

1 **freqpcr: interval estimation of population allele frequency**
2 **based on quantitative PCR $\Delta\Delta Cq$ measures from bulk samples**

3

4 Allele-frequency estimation based on $\Delta\Delta Cq$

5

6 Masaaki Sudo^{*1}, Masahiro Osakabe²

7

8 ¹ Tea Pest Management Unit, Institute of Fruit Tree and Tea Science, NARO: Kanaya Tea Research Station,
9 2769, Shishidoi, Kanaya, Shimada, Shizuoka 428-8501, Japan

10 ² Laboratory of Ecological Information, Graduate School of Agriculture, Kyoto University: Kyoto, Japan

11

12 * To whom correspondence: masaaki@sudori.info

13 ORCID ID:

14 0000-0001-9834-9857 (Masaaki Sudo)

15 Telephone: +81-547-45-4419

16

17 Manuscript types: article

18

19

20 **Abstract**

21 PCR techniques, both quantitative (qPCR) and non-quantitative, have been used to estimate allele frequency
22 in a population. However, the labor required to sample more individuals and handle each sample makes it
23 difficult to quantify rare mutations, such as pesticide-resistance genes at the early stages of resistance
24 development. Pooling DNA from multiple individuals as a “bulk sample” may reduce handling costs. The
25 output of qPCR on a bulk sample, however, contains uncertainty owing to variations in DNA yields from
26 each individual, in addition to measurement error. In this study, we developed a statistical model for the
27 interval estimation of allele frequency via $\Delta\Delta Cq$ -based qPCR analyses of multiple bulk samples taken from a
28 population. We assumed a gamma distribution as the individual DNA yield and developed an R package for
29 parameter estimation, which was verified with real DNA samples from acaricide-resistant spider mites, as
30 well as a numerical simulation. Our model resulted in unbiased point estimates of the allele frequency
31 compared with simple averaging of the $\Delta\Delta Cq$ values, and their confidence intervals suggested collecting
32 more samples from individuals and pooling them may produce higher precision than individual PCR tests
33 with moderate sample sizes.

34 **Keywords:** Real-time polymerase chain reaction, group testing, confidence interval, maximum likelihood
35 estimation, R language

36

37 Introduction

38 Estimating the frequency of certain alleles in populations is one of the key techniques not only in population
39 genetics and molecular ecology, but also in agricultural and regulatory sciences (Falconer 1960; Kim et al.
40 2011; Yamamura and Hino 2007). In applied entomology, field monitoring has been performed to detect
41 resistance genes of arthropod pests to pesticides and genetically modified (GM) insecticidal plants, such as
42 *Bt* crops (Andow and Alstad 1998; Sonoda et al. 2017).

43 Entomologists have traditionally estimated resistance allele frequencies via bioassays (Gould et al. 1997;
44 Li et al. 2016; Tabashnik et al. 2000), in which insects directly collected from fields or their offspring reared
45 in laboratories are exposed to chemical compounds of interest to obtain measurements, such as mortality
46 rate. However, bioassays have drawbacks associated with the treatment of living organisms. It is usually
47 labor-intensive and time-consuming. Although the resistance level can be directly measured using a bioassay
48 as the mortality of tested individuals, additional information including the dominance of the resistance gene
49 is required to estimate the allele frequency.

50 In accordance with the development of genome-wide association studies on resistance genes (Frensch-
51 Constant 2013; Snoeck et al. 2019; Sugimoto et al. 2020), molecular diagnostics have rapidly developed in
52 recent years (Donnelly et al. 2016; Samayoa et al. 2015; Toda et al. 2017). To quantify the resistance-
53 associated point mutation at the population scale, the most fundamental molecular technique is an individual-
54 based polymerase chain reaction (PCR) analysis (Toda et al. 2017). If the alleles are distributed randomly in
55 the target population, a simple binomial assumption enables us to estimate the population allele frequency
56 and its confidence interval. However, it may not be realistic to extract and analyze DNA individually,
57 especially when dealing with many samples from multiple sites or when we need to estimate mutation
58 frequency, which is rare in the population (below 1%), as is often the case in the early phase of resistance
59 development.

60 Although rearing living insects is no longer needed, molecular diagnostics still require a silver bullet to
61 reduce the time and cost of handling multiple samples while guaranteeing estimation precision and accuracy.
62 The use of a “bulk sample” (i.e., pooling multiple individual samples and processing a single DNA extract),
63 in coordination with statistical methods, such as group testing, may help. Quantitative PCR (qPCR), based on
64 real-time PCR, is also used for the point estimation of allele frequency (Germer et al. 2000). Osakabe *et al.*
65 (2017) and Maeoka *et al.* (2020) developed diagnostic methods for acaricide resistance in the two-spotted
66 spider mite, *Tetranychus urticae* Koch (Acari: Tetranychidae), where they used a bulk sample to measure the
67 frequency of the resistant point mutation in field mite populations. To calculate the point estimate, these
68 studies compared the relative quantity of the resistance allele with an internal reference (housekeeping gene)
69 in the sample, which is known as the $\Delta\Delta Cq$ method (Livak and Schmittgen 2001).

70 In this study, we propose a statistical procedure to obtain the interval estimate of allele frequency using
71 $\Delta\Delta Cq$ -based qPCR analyses over multiple bulk samples taken from a population. We first introduced the
72 random error structure to approximate the amounts of the two alleles (wild-type and mutant) and their ratios
73 in the bulk DNA sample. Thereafter, we formulated how the relative amounts of the two alleles in a sample
74 solution resulted in the Cq measurements through qPCR analysis. Finally, we combined the models for
75 individual DNA yields and the model for $\Delta\Delta Cq$ -based qPCR analysis. We developed a maximum likelihood
76 estimation (MLE) procedure to estimate an allele frequency implemented using the R language. The package
77 source is available on the Internet (<https://github.com/sudoms/freqqpcr>).

78 Model

79 *Approximation of allele quantities contained in a bulk DNA sample*

80 When DNA is directly extracted from the whole body of a living organism, the DNA yield is roughly
81 proportional to its body weight (Chen et al. 2010). For insects, the intra-population frequency distribution of
82 body weight is often approximated using a unimodal and right-skewed continuous distribution, typically
83 lognormal or gamma distribution (May 1976; Rakovski et al. 2011; Knapp 2016). A study suggested that
84 body weights are distributed lognormally in many non-social insect species (Gouws et al. 2011).

85 In this study, we adopted a gamma distribution, instead of a lognormal, to approximate the DNA amount
86 per individual organism for two reasons. First, it is difficult to distinguish which distribution a real
87 population obeys when the sample size is small. They are considered interchangeable (Wiens 1999; Kundu
88 and Manglick 2005). Second, the sum and proportion of independent gamma distributions have closed forms
89 under certain conditions. Assuming, let X ($X \geq 0$) be the DNA yield per single locus per individual:

$$90 \text{Ga}(X|k, \theta) = \frac{1}{\Gamma(k)} \left(\frac{1}{\theta}\right)^k X^{k-1} \exp\left(-\frac{X}{\theta}\right),$$

91 *Eq. 1*

92 where $\Gamma(\cdot)$ denotes the gamma function. The parameters k and θ ($k, \theta > 0$) are the shape and scale
93 parameters of the gamma distribution, respectively. The mean is given by $k\theta$.

94 Using Eq. 1, let us consider the amounts of allelic DNA in the sample extracted from multiple individuals
95 at once, hereafter referred to as “a bulk sample.” Table 1 lists the variables and parameters of the model
96 structure. For simplicity, we model the case of haploidy in the main text. Appendix A1 describes the
97 approximated formulation for diploids. Now, we have n insects, of which m ($m = 0, 1, \dots, n$) are the
98 genotypes resistant to an insecticide (hereafter denoted by R). The rest $n - m$ had S, the susceptible allele.
99 When we capture insects from a wild population, the size of n is obvious, but m is usually unknown.
100 Assuming random sampling from an infinite population with the R allele at the frequency p , m follows a
101 binomial distribution:

$$102 \text{Bin}(m|n, p) = \frac{n!}{m!(n-m)!} p^m (1-p)^{n-m}.$$

103 *Eq. 2*

104 When the bulk sample contains at least one resistant individual, $X_R = \sum_{i=1}^m X_i$ denotes the total R content. If
105 there is no systematic error in the efficiency of DNA extraction between the genotypes, and if X_i , the
106 individual DNA yield obeys the gamma distribution of Eq. 1, then X_R follows the gamma distribution with
107 the shape parameter mk and scale parameter θ based on the reproductive property. Conversely, the amount
108 of S allele is denoted by $X_S = \sum_{i=m+1}^n X_i$, which follows the gamma distribution with $(n-m)k$ and θ .

$$109 X_S \sim \text{Ga}((n-m)k, \theta),$$

$$110 X_R \sim \text{Ga}(mk, \theta).$$

111 *Eq. 3*

112 When X_R and X_S independently follow gamma distributions with the same scale parameter, the observed
113 allele frequency $Y_R = X_R/(X_S + X_R)$ follows a beta distribution with the shape parameters mk and
114 $(n - m)k$:

$$115 \text{Beta}(Y_R|mk, (n - m)k) = \frac{Y_R^{mk-1}(1 - Y_R)^{(n-m)k-1}}{B(mk, (n - m)k)},$$

116 Eq. 4

117 where $B(\cdot)$ is a beta function. This error structure was originally developed to model allele frequencies
118 measured via quantitative sequencing (Sudo et al. in press). In quantitative sequencing, unlike qPCR, we
119 cannot directly observe the quantities of template DNA (X_R and X_S). Instead, the output from the Sanger
120 sequencer is reflected as Y_R . Although Sudo *et al.* (in press) used Eq. 4 to approximate DNA yield
121 distribution in dead insect bodies on a trap, that is, considering variations in body weight plus post-mortem
122 DNA degradation, it is also applicable to DNA from fresh bodies.

123 ***Relative quantification of DNA by real-time PCR: $\Delta\Delta Cq$ and RED- $\Delta\Delta Cq$ methods***

124 *Relationship between the template DNA amount and qPCR measure*

125 In real-time qPCR, the target molecule is amplified at a nearly constant efficiency until it exhausts
126 nucleotides (dNTPs) to synthesize the new DNA strand. After amplification cycles with an appropriate
127 primer set, the abundance of the initial template DNA was measured as Cq: quantification cycle (Bustin et al.
128 2009), also known as cycle threshold (Ct). According to Livak and Schmittgen (2001), we assume an ideal
129 amplification, where the threshold X_Θ is set within the early exponential amplification phase:

$$130 X_\Theta = X_0 \times (1 + \eta)^\tau.$$

131 Eq. 5

132 Here, X_0 and $1 + \eta$ ($\eta > 0$) denote the initial amount of template DNA and its amplification efficiency,
133 respectively. Standard PCR protocols are designed so that η obtain the range 80% to 120% i.e., doubling in
134 each cycle. The size of Cq, τ , is then defined as:

$$135 \tau = \frac{\ln(X_\Theta) - \ln X_0}{\ln(1 + \eta)}.$$

136 Eq. 6

137 *Relative quantification of template DNA between experimental levels: $\Delta\Delta Cq$ method*

138 The $\Delta\Delta Cq$ (Ct) method (Livak 1997) is the most common method for relative quantification using qPCR. In
139 a typical scenario, an experiment is conducted at two levels (i.e., treated *versus* control [calibrator]) and
140 complementary cDNA libraries are obtained reflecting different gene expression levels at a single target
141 locus (hereafter abbreviated as “TG” or simply T). It is possible to directly compare the cDNA quantities
142 measured by qPCR if a primer set is available to amplify the TG locus. However, there is no guarantee that
143 the samples with different treatments have the same cDNA preparation efficiency.

144 Hence, an internal reference, which is dispensed in accordance with the sample in question, should be
145 included in relative quantification, such as the $\Delta\Delta Cq$ method. The corresponding primer set usually targets
146 the locus of a housekeeping gene (hereafter abbreviated as “HK” or H), a gene that shows a constant
147 expression level and is thus considered the same concentration between treatments. If the experiment had

148 two levels, we amplified at least four samples (two levels, two primer sets for TG and HK loci, ignoring
149 technical replicates). ΔCq is then defined as the difference of the Cq values of “TG – HK” for each treatment
150 level, which is equivalent to the abundance of target cDNA offset by housekeeping gene (= TG / HK) in each
151 sample (Scheffe et al. 2006). Finally, we obtained $\Delta\Delta Cq = \Delta Cq^{\text{treated}} - \Delta Cq^{\text{control}}$ from the Cq measures.
152 Derived from Eq. 6, $2^{-\Delta\Delta Cq}$ gives the relative abundance of template DNA between the treatment levels
153 (Livak and Schmittgen 2001; Pfaffl 2012) ($1 + \eta = 2$ was presupposed there).

154 *Allele frequency estimation from a single bulk sample: RED- $\Delta\Delta Cq$ method*

155 The original $\Delta\Delta Cq$ method compares the quantities of (c)DNA between samples to determine the relative
156 expression levels of the genes of interest. Osakabe *et al.* (2017) expanded it and proposed the “RED-
157 $\Delta\Delta Cq$ method” (RED stands for restriction enzyme digestion), a derivative method that can measure the
158 allele frequency from a single sample solution, to diagnose the regional resistance prevalence of the two-
159 spotted spider mite, *Tetranychus urticae* Koch (Acari: Tetranychidae), to the acaricide etoxazole, which is
160 conferred by an amino acid substitution in chitin synthase 1 (*CHS1*; I1017F) (Van Leeuwen et al. 2010).

