\begin{aligned}
& f_{1,1}^{(0)} g_1^{(0)} \text{Bin} \left(x_{1,1}; v_{1,1}^t, \left(\frac{2-1}{2} \varepsilon_1^{(0)} + \frac{1}{2} (1 - \varepsilon_1^{(0)}) \right) \right) \quad (s_1^t = 1) \\
& = (0.60)(0.309) \times \text{Bin}(0; 4, 0.5) \\
& = 0.012.
\end{aligned}$$

$$\begin{aligned}
& f_{1,2}^{(0)} g_2^{(0)} \text{Bin} \left(x_{1,1}; v_{1,1}^t, \left(\frac{2-2}{2} \varepsilon_1^{(0)} + \frac{2}{2} (1 - \varepsilon_1^{(0)}) \right) \right) \quad (s_1^t = 2) \\
& = (0.60)(0.104) \times \text{Bin}(0; 4, 0.776) \\
& = 0.2\text{E} - 05.
\end{aligned}$$

The sum of these terms is $0.128 + 0.012 + 0.00002 = 0.140$.

Thus,

$$\tau_{0,1}^{(0)} = \frac{0.128}{0.140} = 0.914,$$

$$\tau_{1,1}^{(0)} = \frac{0.012}{0.140} = 0.086,$$

$$\tau_{2,1}^{(0)} = \frac{0.000}{0.140} = 0.000.$$

Based on these calculations, we conclude that the most likely genotype code for individual 1 is 0.

For Individual 10, we have:

$$\begin{aligned}
& f_{0,0}^{(0)} g_0^{(0)} \text{Bin} \left(x_{1,10}; v_{1,10}^t, \left(\frac{2-0}{2} \varepsilon_0^{(0)} + \frac{0}{2} (1 - \varepsilon_0^{(0)}) \right) \right) \quad (s_1^t = 0) \\
& = (0.40)(0.587) \times \text{Bin}(2; 3; 0.000) \\
& = 0.000.
\end{aligned}$$

$$\begin{aligned}
& f_{0,1}^{(0)} g_1^{(0)} \text{Bin} \left(x_{1,10}; v_{1,10}^t, \left(\frac{2-1}{2} \varepsilon_0^{(0)} + \frac{1}{2} (1 - \varepsilon_0^{(0)}) \right) \right) \quad (s_1^t = 1) \\
& = (0.40)(0.309) \times \text{Bin}(2; 3; 0.500) \\
& = 0.012.
\end{aligned}$$

$$\begin{aligned}
& f_{0,2}^{(0)} g_2^{(0)} \text{Bin} \left(x_{1,10}; v_{1,10}^t; \left(\frac{2-2}{2} \varepsilon_0^{(0)} + \frac{2}{2} (1 - \varepsilon_0^{(0)}) \right) \right) \quad (s_1^t = 2) \\
& = (0.40)(0.104) \times \text{Bin}(2; 3; 1.000) \\
& = 0.000.
\end{aligned}$$

The sum of these terms is 0.012.

Thus,

$$\tau_{0,10}^{(0)} = \frac{0.000}{0.012} = 0.000,$$

$$\tau_{1,10}^{(0)} = \frac{0.012}{0.012} = 1.000,$$

$$\tau_{2,10}^{(0)} = \frac{0.000}{0.012} = 0.000.$$

These results are consistent with the parameter estimates. Note that the error estimate for controls is $\varepsilon_0^{(0)} = 0.000$. Because the observed number of variants is 2, we know with certainty that the genotype must be a heterozygote. When there is no sequence error, it is impossible for either homozygote to have any number of counts other than 0 or $v_{m,k}^t$.

Table A7. Posterior probabilities of a given genotype for all individuals in Table 1 for Iteration 0 (alternative hypothesis).