161 The RED- $\Delta\Delta Cq$ method also utilized $\Delta\Delta Cq$ as a proxy for relative quantity, but the Cq measurements were
162 all taken from a single bulk sample, which was collected from a population in which each individual
163 possesses R or S. The calibrator was an intact sample containing total DNA (= $X_R + X_S$) on the target locus.
164 The sample in question was the same DNA extract, but digested with restriction endonucleases prior to
165 qPCR analysis. The restriction site is designed to recognize the S allele on the target locus so that the
166 operation digests the major part of S (denoted by $1 - z$: z is a small, but positive variable giving the residual
167 rate). Consequently, we obtained the template amount $X_R + zX_S$ at the target locus after digestion.

168 The samples before and after digestion were also amplified using the HK primer set as an internal
169 reference. In the etoxazole-R diagnosis by Osakabe *et al.* (2017), glyceraldehyde-3-phosphate
170 dehydrogenase (*GAPDH*) was used. Taken together, the single bulk sample results in a quartet of Cq
171 measurements differentiating at the target loci (*CHS1* and *GAPDH*) \times restriction enzyme digestion
172 (undigested and digested).

173 We now formulate the allele frequencies. Let X^{HW} and X^{TW} represent the total amounts of the template
174 DNA at the housekeeping (H) and target (T) loci included in the sample without digestion, the state denoted
175 by W.

$$\begin{aligned} X^{\text{HW}} &= X_S + X_R, \\ X^{\text{TW}} &= \delta_T(X_S + X_R). \end{aligned}$$

177 *Eq. 7*

178 The coefficient δ_T ($\delta_T > 0$) provides the relative content of the target gene to the housekeeping gene in
179 genomic DNA (the difference in the DNA extraction efficiencies is also included). After digestion (state D),
180 X^{HD} and X^{TD} denote the DNA amounts at the H and T loci, respectively:

$$\begin{aligned} X^{\text{HD}} &= \delta_B(X_S + X_R), \\ X^{\text{TD}} &= \delta_B\delta_T(zX_S + X_R). \end{aligned}$$

182 *Eq. 8*

183 The common coefficient δ_B ($\delta_B > 0$) provides the rate of certain locus-independent changes in the quantities
184 of template DNA accompanying the restriction enzyme treatment.

185 As a result of qPCR, the Cq quartet, τ^{HW} , τ^{TW} , τ^{HD} , and τ^{TD} were obtained. From Eq. 6,

$$\begin{aligned} \tau^{\text{HW}} &= \frac{\ln(X_{\Theta}) - \ln(X_S + X_R)}{\ln(1 + \eta)} + \varepsilon_c, \\ \tau^{\text{TW}} &= \frac{\ln(X_{\Theta}) - \ln\delta_T - \ln(X_S + X_R)}{\ln(1 + \eta)} + \varepsilon_c, \\ \tau^{\text{HD}} &= \frac{\ln(X_{\Theta}) - \ln\delta_B - \ln(X_S + X_R)}{\ln(1 + \eta)} + \varepsilon_c, \\ \tau^{\text{TD}} &= \frac{\ln(X_{\Theta}) - \ln\delta_B - \ln\delta_T - \ln(zX_S + X_R)}{\ln(1 + \eta)} + \varepsilon_c. \end{aligned}$$

188 Eq. 9

189 The actual Cq data contain measurement errors in addition to uncertainty due to experimental operations,
190 such as sample dispensation or PCR amplification. We express these using the common error term
191 $\varepsilon_c \sim \text{N}(0, \sigma_c^2)$, following the normal distribution of mean = 0 and variance = σ_c^2 in the scale of raw Cq values.
192 The validity of this error structure is verified later.

193 The two ΔCq values were then defined as $\Delta\tau^{\text{W}} = \tau^{\text{TW}} - \tau^{\text{HW}}$ and $\Delta\tau^{\text{D}} = \tau^{\text{TD}} - \tau^{\text{HD}}$, respectively. Their
194 $\Delta\Delta\text{Cq}$ are:

$$\Delta\Delta\tau = \Delta\tau^{\text{D}} - \Delta\tau^{\text{W}} = -\frac{\ln\left(\frac{zX_S + X_R}{X_S + X_R}\right)}{\ln(1 + \eta)} + \varepsilon, \quad \varepsilon \sim \text{N}(0, 4\sigma_c^2).$$

196 Eq. 10

197 From Eq. 10, the expected value of $(zX_S + X_R)/(X_S + X_R)$ is calculated as $(1 + \eta)^{-\Delta\Delta\tau}$.

198 The point estimate of the resistance allele frequency, \hat{Y}_R , is defined as $X_R/(X_R + X_S)$ for each bulk sample.
199 When z is much smaller than \hat{Y}_R , the quantity $(zX_S + X_R)/(X_S + X_R) = \hat{Y}_R + z(1 - \hat{Y}_R)$ itself can
200 approximate the frequency, which will be the case with enough digestion time before qPCR. However, the
201 use of the point estimate may introduce a problem in that the size of \hat{Y}_R often exceeds 1 when the R
202 frequency is high and there is a larger error in the Cq measurement (also see the result of Experiment 2).

203 Although the value of $1 + \eta$ may vary on the primer sets, both target and housekeeping loci share the
204 same amplification efficiency in Eq. 9. This is because practical PCR protocols were designed to be $1 + \eta \cong$
205 2. We can also approximately cancel the effect of heterogeneous amplification efficiencies by fitting the size
206 of δ_T the sample sets with known allele ratios (Experiment 1).

207 *Measurement of $\Delta\Delta\text{Cq}$ using allele-specific primer sets*

208 While the RED- $\Delta\Delta\text{Cq}$ method enabled us to measure allele frequency from the bulk sample, enzyme
209 availability is a prerequisite to digest the S-allele-specific restriction site at the target locus. A longer
210 digestion period (3 h) was also required to quantify etoxazole resistance in the protocol by Osakabe *et al.*
211 (2017).

212 Maeoka *et al.* (2020) demonstrated that a general $\Delta\Delta\text{Cq}$ method without restriction enzyme treatment
213 could be used for allele-frequency measurement if a specific primer set was designed to amplify only the R
214 allele at the target locus. Similar to the RED- $\Delta\Delta\text{Cq}$ method, DNA samples with unknown mixing ratios were
215 dispensed and amplified using primer sets corresponding to TG and HK loci, respectively. Unlike the RED-
216 $\Delta\Delta\text{Cq}$ method, the control sample was not taken from the test sample solution, but was prepared as a DNA

217 solution containing 100% R, hereafter denoted as U (= pUre R line). Then, X^{HU} and X^{TU} denote the template
218 DNA quantities ready for subsequent PCR amplification:

$$\begin{aligned} X^{\text{HU}} &= X'_{\text{R}}, \\ X^{\text{TU}} &= \delta_{\text{T}} X'_{\text{R}}. \end{aligned}$$

220 Eq. 11

221 Though the definition of δ_{T} is the same as Eq. 7, the quantity is denoted by X'_{R} instead of $X_{\text{S}} + X_{\text{R}}$ because it
222 no longer originates from the R portion of the test sample itself (i.e., not internal).

223 For the test sample (denoted as V), the template DNA quantities amplified at the housekeeping (X^{HV}) and
224 target (X^{TV}) loci are expressed as follows:

$$\begin{aligned} X^{\text{HV}} &= X_{\text{S}} + X_{\text{R}}, \\ X^{\text{TV}} &= \delta_{\text{T}}(zX_{\text{S}} + X_{\text{R}}). \end{aligned}$$

226 Eq. 12

227 In the PCR process of the modified $\Delta\Delta\text{Cq}$ method, the small positive number z provides the template
228 quantity of S, which is non-specifically amplified even with the specific primer set, which is designed to
229 amplify only the R allele at the target locus. As the primer set for the housekeeping gene was non-specific,
230 both X^{HU} and X^{HV} were fully amplified. Assuming that all four template DNAs are amplified with efficiency
231 $1 + \eta$, we define the two ΔCq values as $\Delta\tau^{\text{U}} = \tau^{\text{TU}} - \tau^{\text{HU}}$ and $\Delta\tau^{\text{V}} = \tau^{\text{TV}} - \tau^{\text{HV}}$. Finally, their $\Delta\Delta\text{Cq}$ values
232 are $\Delta\Delta\tau = \Delta\tau^{\text{V}} - \Delta\tau^{\text{U}}$, which yields a formula identical to Eq. 10.

233 ***Interval estimation of allele frequency and experimental parameters based on*** 234 ***qPCR over multiple bulk samples***

235 Finally, we consider the likelihood model to obtain the interval estimate of the allele frequency based on the
236 (RED-) $\Delta\Delta\text{Cq}$ analysis over multiple bulk samples. Assume that the population has the R allele at the
237 frequency p from which N bulk samples are taken. The h th sample ($h = 1, 2, 3, \dots, N$) consists of n_h haploid
238 individuals, of which m_h are resistant mutants. As shown in Eq. 9, each Cq value is determined not only by
239 the DNA quantities, which are denoted as $X_{h,\text{R}}$ and $X_{h,\text{S}}$ for each sample, but also by parameters such as δ_{T} or
240 σ_c^2 accompanying the experimental operation. We can simultaneously estimate these if we have multiple
241 bulk samples, for which the likelihood function of obtaining the Cq values under the parameters is defined.

242 Although it was possible to define a joint likelihood for each Cq quartet, or we could define the likelihood
243 of a single $\Delta\Delta\text{Cq}$ value, we propose the joint likelihood for the two ΔCq values, $\Delta\tau_h^{\text{W}} = \tau_h^{\text{TW}} - \tau_h^{\text{HW}}$ and
244 $\Delta\tau_h^{\text{D}} = \tau_h^{\text{TD}} - \tau_h^{\text{HD}}$, for the convenience of numerical calculation:

$$\begin{aligned} \Delta\tau_h^{\text{W}} &\sim \text{N}\left(-\frac{\ln\delta_{\text{T}}}{\ln(1+\eta)}, 2\sigma_c^2\right), \\ \Delta\tau_h^{\text{D}} &\sim \text{N}\left(-\frac{\ln\delta_{\text{T}} + \ln\left(\frac{zX_{h,\text{S}} + X_{h,\text{R}}}{X_{h,\text{S}} + X_{h,\text{R}}}\right)}{\ln(1+\eta)}, 2\sigma_c^2\right). \end{aligned}$$

247 Eq. 13

248 Although Eq. 13 is defined for the RED- $\Delta\Delta\text{Cq}$ method, it is also applicable to the $\Delta\Delta\text{Cq}$ method by Maeoka
249 *et al.* (2020) by substituting $\Delta\tau_h^{\text{W}}$ and $\Delta\tau_h^{\text{D}}$ to $\Delta\tau_h^{\text{U}} = \tau_h^{\text{TU}} - \tau_h^{\text{HU}}$ and $\Delta\tau_h^{\text{V}} = \tau_h^{\text{TV}} - \tau_h^{\text{HV}}$, respectively.

250 *Formulation of likelihood based on gamma or beta distribution*

251 Using the relationship between m_h , n_h , and p in Eq. 2, we proceed to the likelihood function defined as the
 252 probability of observing the set of $\Delta\tau_h^W$ and $\Delta\tau_h^D$ under the given values of p , n_h , and other experimental
 253 parameters. In Eq. 13, $\Delta\tau_h^W$ is not affected by the R : S ratio in the bulk sample; it is only affected by the
 254 experimental parameters, δ_T , η , and σ_c^2 . In addition, by taking the differences, there is no need to estimate as
 255 X_θ and δ_B appear in Eq. 9.

256 Conversely, we must consider the amount of DNA in the bulk sample to calculate the probability of
 257 obtaining $\Delta\tau_h^D$. When the size of m_h is specified under the binomial assumption, the quantities of DNA in the
 258 h th bulk sample, $X_{h,R|m_h}$ and $X_{h,S|m_h}$, can independently take any positive values following the gamma
 259 distribution of Eq. 3, and their proportions $Y_{h,R|m_h} = X_{h,R|m_h} / (X_{h,R|m_h} + X_{h,S|m_h})$ are Beta($m_h k$, ($n_h -$
 260 m_h) k) as shown in Eq. 4. If the sample contains only S or R, then $X_{h,R|m_h=0} = 0$ or $X_{h,S|m_h=n_h} = 0$ is
 261 guaranteed.

262 The likelihood function for the observed ΔC_q values on the h th bulk sample L_h is defined as follows:

$$263 \quad L_h = P(\Delta\tau_h^W | \delta_T, \eta, \sigma_c^2) \sum_{m_h=0}^{n_h} [\text{Bin}(m_h | n_h, p) P(\Delta\tau_h^D | m_h, \delta_T, z, \eta, \sigma_c^2)],$$

$$264 \quad P(\Delta\tau_h^D | m_h, \delta_T, z, \eta, \sigma_c^2) = \begin{cases} N\left(-\frac{\ln(z\delta_T)}{\ln(1+\eta)}, 2\sigma_c^2\right) & (m_h = 0) \\ \psi_G \text{ or } \psi_B & (m_h = 1, 2, \dots, n_h - 1) \\ N\left(-\frac{\ln\delta_T}{\ln(1+\eta)}, 2\sigma_c^2\right) & (m_h = n_h) \end{cases}$$

265 Eq. 14

266 We must consider not only the possible cases of m_h , but also the entire range of the DNA amounts. If we use
 267 the gamma distributions, for every case $m_h = 1, 2, \dots, n_h - 1$, we need to calculate the double integration for
 268 ψ_G , the probability of obtaining $\Delta\tau_h^D$ under the whole region of $X_{h,R|m_h} = r$ and $X_{h,S|m_h} = s$ for the interval
 269 $\{D: 0 \leq r < \infty, 0 \leq s < \infty\}$.

$$270 \quad \psi_G = \iint_D N\left(-\frac{\ln\delta_T + \ln\left(\frac{zs+r}{s+r}\right)}{\ln(1+\eta)}, 2\sigma_c^2\right) \text{Ga}(r | m_h k, \theta) \text{Ga}(s | (n_h - m_h)k, \theta) dr ds.$$

271 Eq. 15

272 The common scale parameter of the gamma distributions, θ , is not identifiable from the data, although we
 273 can substitute arbitrary values $\theta = 1$ for it because it is canceled in $\Delta\tau_h^D$ as a quotient.

274 Since the computational burden for the double integration is large, we simplified the likelihood model
 275 with the beta distribution. As shown in Eq. 4, the proportion $Y_{h,R} = X_{h,R} / (X_{h,R} + X_{h,S})$ is as follows:
 276 Beta($m_h k$, ($n_h - m_h$) k). Then, the probability of obtaining $\Delta\tau_h^D$ is replaced with, ψ_B defined as follows:

$$277 \quad \psi_B = \int_0^1 N\left(-\frac{\ln\delta_T + \ln(z + y(1-z))}{\ln(1+\eta)}, 2\sigma_c^2\right) \text{Beta}(y | m_h k, (n_h - m_h)k) dy.$$

278 Eq. 16

279 We provide an R function “freqpcr()” to estimate the parameters p , k , δ_T , and σ_c simultaneously when the
280 set of Cq measurements (τ_h^{HW} , τ_h^{TW} , τ_h^{HD} , and τ_h^{TD}) and n_h are given for each of the N bulk samples. The
281 package source is available at <https://github.com/sudoms/freqpcr>. The default is “beta = TRUE,” where the
282 beta distribution model of Eq. 16 was used instead of gamma. Regardless of the algorithms, the asymptotic
283 confidence intervals are calculated using the inverse of the Hessian matrix evaluated at the last iteration. The
284 functions `nlm()` of R and `cubintegrate()` in the R package “cubature” (Narasimhan et al. 2019) are used for
285 the iterative optimization and the double integration, respectively.

286 ***Identification of auxiliary parameters using DNA samples with known allele-*** 287 ***mixing ratios***

288 The likelihood introduced above ensures that we can estimate the sizes of p and k together with other
289 experimental parameters, δ_T and σ_c , if we have conducted a (RED-) $\Delta\Delta$ Cq analysis on multiple bulk samples.
290 However, the size of z is not identified and must be specified as a fixed parameter. The amplification
291 efficiency, η , is estimated in theory over the iterative calculation of Eq. 13, but in fact, simultaneous
292 estimation sometimes fails when η is set as unknown.

293 Therefore, the experimenter should identify the sizes of these auxiliary parameters. To estimate their
294 plausible sizes, one can conduct (RED-) $\Delta\Delta$ Cq analysis using DNA solutions with known allele ratios; for
295 instance, DNA can be extracted from each of the pure breeding lines of S and R and mix the solutions at
296 multiple ratios, or make a dilution series of R by S. As the ratio of X_R to X_S is strictly fixed, Eq. 9 is directly
297 applicable to express the relationship between DNA quantities and the four Cq measurements. The R
298 functions `knownqpcr()` and `knownqpcr_unpaired()` appearing in the package provide the maximum
299 likelihood estimation for δ_B , δ_T , σ_c , z , and η . These values can be used as fixed parameters in the `freqpcr()`
300 function.

301 Another objective of the analysis with known-ratio samples is to test the homoscedasticity of the qPCR
302 data at the scale of Cq measures. Regarding the relationship between the R allele frequency and the
303 corresponding $2^{-\Delta\Delta Cq}$ measures (the approximate point estimate of the frequency), Osakabe *et al.* (2017)
304 demonstrated linearity using a sample series of *T. urticae* DNA with multiple mixing ratios on CHS1
305 (I1017F). In the next section, we recycled the same data to compare whether the Cq measurements in the
306 RED- $\Delta\Delta$ Cq analysis obey the homoscedasticity in the scale of $\Delta\Delta$ Cq or $(1 + \eta)^{-\Delta Cq}$.

307 **Materials and laboratory methods**

308 ***Experiment 1: estimation of auxiliary parameters and verification of*** 309 ***homoscedasticity in Cq measurements based on mite DNA samples with known*** 310 ***allele-mixing ratios***

311 *Experimental setup*

312 In the experiment by Osakabe *et al.* (2017), the resistant mite strain (SoOm1-etoR strain) originated from a
313 field population collected in Omaezaki City, Shizuoka, Japan (34.7°N, 138.1°E) in January 2012. The
314 susceptible strain was obtained from Kyoyu Agri Co., Ltd. (Kanagawa, Japan) (Kyoyu-S strain). For each
315 strain, two pairs of females and males were used separately. Each pair was allowed to mate and oviposit on a
316 kidney bean leaf square (2 × 2 cm) for four days. The mites were then confirmed to be homozygous on the

317 CHS1 locus using sequence analysis. Genomic DNA extracted from the offspring of each pair was used for
318 qPCR analysis. For each pair, the DNA extracts were prepared twice, each of which was a mixture from 50
319 adult females homogenized together, that is, four extracts (replicates) for each strain.

320 To verify the validity of the RED- $\Delta\Delta Cq$ method, qPCR analysis was performed with heterogeneous DNA
321 solutions with 10 mixing ratios of $X_R/(X_R + X_S) = \{0, 0.001, 0.005, 0.01, 0.05, 0.1, 0.25, 0.5, 0.75, 1\}$. The
322 net DNA concentration of each mixed solution was adjusted to $1 \text{ ng } \mu\text{l}^{-1}$, from which 15 ng was dispensed
323 into each of the two tubes. Only one was digested with the restriction enzymes before qPCR. For digestion,
324 the samples were treated with a mixture of two enzymes, *MluC* I (10 units) and *Taq*^{qi}I (20 units; New
325 England BioLabs, Ipswich, MA, USA), at 37 °C for 3 h, followed by incubation at 65 °C for 3 h. This is due
326 to the polymorphism of the CHS1 loci; the 1017 codon of *T. urticae* displays ATT (Kyoyu-S strain) or TTT
327 (SoOm1-etoR) sequences, whereas the upstream 1016 codon displays a synonymous TCG or TCA
328 independent of the strains (Van Leeuwen et al. 2012). Therefore, we need to digest both TCGATT (underline
329 shows the restriction site of *Taq*^{qi}I) and TCAATT (*MluC* I) to diminish the entire S allele.

330 qPCR analysis using the intercalator method was performed using the LightCycler Nano System (Roche
331 Diagnostics, Basel, Switzerland) with SYBR Fast qPCR Mix (Takara, Kusatsu, Japan) as described
332 previously (Osakabe et al. 2017). The primer sets were tu03*CHS1* (forward: 5'-
333 GGCCTGCTTCATCCACAAG-3' and reverse: 5'-GTGTTCCCAAGTAACAACGTTTC-3') and
334 tu25*GAPDH* (forward: 5'-GCACCAAGTGCTAAAGCATGGAG-3' and reverse: 5'-
335 GAACTGGAACACGGAAAGCCATAC-3').

336 *Statistical analysis*

337 The maximum likelihood estimation of δ_B , δ_T , σ_c , z , and η was conducted with the “knownqpcr_unpaired”
338 function of the freqpcr package. The raw Cq data are available as ESM 1 along with a step-by-step guide for
339 statistical analyses (ESM 2). Due to the limitation of the handling capacity of the thermal cycler, qPCR
340 analysis was not conducted on undigested samples of the nine mixing ratios other than $X_R/(X_R + X_S) = 1$
341 (i.e., pure R solution). Thus, in each replicate, Osakabe *et al.* (2017) used the observed $\Delta\tau^W$ value when the
342 ratio = 1 for other ratios to calculate the conventional $\Delta\Delta Cq$ indices. As we have shown in Eq. 9, this
343 operation does not affect the point estimates of p , although the size of the Cq measurement error (σ_c) will be
344 underestimated if we recycle the observed Cq value multiple times. The “knownqpcr_unpaired” function was
345 developed to deal with such incomplete data (i.e., the observations of τ^{HW} , τ^{TW} , τ^{HD} , and τ^{TD} have different
346 data lengths). If the four Cq measurements are available for all samples, then “knownqpcr” can be used.

347 Regarding the relationship between the true mixing ratio and the RED- $\Delta\Delta Cq$ measures in the sample, the
348 linearity was analyzed using a linear model via the function “lm” running on R version 3.6.1 (R Core Team
349 2019), where the response variables were put into the model at the scale of Cq or $(1 + \eta)^{-\Delta\Delta Cq}$. Based on the
350 linear models, we tested heteroscedasticity using the Breusch-Pagan test via the bptest() function of the R
351 library “lmtest” (Hothorn et al. 2019).

352 ***Experiment 2: evaluation of the simultaneous estimation method with randomly 353 generated data***

354 Since the experiment by Osakabe *et al.* (2017) used a sample series with strict mixing ratios, the effect of
355 individual differences in DNA yield was not evaluated. Instead, we conducted a numerical experiment to
356 verify the accuracy of the simultaneous parameter estimation under uncertainty in the individual DNA yield.
357 The frequency of the R allele in the population, p , was set to 0.01, 0.05, 0.1, 0.25, 0.5, or 0.75.

358 For the sampling strategy, N bulk samples (the parameter ‘ntrap’ in the R source code), each comprising
359 of n individuals (n was fixed among the samples: the parameter ‘npertrap’ in the code), were generated by
360 random sampling from a wild population of a haploid organism. To assess how the estimation interval
361 responds to the sample sizes, we evaluated the combination of $N = \{2, 4, 8, 16, 32, 64\}$ and $n = \{4, 8, 16, 32,$
362 $64\}$, though the combinations with $Nn > 128$ were excluded (Nn corresponds to ‘ntotal’ in the code). The
363 DNA quantities (X_R and X_S) contained in each bulk sample were generated as random numbers that followed
364 the gamma distributions of Eq. 3. To cover a plausible variability range of the DNA yield, the gamma shape
365 parameter was varied as $k = \{1, 3, 9, 27\}$. Depending on the size of k , the gamma scale parameter was set at
366 $\theta = 1 \times 10^{-6}/k$ to fix the mean of the individual DNA yield to 1×10^{-6} . The termination threshold for
367 qPCR X_θ was fixed at 1.

368 We fixed the other parameters due to limitations of the computing resources. From the results of
369 Experiment 1, $\delta_T = 1.2$, $\delta_B = 0.24$, $z = 0.0016$, and $\eta = 0.97$ were presupposed. As for the random errors in
370 the PCR amplification process and/or the Cq measurement, $\sigma_c = 0.2$ was assumed regardless of the initial
371 template quantity. For each of the 624 parameter regions, the dummy datasets comprising N bulk samples
372 were generated 1,000 times independently with different random number seeds (i.e., 1,000 replicates), for
373 which the parameter estimation with `freqpcr(..., beta = TRUE)` was run on the R 3.6.1 environment. The
374 simulation code is available in ESM 3.

375 As we also implemented the gamma distribution model as `freqpcr(..., beta = FALSE)`, a numerical
376 experiment with the gamma model was also conducted for the first 250 replicates, and the estimation
377 accuracy was compared between the two assumptions. Furthermore, we also fitted the function with the
378 settings `freqpcr(..., K = 1)`, that is, assuming the gamma shape parameter was fixed at 1 (a.k.a. exponential
379 distribution), in addition to the default simulation with all parameters unknown. Further, the easiest way to
380 estimate p derived from Eq. 10, we averaged the observed $\Delta\Delta Cq$ values for N bulk samples and transformed
381 them as $\hat{p} = (1 + \eta)^{-\overline{\Delta\Delta\tau}}$.

382 Results

383 *Estimation of auxiliary parameters and verification of homoscedasticity*

384 Based on the Cq measures, the auxiliary parameters were estimated based on the RED- $\Delta\Delta Cq$ analysis of the
385 I1017F mutation of *T. urticae*. As for the initial quantity of template DNA (the parameter “meanDNA” on
386 the R code; defined as X/X_θ), the maximum-likelihood estimate was 1.256×10^{-6} (95% confidence
387 interval: 7.722×10^{-7} to 2.041×10^{-6}). The relative quantity of the target gene to the housekeeping gene
388 δ_T (targetScale) was estimated to be 1.170 (95% CI: 1.069–1.280). The locus-independent change rate in the
389 template quantity accompanying the restriction enzyme treatment δ_B (baseChange) was 0.2361 (95% CI:
390 0.2040 to 0.2731). The measurement error in the scale of Cq σ_c (SD) was 0.2376 (95% CI: 0.2050 to
391 0.2755). The residue rate of the S allele after digestion z (zeroAmount) was 0.001564 (95% CI: 0.001197–
392 0.002044). The efficiency of amplification per PCR cycle η (EPCR) was 0.9712 (95% CI: 0.9231–1.022).

393 In the RED- $\Delta\Delta Cq$ analysis of the etoxazole resistance of *T. urticae*, the relationship between the true R
394 allele frequency ($Y_R = X_R/(X_R + X_S)$ in the sample) and the corresponding Cq measures exhibited higher
395 homoscedasticity in the scale of the measured $\Delta\Delta Cq$ values rather than in $(1 + \eta)^{-\Delta\Delta Cq}$, the transformation to
396 \hat{Y}_R (Fig. 1). The linear regression of the $\Delta\Delta Cq$ values on $-\ln[0.001564 \times (1 - Y_R) + Y_R]/\ln(1 + 0.971)$
397 showed high linearity (intercept = -0.07694, coefficient = 1.025, adjusted $R^2 = 0.9936$). The
398 homoscedasticity of the coefficient of determination was not rejected at the 5% level of significance

399 (Breusch-Pagan test: $BP = 3.1577$, $df = 1$, $p = 0.07557$) (Fig. 1A). Conversely, the linear regression of
400 $1.971^{-\Delta\Delta Cq}$ on $[0.001564 \times (1 - Y_R) + Y_R]$ showed a slightly lower linearity (intercept = -0.008625 ,
401 coefficient = 1.092 , adjusted $R^2 = 0.9709$). The Breusch-Pagan test was highly significant ($BP = 13.978$, $df =$
402 1 , $p = 0.0001849$), rejecting homoscedasticity (Fig. 1B). These results suggest that it is easier to model the
403 error structure of the RED- $\Delta\Delta Cq$ method on the scale of Cq values (logarithm) rather than frequency (linear
404 scale).