Individual	$\tau_{0,k}^{(0)}$	$\tau_{1,k}^{(0)}$	$\tau_{2,k}^{(0)}$
1	0.916	0.083	0.001
2	0.753	0.236	0.011
3	0.441	0.480	0.078
4	0.303	0.511	0.186
5	0.873	0.122	0.004
6	0.650	0.317	0.033
7	0.968	0.032	0.000
8	0.968	0.032	0.000
9	0.000	0.270	0.730
10	0.000	1.000	0.000

We compute the 1st Iteration Step estimates for the genotype frequencies as: $g_{j_1^t}^{t(1)} = \frac{\sum_{k=1}^N \tau_{j_1^t, k}^{(0)}}{N}$. As

in the null situation, we sum up the first column of Table A7 and divide by $N = 10$. We get:

$$g_0^{(1)} = \frac{5.872}{10} = 0.587. \text{ Similarly, } g_1^{(1)} = \frac{3.084}{10} = 0.308, g_2^{(1)} = \frac{1.0441}{10} = 0.104.$$

To compute the error parameter updates, we use the formula (A16). The interested reader can verify that the component values necessary to compute the updated error parameter estimates are the ones provided in Table A8 . Using the columns total, we compute that the updated error-parameter estimates for Iteration Step $r = 1$ are:

$$\varepsilon_0^{(1)} = \frac{0.000}{9.936} = 0.000,$$

$$\varepsilon_1^{(1)} = \frac{3.350}{14.951} = 0.224.$$

One can check that the log-likelihood tolerance condition is met at Iteration Step $r = 53$.

At that step, the maximum log-likelihood under the alternative is -18.8506. Thus, the value of the test statistic is $LTT_{ae,NGS} = 2 \ln[-18.851 - (-19.483)] = 1.265$. The corresponding p-value using a central chi-square distribution as the null distribution is 0.261.

Table A8. Numerator and denominator terms from formula (A16).

Individual	Numerator-Term ($\varepsilon_0^{(1)}$)	Denominator-Term ($\varepsilon_0^{(1)}$)	Numerator-Term ($\varepsilon_1^{(1)}$)	Denominator-Term ($\varepsilon_1^{(1)}$)
1	0	0	0.005	3.668
2	0	0	0.786	3.055
3	0	0	1.040	2.079
4	0	0	0.792	1.467
5	0	0	0.011	2.632
6	0	0	0.717	2.050
7	0	3.873	0	0
8	0	3.873	0	0
9	0	2.190	0	0
10	0	0	0	0
Total	0	9.936	3.350	14.951

Genotype probability calculations for SKAT and CMAT

The genotype probabilities are calculated from the posterior probabilities using Bayes' Rule.

Given CV count data $x_k = (x_{1,k}, \dots, x_{M,k})$, the genotype probabilities are obtained as follows.

At each indexed site m ($m = 1, 2, \dots, M$), we have:

$$\Pr(G_{m,k} = j | x_{m,k}, i_k^t, v_{m,k}^t) = \frac{\Pr(G_{m,k} = j) \Pr(x_{m,k} | G_{m,k} = j, i_k^t, v_{m,k}^t)}{\sum_{l=0}^2 \Pr(G_{m,k} = l) \Pr(x_{m,k} | G_{m,k} = l, i_k^t, v_{m,k}^t)}. \quad (A17)$$

We may include the sequence error parameter ε in equation (A17) by rewriting it as:

$$\Pr(G_{m,k} = j | x_{m,k}, i_k^t, v_{m,k}^t) = \frac{\Pr(G_{m,k} = j) \Pr(x_{m,k} | G_{m,k} = j, \varepsilon_{i_k^t}, v_{m,k}^t)}{\sum_{l=0}^2 \Pr(G_{m,k} = l) \Pr(x_{m,k} | G_{m,k} = l, \varepsilon_{i_k^t}, v_{m,k}^t)}. \quad (A18)$$

Since the equation (A18) requires the probability $\Pr(G_{m,k} = j)$ and the sequence error parameters $\varepsilon_{i_k^t}$, we estimate these values internally using the $LTT_{ae,NGS}$ statistic.