405 ***Evaluation of the simultaneous estimation method with randomly generated data***

406 Among the 624 parameter regions of the numerical simulation with 1,000 replicates (250 for the gamma
407 model), the total success rate of the interval estimation p using `freqpqr(..., beta = TRUE)` was 70.6% and
408 94.5% when all parameters were unknown, and when the gamma shape parameter was fixed as $k = 1$,
409 respectively. The “success rate” here indicates the probability when the function returns certain values other
410 than NA (i.e., the diagonal of the Hessian was not negative): no guarantee that the estimated confidence
411 interval was accurate. The estimation success for the Cq measurement error, σ_c , was 69.6% and 97.6% in the
412 beta-distribution model with unknown k and $k = 1$, respectively. The relative quantity of the target gene, δ_T ,
413 was 68.1% and 96.1%, respectively. However, the estimation success of k was 59.9% with the beta
414 distribution model, showing a lower performance than the other parameters. This result implies that the
415 likelihood is insensitive to the size of k . Conversely, the estimation of p is robust to the size of k , as we show
416 later in this section.

417 The estimation success of `freqpqr()` largely depended on the total sample size (Nn corresponding to the
418 facet ‘ntotal’ in the figures), as well as the level of p (Figure S1 and S2 for the beta and gamma models, with
419 all parameters unknown). In each parameter region, the quantity $\text{Bin}(0|Nn, p)$ generally gives the probability
420 that the whole sample contains no R individuals. When Nn is larger enough, $Nn > 3/p$ is approximately the
421 requirement for the total sample size to contain at least one R individual with 95% confidence, called the
422 “rule of three” (Eypasch et al. 1995). The gray backgrounds in the facets of Figures 2–4 and S1–S6 signify
423 the regions where the total sample sizes are smaller than the thresholds (e.g., 60 individuals are required
424 when $p = 0.05$). As shown in Figure S1, the parameter estimation often failed when Nn did not meet the
425 rule of three. Once we exclude the parameter regions of $Nn \leq 3/p$, the estimation success rate of p with
426 `freqpqr(..., beta = TRUE)` improved to 84.3% and 99.9% with all parameters unknown and assuming $k = 1$,
427 respectively.

428 As for the estimation accuracy of p , the `freqpqr()` function assuming beta distribution provides an unbiased
429 estimator. Figures 2 and 3 show the estimated sizes of p using the beta model with all parameters unknown
430 and assuming $k = 1$, respectively. Both settings demonstrated that the estimator converged to the true R
431 frequency; the upper/lower bounds of the estimated 95% confidence intervals (yellow/blue boxes in each
432 plot) became narrower as we increased the total sample sizes (Nn) or included more bulk DNA samples (N).
433 Fixing the size of the gamma shape parameter to $k = 1$ scarcely affected the point estimates and intervals of
434 p , as long as $Nn > 3/p$ is satisfied (Figure 3). However, if every individual was analyzed separately, the
435 interval estimation was only possible when k was fixed (see the regions of “sample division = ntotal” cases
436 in Figure 2).

437 When we used the gamma distribution model, the interval estimation of p was also possible and unbiased
438 (Figure S3). However, when we defined the point estimator of p as a simple average, that is, $\hat{p} =$
439 $(1 + \eta)^{-\overline{\Delta\Delta Cq}}$, it was strongly underestimated as the samples were more divided (N/Nn was large)
440 (Figure 4). The upper bound of 95% CI often violated 1, suggesting that the “simple average of $\Delta\Delta Cq$ ” \pm
441 1.96 SE is inadequate for the interval estimation based on the RED- $\Delta\Delta Cq$ method.

442 Although the `freqpcr()` function with the gamma and beta distributions both showed an unbiased
443 estimation of p , the gamma model was disadvantageous regarding calculation time and the number of
444 iterations before convergence. The time varied largely in the model settings and sample sizes (Figures S4–
445 S6). Amongst the settings we tried, beta model with fixed k was the fastest; it converged within a few
446 seconds in most parameter regions (median and 75 percentile: 0.32 and 0.69 s: Figure S6). It was three and
447 >10 times faster than the beta (0.91 and 2.4 seconds: Figure S4) and gamma (3.0 and 15 s: Figure S5) model,
448 respectively with all parameters unknown. The calculation time increased as the dataset size increased - Nn
449 and the sample was more divided (larger N/Nn) in the beta distribution model, because the marginal
450 likelihood was calculated for each bulk sample. Conversely, the gamma distribution model (Figure S5)
451 requires increased calculation time as the size of each bulk sample becomes larger (larger n_h). This was
452 considered because the combination of $\text{Bin}(m_h | n_h, p)$ exploded when n_h was large.

453 Regarding the estimation accuracy of the shape parameter, k , it was underestimated as the real size of the
454 parameter increased (e.g., $k = 27$) when we used the gamma distribution model (Figure S7B). Since the
455 iterative fitting of the parameter in `freqpcr()` always starts internally from $k = 1$ (this was determined due to
456 the calculation stability), this bias suggests the likelihood function of ψ_G (Eq. 15), with little information on
457 the size of k compared with p . Then, k tends to stay at its initial value, suggesting that the gamma model is
458 not suitable for the simultaneous estimation of p and k . Unlike the gamma version, the fitting of k with
459 `freqpcr(beta = TRUE)` was satisfactory when we divided the total samples into more bulk samples (larger
460 N/Nn), although the initial value dependence was still observed, especially when p was small (Figure S7A).
461 This may be because the estimation of k via $\text{Beta}(m_h k, (n_h - m_h)k)$ in Eq. 16 is comparable with
462 measuring the overdispersion of $Y_{h,R|m_h}$, which is only possible when multiple bulk samples contain both R
463 and S alleles.

464 Discussion

465 In the present study, we developed a statistical model to estimate the population allele frequency based on
466 qPCR across multiple bulk samples. There have been problems with the conventional point estimator of the
467 allele frequency by averaging the observed $\Delta\Delta\text{Cq}$ values $\hat{p} = (1 + \eta)^{-\overline{\Delta\Delta\tau}}$. It sometimes exceeds 1
468 when the frequency of the target allele is close to 1. Furthermore, when one tries to quantify the mutant allele
469 rare in the population, most bulk samples contains only the wild type. The conventional \hat{p} is vulnerable to
470 many zero samples, which makes the frequency estimation more difficult when p is small. To circumvent
471 these problems, our interval estimation explicitly models the number of individuals contained in each bulk
472 sample (the binomial assumption) as well as the individual DNA yields (the gamma assumption), thereby
473 obtaining the interval estimate over the entire range $0 < p < 1$.

474 The explicit modeling of individuals also allows sample division to various degrees, which helps us to
475 balance our sampling strategy on the cost-precision tradeoff. We can achieve higher precision (narrower
476 confidence interval) by increasing the total sample size, $\sum_{h=1}^N n_h$ although it also increases the costs
477 associated with sample collection and laboratory work, including library preparation and PCR analysis.
478 Recent advances in molecular diagnosis have relieved sampling costs. We can now extract DNA from dead
479 insect bodies obtained from sticky traps (Uesugi et al. 2016). Nevertheless, a larger sample size still imposes
480 a larger handling cost if we analyze the collected individuals individually via non-quantitative PCR.

481 The combination of mass trapping and bulk qPCR analysis solves the latter by collecting more individuals
482 and pooling them. This can result in higher precision with less work than individual PCR. For instance, we
483 sampled 16 individuals from the population with an allele frequency of $p = 0.05$ and analyzed two

484 individuals at once in the numerical experiment (Figure 2: facet of $n_{\text{total}} = 16$, sample division = 8). The
485 lower and upper bounds of the 95% confidence interval p were estimated to be 0.0087 and 0.34, respectively,
486 using `freqpqr(..., beta = TRUE)` (as the medians of the 1,000 independent trials). We also simulated the case
487 of $n_{\text{total}} = 64$ and sample division = 4 (i.e., analyzed 16 individuals together). The upper and lower bounds
488 were 0.015 and 0.15, respectively. Thus, we improved the precision of the interval estimate with half the
489 handling effort.

490 Also in non-quantitative PCR, sample pooling has been considered as a tool to detect (c)DNA rare in the
491 population with practical labor, sometimes as a high throughput pre-screening of a number of samples e.g. in
492 clinical examination (Taylor et al. 2010; Yelin et al. 2020). In some fields, such as plant quarantine, it is
493 important to guarantee that a product is not contaminated with pests or unapproved GM seeds at a certain
494 consumer risk. As the assumed frequency range is extremely low ($p \approx 0.001$), frequency estimation is not
495 realistic (3,000 seeds are needed to meet the “rule of three” when $p = 0.001$) and is not required for the
496 current inspection routine. Thus, group testing based on non-quantitative PCR has been conducted in these
497 fields (Yamamura et al. 2019). Yamamura and Hino (2007) proposed a procedure to estimate the upper limit
498 of the population allele frequency, in which they used the proportion of bulk samples detected as “positive.”

499 Overall, there has been a gap in methodology between the frequency estimation based on the individual
500 PCR and the non- or semi-quantitative PCR based on the non-quantitative bulk PCR. Although it provides
501 the highest estimation precision following binomial distribution, the former is only available at a higher p ; it
502 becomes labor-intensive once we try to quantify rare alleles. The latter can be applied to a lower range of p ,
503 but the precision is generally low or even non-quantitative. Bridging the gap, our qPCR-based procedure
504 offers an allele-frequency estimation in the mid-low range ($p = 0.01$ to 0.25), which is considered a critical
505 range for decision making in some fields like pesticide resistance management (Takahashi et al. 2017; Sudo
506 et al. 2018).

507 Although this study focused on resistance genes, the likelihood model in Eq. 13 can be used for other
508 qPCR protocols based on this $\Delta\Delta C_q$ method. If both the specific and nonspecific primer sets are available to
509 amplify the mutant and “wild type + mutant” alleles at the target locus, they can be used for the test and
510 control samples equivalent to X^{TV} in Eq. 12 and X^{TU} in Eq. 11, respectively. However, there is a caveat in
511 determining which allele should be amplified with a specific primer set and which affects the estimation
512 accuracy due to the intrinsic nature of $(1 + \eta)^{-\Delta\Delta\tau}$. As shown, the 95% confidence intervals were broader
513 when $p = 0.75$ than when $p = 0.25$ (Figure 2), the accuracy was not symmetric around 0.5, but more accurate
514 when the frequency was low. That is, one should design a specific primer set to amplify the allele that would
515 be rare in the population to improve the signal-to-noise ratio.

516 The maximum likelihood estimation with `freqpqr()` relies on the assumption that the quantities of the S and
517 R alleles in each bulk sample independently follow gamma distribution and that their quotient is expressed
518 using beta distribution. Fixing the size of the gamma shape parameter k further accelerated the optimization,
519 which was owing to the robustness of p to the size of k . However, once the size of k was fixed much larger
520 than the actual size of the gamma shape parameter (i.e., the individual DNA yield was regarded as almost a
521 fixed value), the iterative optimization using the `nlm()` function sometimes returned an error. Therefore, one
522 should start with a smaller shape parameter e.g., $k = 1$ (the exponential distribution: Figure 3), which is
523 currently the default setting of the `freqpqr` package.

524 In qPCR applications for diagnostic use, $\Delta\Delta C_q$ is often used with calibration. One of the popular methods
525 is the involvement of technical replicates; each sample is dispensed and analyzed using qPCR multiple times,
526 which cancels the C_q measurement error. The measurement error obeys a homoscedastic normal distribution
527 in the C_q scale, as shown in Experiment 1. Thus, a simple solution is to average the C_q values measured for

528 every bulk sample before the estimation with `freqpcr()`, although the estimated size of σ_c changes from its
529 original definition in Eq. 9. However, it is trivial if the number of technical replicates is unified between bulk
530 samples.

531 Moreover, the comparison of Cq values is sometimes conducted on more than one internal reference
532 because there is no guarantee that the expression level of a “housekeeping gene” is always constant
533 (Vandesompele et al. 2002). Future updates of `freqpcr()` will handle multiple internal references. As long as
534 qPCR is used to estimate population allele frequency, the use of statistical inferences on the bulk samples, as
535 presented in this study, will continue to be a realistic option for regional allele monitoring and screening for
536 practitioners, such as those in agricultural, food security, and public health sectors.

537 **Acknowledgements**

538 We appreciate Dr. Kohji Yamamura and Dr. Takehiko Yamanaka for earlier discussion on the gamma
539 assumption of the individual DNA yield. The work was supported by a grant from the Ministry of
540 Agriculture, Forestry, and Fisheries of Japan (Genomics-based Technology for Agricultural Improvement):
541 PRM05 to MO and PRM07 to MS.

542 **Conflict of Interest**

543 None declared.

544 **Data Accessibility**

545 The R package source is available at <https://github.com/sudoms/freqpcr>. The output data of the numerical
546 experiment are available at <https://figshare.com/collections/freqpcr/5258027>. The source code for the figures,
547 including the mite dataset from Osakabe *et al.* (2017), are available as electronic supplementary materials.

548

549 ESM 1

550 Verification of the RED- $\Delta\Delta$ Cq method: raw dataset used by Osakabe et al. (2017).

551 ESM 2

552 R source code for Experiment 1 (Figure 1), including a brief guide to the “freqpcr” package.

553 ESM 3

554 R source code for the numerical simulation of `freqpcr()` (Experiment 2), and the codes for Figures 2 and
555 after.

556

557 **References**

558 Andow, D. A., and D. N. Alstad. 1998. F2 screen for rare resistance alleles. *Journal of Economic*
559 *Entomology* 91:572–578.

- 560 Bustin, S. A., V. Benes, J. A. Garson, J. Hellemans, J. Huggett, M. Kubista, R. Mueller, et al. 2009. The
561 MIQE guidelines: minimum information for publication of quantitative real-time PCR experiments.
562 *Clinical chemistry* 55:611–622.
- 563 Chen, H., M. Rangasamy, S. Y. Tan, H. Wang, and B. D. Siegfried. 2010. Evaluation of five methods for
564 total DNA extraction from western corn rootworm beetles. *PLoS one* 5:e11963.
- 565 Donnelly, M. J., A. T. Isaacs, and D. Weetman. 2016. Identification, validation, and application of molecular
566 diagnostics for insecticide resistance in malaria vectors. *Trends in parasitology* 32:197–206.
- 567 Eypasch, E., R. Lefering, C. K. Kum, and H. Troidl. 1995. Probability of adverse events that have not yet
568 occurred: a statistical reminder. *Bmj* 311:619–620.
- 569 Falconer, D. S. 1960. *Introduction to quantitative genetics*. Oliver And Boyd; Edinburgh; London.
- 570 French-Constant, R. H. 2013. The molecular genetics of insecticide resistance. *Genetics* 194:807–815.
- 571 Germer, S., M. J. Holland, and R. Higuchi. 2000. High-throughput SNP allele-frequency determination in
572 pooled DNA samples by kinetic PCR. *Genome research* 10:258–266.
- 573 Gould, F., A. Anderson, A. Jones, D. Sumerford, D. G. Heckel, J. Lopez, S. Micinski, et al. 1997. Initial
574 frequency of alleles for resistance to *Bacillus thuringiensis* toxins in field populations of *Heliothis*
575 *virescens*. *Proceedings of the National Academy of Sciences* 94:3519–3523.
- 576 Gouws, E. J., K. J. Gaston, and S. L. Chown. 2011. Intraspecific body size frequency distributions of insects.
577 *PLoS One* 6:e16606.
- 578 Hothorn, T., A. Zeileis, R. W. Farebrother, C. Cummins, G. Millo, and D. Mitchell. 2019. *lmtree: Testing*
579 *Linear Regression Models*. <https://cran.r-project.org/package=lmtree>
- 580 Kim, S. Y., K. E. Lohmueller, A. Albrechtsen, Y. Li, T. Korneliusen, G. Tian, N. Grarup, et al. 2011.
581 Estimation of allele frequency and association mapping using next-generation sequencing data. *BMC*
582 *bioinformatics* 12:231.
- 583 Knapp, M. 2016. Relative importance of sex, pre-starvation body mass and structural body size in the
584 determination of exceptional starvation resistance of *Anchomenus dorsalis* (Coleoptera: Carabidae). *PLoS*
585 *one* 11:e0151459.
- 586 Kundu, D., and A. Manglick. 2005. Discriminating between the log-normal and gamma distributions. *Journal*
587 *of the Applied Statistical Sciences* 14:175–187.
- 588 Li, G., D. Reising, J. Miao, F. Gould, F. Huang, and H. Feng. 2016. Frequency of Cry1F non-recessive
589 resistance alleles in North Carolina field populations of *Spodoptera frugiperda* (Lepidoptera: Noctuidae).
590 *PLoS one* 11:e0154492.
- 591 Livak, K. J. 1997. Comparative Ct method. Pages 11–15 *in* ABI Prism 7700 sequence detection system User
592 *Bulletin #2*, P/N 4303859 (Vol. 2). Applied Biosystems.
- 593 Livak, K. J., and T. D. Schmittgen. 2001. Analysis of relative gene expression data using real-time
594 quantitative PCR and the 2- $\Delta\Delta$ CT method. *methods* 25:402–408.
- 595 Maeoka, A., L. Yuan, Y. Itoh, C. Saito, M. Doi, T. Imamura, T. Yamaguchi, et al. 2020. Diagnostic
596 prediction of acaricide resistance gene frequency using quantitative real-time PCR with resistance allele-
597 specific primers in the two-spotted spider mite *Tetranychus urticae* population (Acari: Tetranychidae).
598 *Applied Entomology and Zoology* 55:329–335.

- 599 May, R. M. 1976. Models for single populations. Pages 5–29 in *Theoretical ecology*. Blackwell, Oxford.
- 600 Narasimhan, B., M. Koller, S. G. Johnson, T. Hahn, A. Bouvier, K. Kiêu, and S. Gaure. 2019. cubature:
601 Adaptive Multivariate Integration over Hypercubes. <https://cran.r-project.org/package=cubature>
- 602 Osakabe, M., T. Imamura, R. Nakano, S. Kamikawa, M. Tadatsu, Y. Kunimoto, and M. Doi. 2017.
603 Combination of restriction endonuclease digestion with the $\Delta\Delta C_t$ method in real-time PCR to monitor
604 etoxazole resistance allele frequency in the two-spotted spider mite. *Pesticide biochemistry and*
605 *physiology* 139:1–8.
- 606 Pfaffl, M. W. 2012. Quantification strategies in real-time polymerase chain reaction. *Quantitative real-time*
607 *PCR. Appl Microbiol* 53–62.
- 608 R Core Team. 2019. R version 3.6.1.
- 609 Rakovski, C., D. J. Weisenberger, P. Marjoram, P. W. Laird, and K. D. Siegmund. 2011. Modeling
610 measurement error in tumor characterization studies. *BMC bioinformatics* 12:284.
- 611 Samayoa, L. F., R. A. Malvar, B. A. Olukolu, J. B. Holland, and A. Butrón. 2015. Genome-wide association
612 study reveals a set of genes associated with resistance to the Mediterranean corn borer (*Sesamia*
613 *nonagrioides* L.) in a maize diversity panel. *BMC plant biology* 15:35.
- 614 Schefe, J. H., K. E. Lehmann, I. R. Buschmann, T. Unger, and H. Funke-Kaiser. 2006. Quantitative real-time
615 RT-PCR data analysis: current concepts and the novel “gene expression’s C T difference” formula.
616 *Journal of molecular medicine* 84:901–910.
- 617 Snoeck, S., A. H. Kurlovs, S. Bajda, R. Feyereisen, R. Greenhalgh, E. Villacis-Perez, O. Kosterlitz, et al.
618 2019. High-resolution QTL mapping in *Tetranychus urticae* reveals acaricide-specific responses and
619 common target-site resistance after selection by different METI-I acaricides. *Insect biochemistry and*
620 *molecular biology* 110:19–33.
- 621 Sonoda, S., K. Inukai, S. Kitabayashi, S. Kuwazaki, and A. Jouraku. 2017. Molecular evaluation of diamide
622 resistance in diamondback moth (Lepidoptera: Yponomeutidae) populations using quantitative
623 sequencing. *Applied entomology and zoology* 52:353–357.
- 624 Sudo, M., D. Takahashi, D. A. Andow, Y. Suzuki, and T. Yamanaka. 2018. Optimal management strategy of
625 insecticide resistance under various insect life histories: Heterogeneous timing of selection and interpatch
626 dispersal. *Evolutionary applications* 11:271–283.
- 627 Sudo, M., K. Yamamura, S. Sonoda, and T. Yamanaka. in press. Estimating the proportion of resistance
628 alleles from bulk Sanger sequencing, circumventing the variability of individual DNA. *Journal of*
629 *Pesticide Science*.
- 630 Sugimoto, N., A. Takahashi, R. Ihara, Y. Itoh, A. Jouraku, T. Van Leeuwen, and M. Osakabe. 2020. QTL
631 mapping using microsatellite linkage reveals target-site mutations associated with high levels of resistance
632 against three mitochondrial complex II inhibitors in *Tetranychus urticae*. *Insect Biochemistry and*
633 *Molecular Biology* 103410.
- 634 Tabashnik, B. E., A. L. Patin, T. J. Dennehy, Y.-B. Liu, Y. Carriere, M. A. Sims, and L. Antilla. 2000.
635 Frequency of resistance to *Bacillus thuringiensis* in field populations of pink bollworm. *Proceedings of the*
636 *National Academy of Sciences* 97:12980–12984.
- 637 Takahashi, D., T. Yamanaka, M. Sudo, and D. A. Andow. 2017. Is a larger refuge always better? Dispersal
638 and dose in pesticide resistance evolution. *Evolution* 71:1494–1503.

- 639 Taylor, S. M., J. J. Juliano, P. A. Trottman, J. B. Griffin, S. H. Landis, P. Kitsa, A. K. Tshefu, et al. 2010.
640 High-throughput pooling and real-time PCR-based strategy for malaria detection. *Journal of clinical*
641 *microbiology* 48:512–519.
- 642 Toda, S., K. Hirata, A. Yamamoto, and A. Matsuura. 2017. Molecular diagnostics of the R81T mutation on
643 the D-loop region of the $\beta 1$ subunit of the nicotinic acetylcholine receptor gene conferring resistance to
644 neonicotinoids in the cotton aphid, *Aphis gossypii* (Hemiptera: Aphididae). *Applied entomology and*
645 *zoology* 52:147–151.
- 646 Uesugi, R., N. Hinomoto, and C. Goto. 2016. Estimated time frame for successful PCR analysis of
647 diamondback moths, *Plutella xylostella* (Lepidoptera: Plutellidae), collected from sticky traps in field
648 conditions. *Applied entomology and zoology* 51:505–510.
- 649 Van Leeuwen, T., J. Vontas, A. Tsagkarakou, W. Dermauw, and L. Tirry. 2010. Acaricide resistance
650 mechanisms in the two-spotted spider mite *Tetranychus urticae* and other important Acari: a review. *Insect*
651 *biochemistry and molecular biology* 40:563–572.
- 652 Van Leeuwen, T., P. Demaeght, E. J. Osborne, W. Dermauw, S. Gohlke, R. Nauen, M. Grbić, et al. 2012.
653 Population bulk segregant mapping uncovers resistance mutations and the mode of action of a chitin
654 synthesis inhibitor in arthropods. *Proceedings of the National Academy of Sciences* 109:4407–4412.
- 655 Vandesompele, J., K. De Preter, F. Pattyn, B. Poppe, N. Van Roy, A. De Paepe, and F. Speleman. 2002.
656 Accurate normalization of real-time quantitative RT-PCR data by geometric averaging of multiple internal
657 control genes. *Genome biology* 3:research0034–1.
- 658 Wiens, B. L. 1999. When log-normal and gamma models give different results: a case study. *The American*
659 *Statistician* 53:89–93.
- 660 Yamamura, K., and A. Hino. 2007. Estimation of the proportion of defective units by using group testing
661 under the existence of a threshold of detection. *Communications in Statistics—Simulation and*
662 *Computation*® 36:949–957.
- 663 Yamamura, K., J. Mano, and H. Shibaïke. 2019. Optimal definition of the limit of detection (LOD) in
664 detecting genetically modified grains from heterogeneous grain lots. *Quality Technology & Quantitative*
665 *Management* 16:36–53.
- 666 Yelin, I., N. Aharony, E. Shaer-Tamar, A. Argoetti, E. Messer, D. Berenbaum, E. Shafran, et al. 2020.
667 Evaluation of COVID-19 RT-qPCR test in multi-sample pools. *Clinical Infectious Diseases* 71: 2073–
668 2078. doi:10.1093/cid/ciaa531
669

670 Appendix

671 *Appendix A1: Case of Diploidy*

672 Although we considered sampling from haploid organisms, many insects and vertebrates are diploid. Let us
 673 consider that the population of a diploid insect species has the R allele frequency p , from which we collected
 674 n individuals. The bulk sample then consists of m_1 ($m_1 = 0, 1, \dots, n$) individuals of RR homozygotes, $n -$
 675 $m_1 - m_0$ RS heterozygotes, and m_0 ($m_0 = 0, 1, \dots, n$) SS homozygotes ($m_1 + m_0 \leq n$). The joint
 676 probability of obtaining $\{m_1, m_0\}$ obeys the trinomial distribution with probabilities p^2 and $(1 - p)^2$

$$677 \quad \text{Tri}(m_1, m_0 | n, p^2, (1 - p)^2) = \frac{n!}{m_1! m_0! (n - m_1 - m_0)!} \cdot p^{2m_1} \cdot (1 - p)^{2m_0} \cdot (2p - 2p^2)^{(n - m_1 - m_0)}.$$

678 *Eq. 17*

679 The total R allele in the bulk sample comes from two R/R sets contained in the m_1 homozygotes and a single
 680 set of R from the $n - m_1 - m_0$ heterozygotes. However, two S/S sets from m_0 homozygotes and a single S
 681 set from the $n - m_1 - m_0$ heterozygotes constitute the total S body. Note that the yields of R and S from
 682 these heterozygotes would be the same unless there is a genotype-dependent systematic error in the
 683 extraction efficiency.

684 Let us define the amount of DNA copies per genome: the random variable $X_{* \in (S,R) | \text{homo}}$ for the yield of S
 685 or R from the homozygotes, and $X_{* \in (S,R) | \text{hetero}}$ for S or R from the heterozygotes. As in the case of haploidy,
 686 X_R and X_S denote the allele contents in the bulk sample; they are the linear combinations of $X_{* | \text{homo}}$ and
 687 $X_{* | \text{hetero}}$:

$$688 \quad X_R = 2 \times X_{R | \text{homo}} + X_{R | \text{hetero}}, \quad X_S = X_{S | \text{hetero}} + 2 \times X_{S | \text{homo}},$$

$$689 \quad 2 \times X_{R | \text{homo}} \sim \text{Ga}(m_1 k, 2\theta), \quad X_{R | \text{hetero}} \sim \text{Ga}((n - m_1 - m_0)k, \theta),$$

$$X_{S | \text{hetero}} = X_{R | \text{hetero}}, \quad 2 \times X_{S | \text{homo}} \sim \text{Ga}(m_0 k, 2\theta).$$

690 *Eq. 18*

691 *Parameter estimation*

692 There are $n - i + 1$ cases from $m_0 = 0$ to $m_0 = n - i$ when the number of RR homozygotes is given by
 693 $m_1 = i$. The segregation ratio in the bulk sample has $\sum_{i=0}^n (n - i + 1)$ total combinations. For each
 694 combination of n , m_0 , and m_1 , Eq. 18 gives the probability of obtaining the ΔCq measures in Eq. 13.
 695 However, a drawback arises from the constraint of the amounts of R and S possessed by heterozygotes. The
 696 applicability of the likelihood model (Eq. 15 or Eq. 16) depends largely on the independence of X_R and X_S . If
 697 we define the likelihood using Eq. 18 as it was, we must convolve the DNA amounts not on the two-
 698 dimensional parameter space spanned by X_R and X_S , but a three-dimensional space by $X_{R | \text{homo}}$, $X_{S | \text{hetero}} =$
 699 $X_{R | \text{hetero}}$, and $X_{S | \text{homo}}$, which would increase the calculation time by 1,000 to 10,000 times.

700 Therefore, we removed the constraint and assumed that $X_{R | *}$ and $X_{S | *}$ were distributed independently and
 701 identically; that is, instead of the heterozygotes, we captured $n - m_1 - m_0$ individuals of haploid R and
 702 another $n - m_1 - m_0$ individuals of haploid S separately. Regarding homozygotes, we also assumed that we
 703 captured $2m_1$ R haploids and $2m_0$ S haploids instead of m_1 RR and m_0 SS, respectively. Then,

704
$$\begin{aligned} X_{R|homo} &\sim \text{Ga}(2m_1 k, \theta), & X_{R|hetero} &\sim \text{Ga}((n - m_1 - m_0)k, \theta), \\ X_{S|hetero} &\sim \text{Ga}((n - m_1 - m_0)k, \theta) \text{ i. i. d.}, & X_{S|homo} &\sim \text{Ga}(2m_0 k, \theta). \end{aligned}$$

705 *Eq. 19*

706 Finally, we can approximate the DNA amounts of a diploid organism in the bulk sample by simply
707 substituting Eq. 3:

708
$$X_R \sim \text{Ga}((n + m_1 - m_0)k, \theta), \quad X_S \sim \text{Ga}((n - m_1 + m_0)k, \theta).$$

709 *Eq. 20*

710 In addition, at probability $\text{Bin}(0|2n_h, p)$, all (hypothetically haploid) individuals become S or R; in that case,
711 there is no need to convolve the DNA amounts.

712

713

Table and figure captions

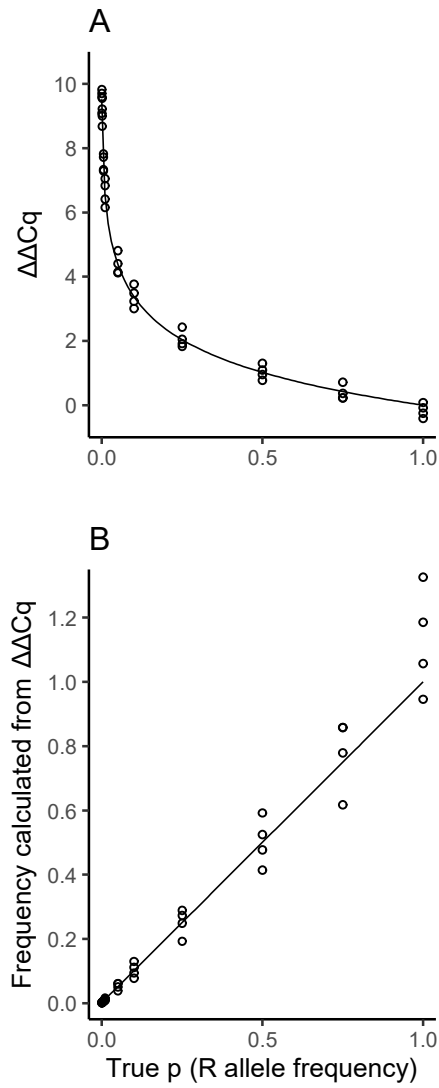
714

Table 1. Description of variables and parameters

Symbol	Description	Range	Arguments in the R package
p	Frequency of the R (resistant) allele in a population	$0 \leq p \leq 1$	P
X_S, X_R	Amounts of DNA belonging to S (susceptible) or R alleles included in a bulk sample	$X_S \geq 0, X_R \geq 0$	—
Y_R	The observed frequency of R in the bulk sample, defined as $X_R/(X_R + X_S)$	$0 \leq Y_R \leq 1$	—
k, θ	Shape and scale parameters of the gamma distribution $\text{Ga}(k, \theta)$	$k > 0, \theta > 0$	K
N	Number of bulk samples taken from a population, each of which consists of n_h individuals ($h = 1, 2, 3, \dots, N$)	$N \in \mathbb{N}$	ntrap
n, n_h	Number of individuals constituting the (h th) bulk sample	$n \in \mathbb{N}$	npertrap
m, m_h	Number of R individuals included in the (h th) bulk sample	$0 \leq m \in \mathbb{Z} \leq n$	m (as an internal variable)

qPCR-related variables and parameters			
η	Per-cycle efficiency in the PCR amplification (as $1 + \eta$)	$\eta > 0$	EPCR
X_0, X_θ	Initial amount of template DNA and the termination threshold of the amplification in the real-time PCR process	$X_0 > 0, X_\theta > 0$	X_θ is fixed 1 in the package
τ	Cq value: the number of PCR amplification cycles before termination	$\tau \in \mathbb{R}$	τ_h^{TW} : target0, τ_h^{TD} : target1, τ_h^{HW} : housek0, τ_h^{HD} : housek1
δ_T	Relative content of the target gene to the internal reference (housekeeping gene)	$\delta_T > 0$	targetScale
δ_B	(In RED- $\Delta\Delta\text{Cq}$ method) the locus-independent change rate of the template DNA quantity accompanying the restriction enzyme treatment.	$\delta_T > 0$	baseChange
z	(In RED- $\Delta\Delta\text{Cq}$ method) residual rate of restriction enzyme digestion, or (in general $\Delta\Delta\text{Cq}$ analyses) portion of the off-target allele amplified in the PCR	$z > 0$	zeroAmount
ε_c	Cq measurement error (standard deviation)	$\varepsilon_c > 0$	sdMeasure

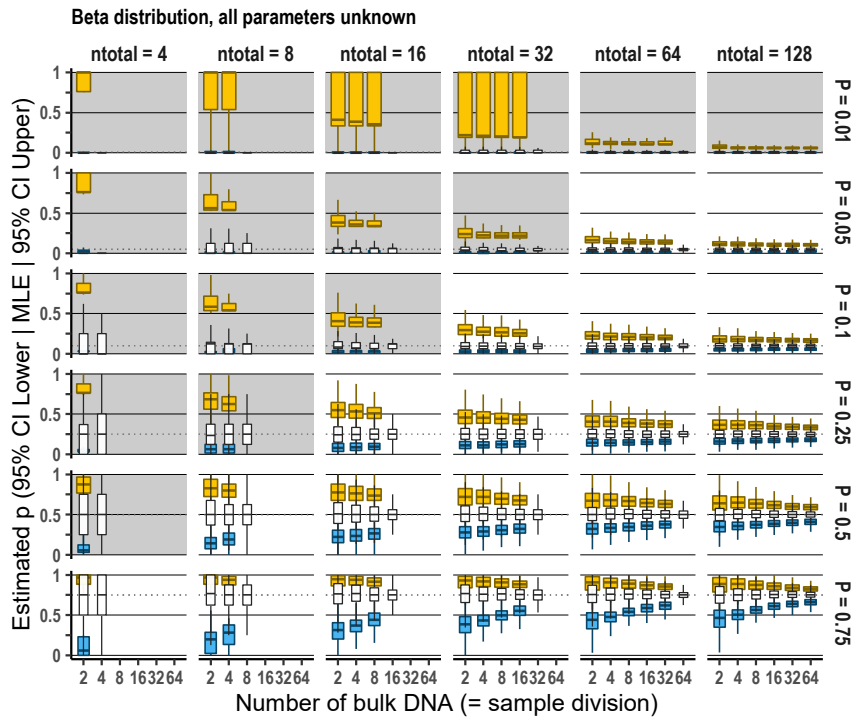
715



716

717 Figure 1 The relationship between the allele frequency in the sample and A: the RED- $\Delta\Delta Cq$ measures, B: the
718 observed frequency calculated as $(1 + \eta)^{-\Delta Cq}$, showing the results of etoxazole resistance in the two-
719 spotted spider mites. The lines are not the regression on the actual Cq measurement (shown as points), but
720 the theoretical relationship between true frequency of the R allele and the quantity defined as A:
721 $-\ln(z + Y_R(1 - z))/\ln(1 + \eta)$ or B: $z + Y_R(1 - z)$, where $Y_R = X_R/(X_R + X_S)$. Parameters are $z =$
722 0.00156 and $\eta = 0.971$.

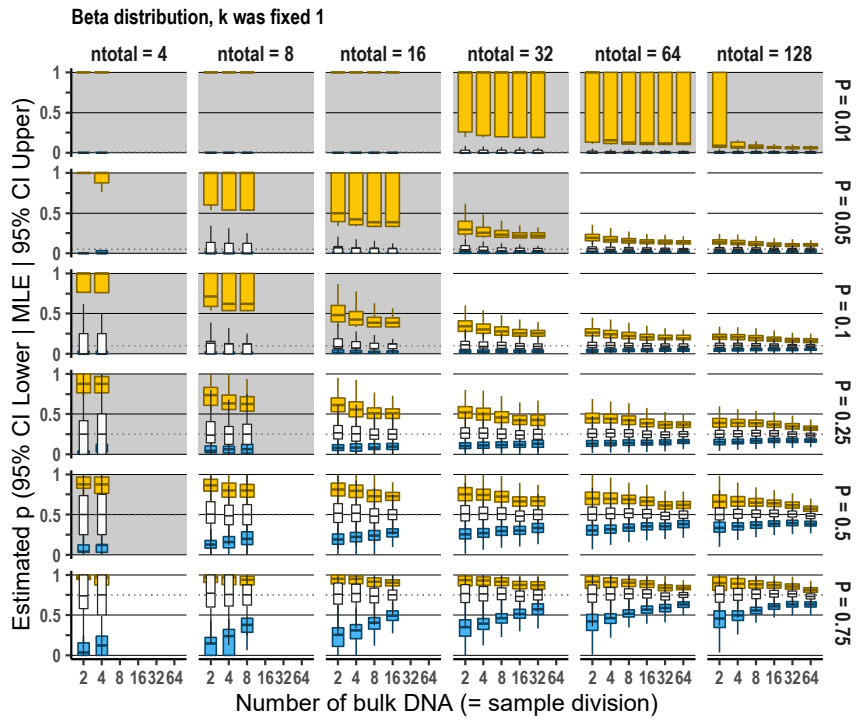
723



724

725 Figure 2 Estimation accuracy of the resistance allele frequency, p , with `freqpcr()` when the beta distribution
 726 was assumed, and all estimable parameters (P , K , `targetScale`, and `sdMeasure`) were set as unknown. The
 727 result of numerical experiments based on 1,000 dummy datasets per parameter region. The x-axes
 728 correspond to the parameter “`ntrap`.” The three box plots (white thin, blue, and yellow wide) in each region
 729 show the maximum likelihood estimates (MLE), lower bound of the 95% CI, and the upper bound,
 730 respectively. In each boxplot, the horizontal line signifies the median of the simulations, hinges of the box
 731 show 25 and 75 percentiles, and the upper/lower whiskers correspond to the $1.5 \times$ interquartile ranges. The
 732 shaded facets show that the total sample sizes ($ntotal$) are smaller than $3/p$.

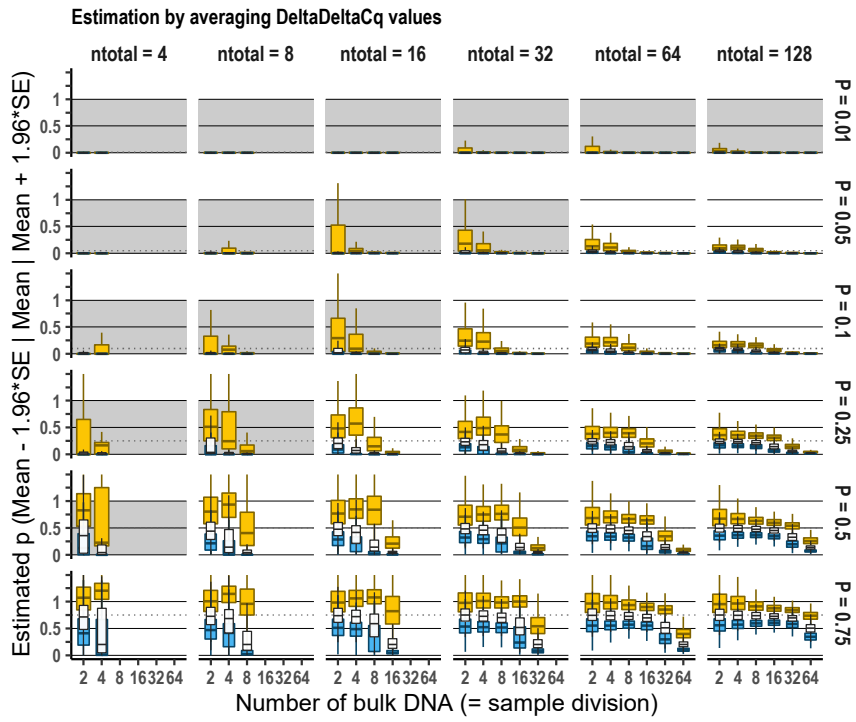
733



734

735 Figure 3 Estimation accuracy of the resistance allele frequency with `freqpcr()` when the beta distribution was
736 assumed, fixing $K = 1$.

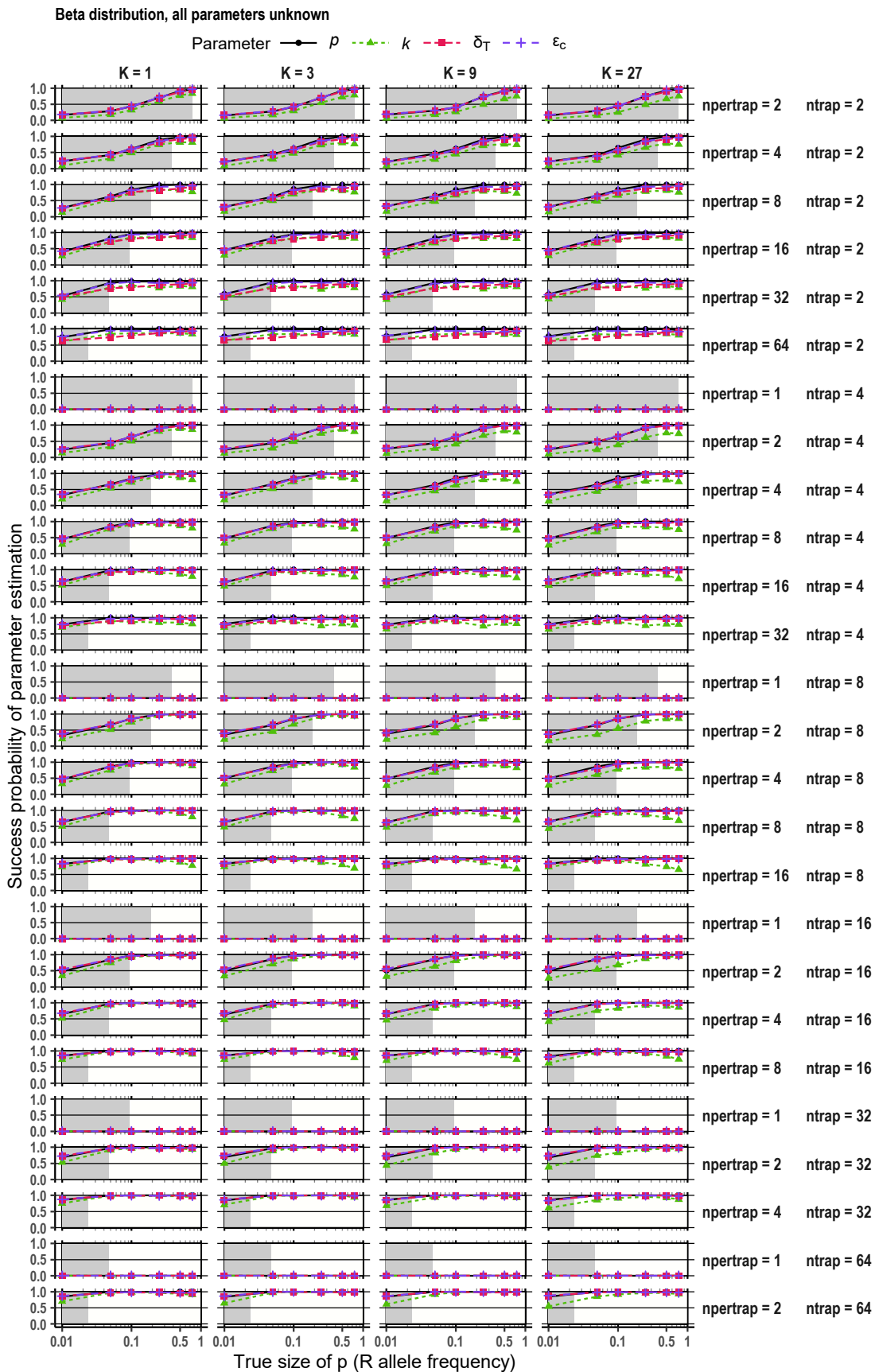
737



738

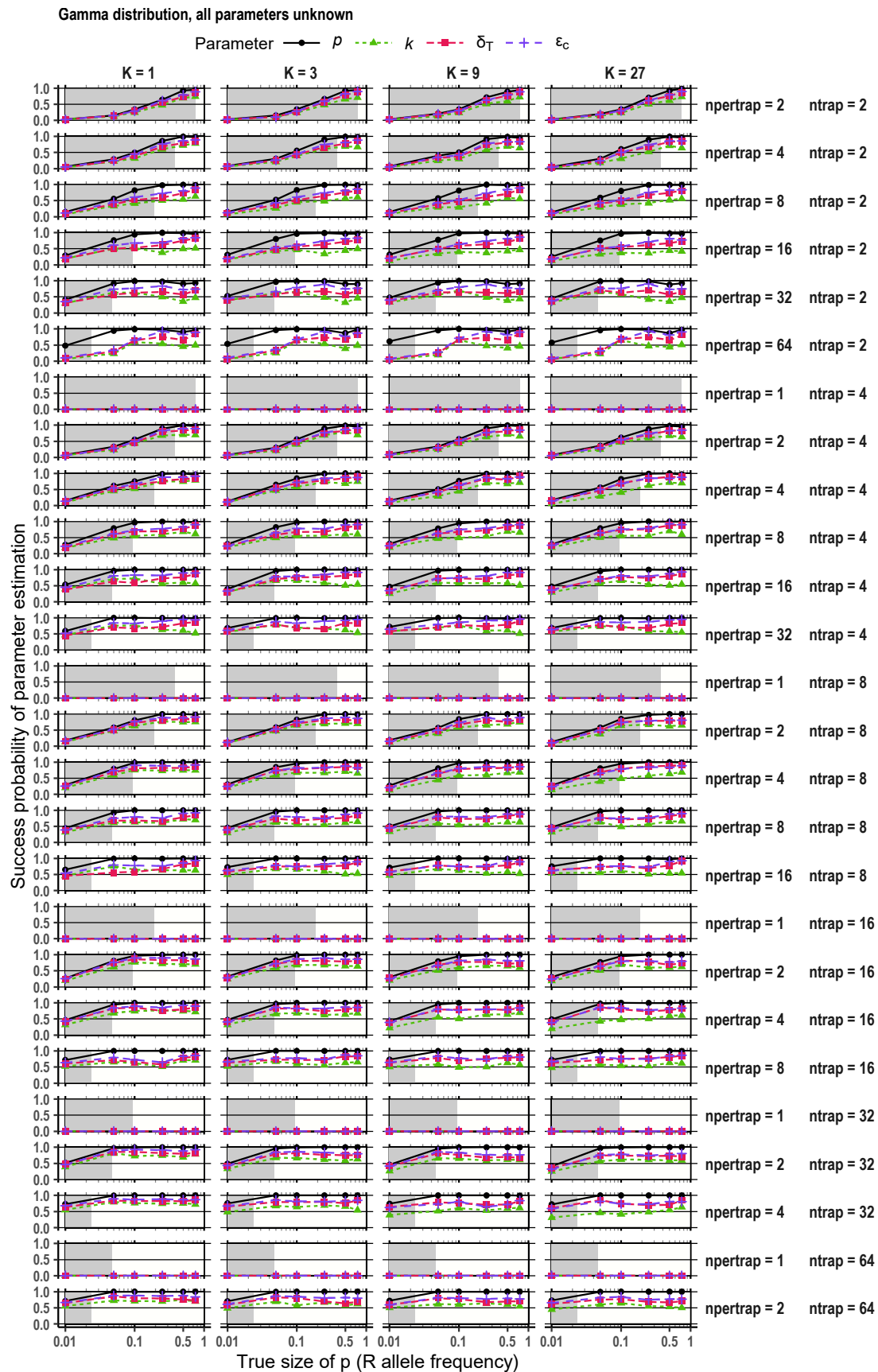
739 Figure 4 Estimation accuracy of the resistance allele frequency by simple averaging of $\Delta\Delta Cq$ measures. The
740 dummy dataset was derived from the numerical experiment of “beta distribution, all parameters unknown.”

741



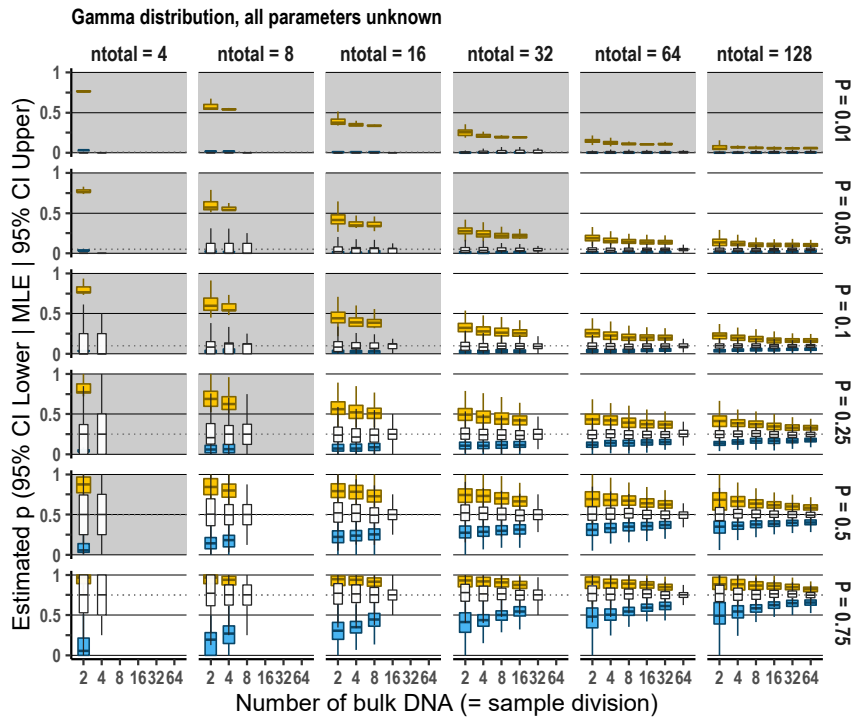
742

743 Figure S1 Probability of estimation success with freqpcr() in each parameter region. The beta distribution
744 was assumed, and all estimable parameters (P , K , targetScale, and sdMeasure) were set as unknown. The
745 shaded boxes in the background show the frequency ranges where the total sample sizes (n_{total}) are smaller
746 than $3/p$.



747

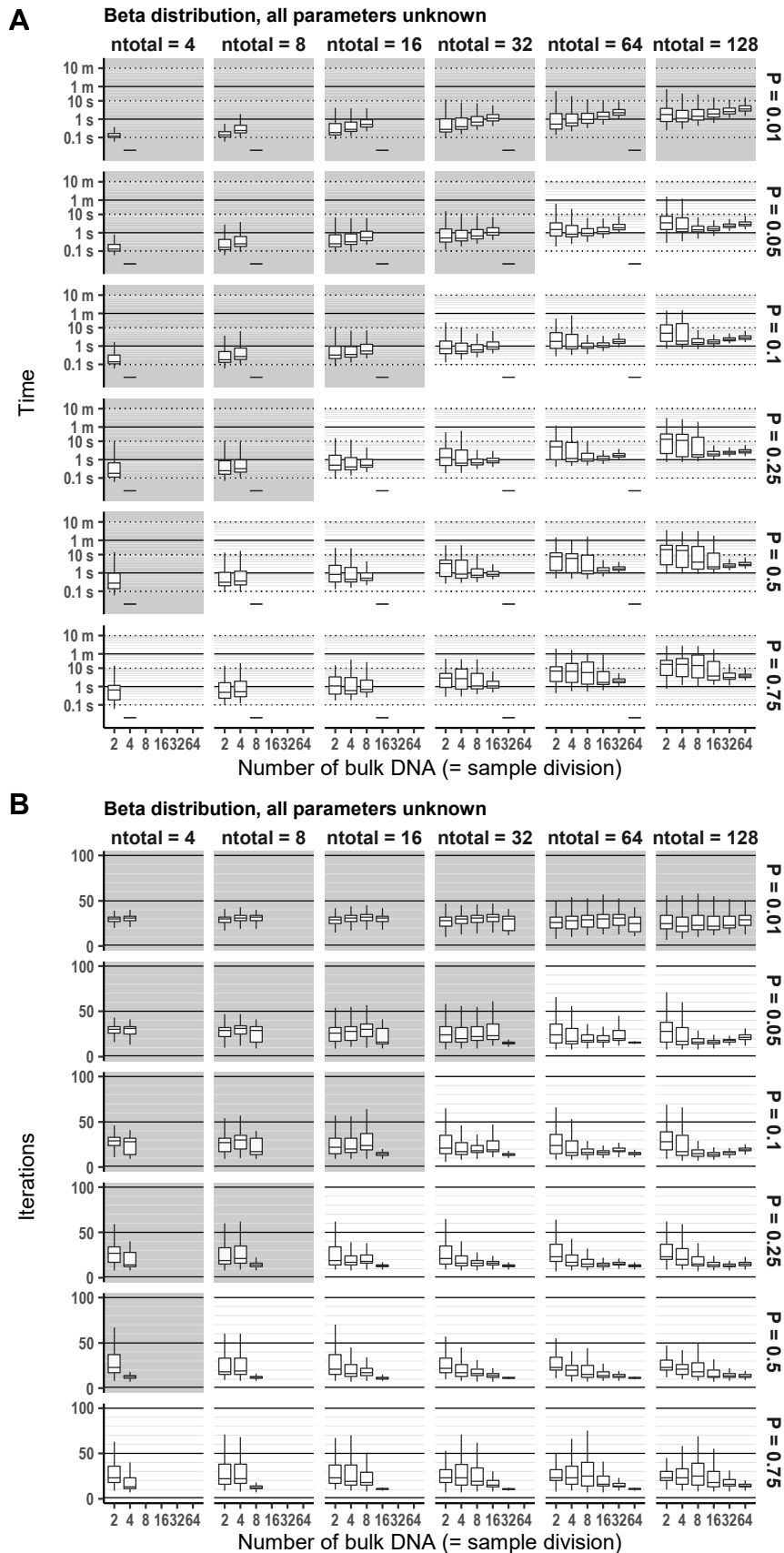
748 Figure S2 Probability of estimation success with freqpcr() in each parameter region. The gamma
749 distributions were assumed, and all estimable parameters were set as unknown. The function often failed to
750 calculate the CIs for k when npertrap (individuals in each bulk sample) were larger, possibly due to the
751 accumulation of numerical calculation error.



752

753 Figure S3 Estimation accuracy of p with `freqpcr()` when gamma distributions were assumed and all estimable
754 parameters were set as unknown. The shaded facets show that the total sample sizes (n_{total}) are smaller than
755 $3/p$.

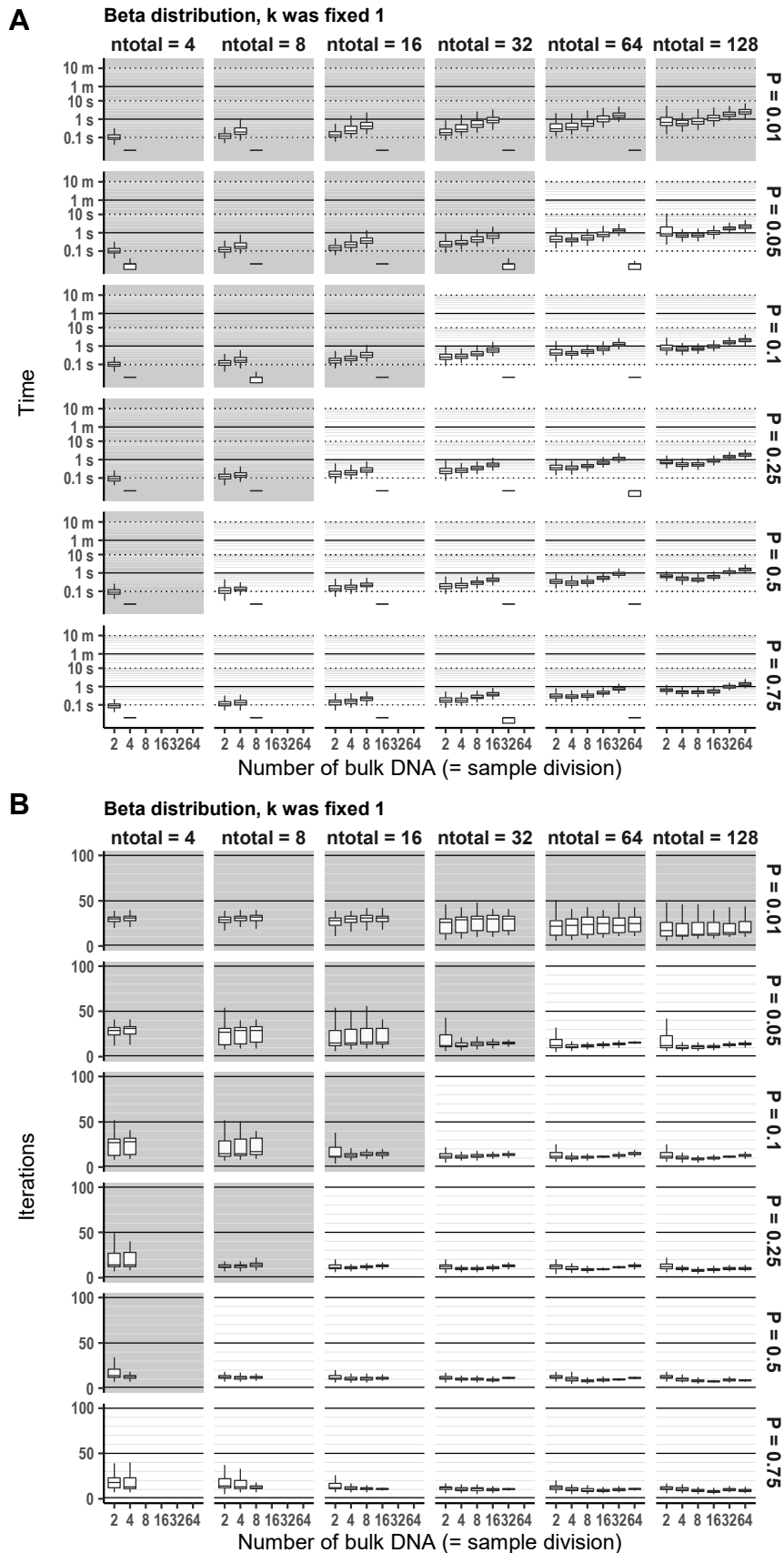
756



757

758 Figure S4 A: Calculation time and B: number of iterations until the `freqper()` function converges. The beta
759 distribution was assumed, and all estimable parameters were set as unknown.

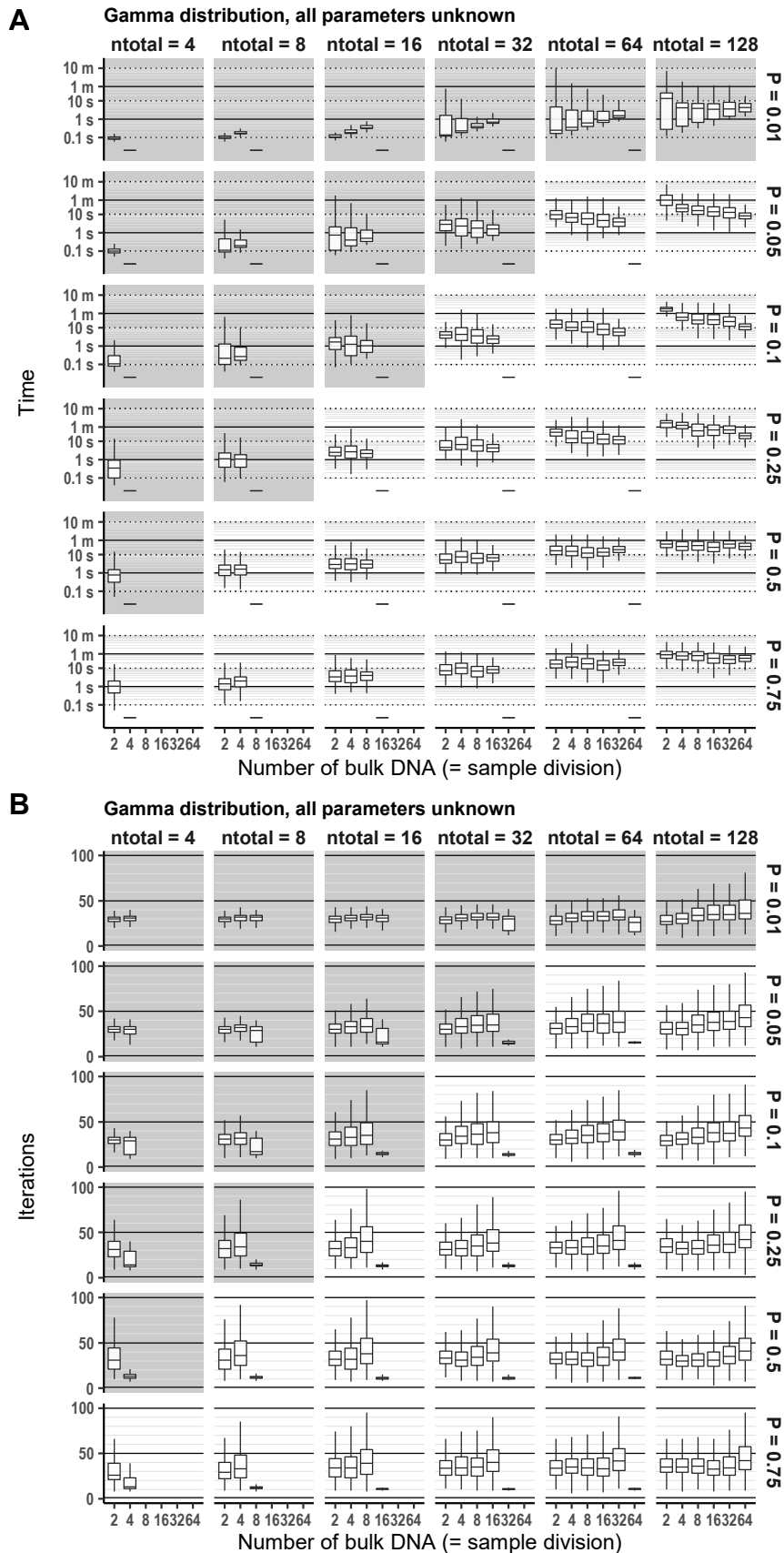
760



761

762 Figure S5 A: Calculation time and B: number of iterations until the `freqper()` function converges. The beta
 763 distribution was assumed, fixing the gamma shape parameter $K = 1$.

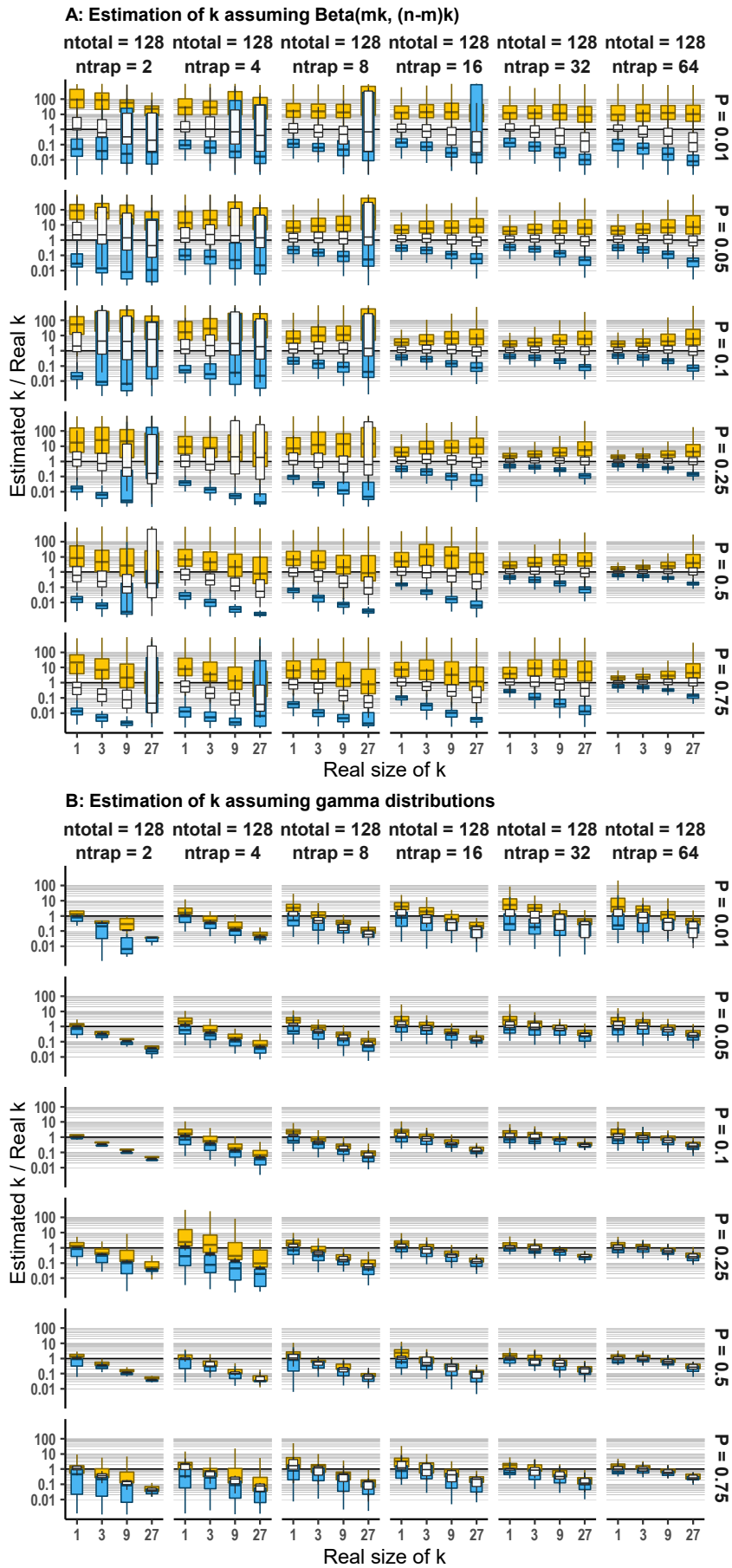
764



765

766 Figure S6 A: Calculation time and B: number of iterations until the `freqpcr()` function converges, assuming
767 gamma distributions. All estimable parameters were set as unknown.

768



769

770 Figure S7 Estimation accuracy of k (the gamma shape parameter) in the simulation, showing the maximum
771 likelihood estimate by `freqpccr()` divided by the actual parameter size.