

1 **freqpcr: estimation of population allele frequency using**  
2 **qPCR  $\Delta\Delta Cq$  measures from bulk samples**

3 Short running title: Allele frequency estimation based on  $\Delta\Delta Cq$

4  
5 Masaaki Sudo<sup>\*1</sup>, Masahiro Osakabe<sup>2</sup>

6 <sup>1</sup> Tea Pest Management Unit, Institute of Fruit Tree and Tea Science, NARO: Kanaya Tea Research  
7 Station, 2769, Shishidoi, Kanaya, Shimada, Shizuoka 428-8501, Japan

8 <sup>2</sup> Laboratory of Ecological Information, Graduate School of Agriculture, Kyoto University: Kyoto,  
9 Japan

10  
11 \* To whom correspondence: [masaaki@sudori.info](mailto:masaaki@sudori.info)

12 Telephone: +81-547-45-4419

13  
14 ORCID ID:

15 0000-0001-9834-9857 (Masaaki Sudo)

16 0000-0002-2246-3431 (Masahiro Osakabe)

17  
18 Manuscript type: Resource article: Molecular and Statistical Advances

## 20 Abstract

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

35

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

38

## 39 Introduction

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

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

53 In accordance with the development of genome-wide association studies on resistance genes  
54 (Frensch-Constant, 2013; Snoeck et al., 2019; Sugimoto et al., 2020), rapid advancements have  
55 recently been made in molecular diagnostics (Donnelly et al., 2016; Samayoa, et al., 2015; Toda et al.,  
56 2017). To quantify resistance-associated point mutations at the population scale, the most fundamental  
57 molecular technique is an individual-based polymerase chain reaction (PCR) analysis (Toda et al.,  
58 2017). Moreover, quantitative PCR (qPCR), based on real-time PCR, is also used for the point  
59 mutation of allele frequencies (Germer et al., 2000). If the alleles are distributed randomly in the target  
60 population, a simple binomial assumption enables us to estimate the population allele frequency and  
61 its confidence interval. However, collecting, processing, and analyzing multiple DNA samples may  
62 not be feasible, particularly when dealing with numerous samples from multiple sites, or when estimation

63 of a rare (<1%) mutation frequency is required for a given population, as is often the case in the early  
64 phase of resistance development.

65 Although rearing living insects is no longer necessary, the field of molecular diagnostics is still  
66 lacking a metaphorical silver bullet capable of reducing the required time and cost associated with  
67 handling multiple samples, while guaranteeing estimation precision and accuracy. The use of a “bulk  
68 sample” (i.e., pooling multiple individual samples and processing a single DNA extract), in  
69 coordination with statistical methods, such as group testing, may address some of these challenges. In  
70 fact, Osakabe *et al.* (2017) and Maeoka *et al.* (2020) developed diagnostic methods for detecting  
71 resistance to the acaricide, etoxazole, in the two-spotted spider mite, *Tetranychus urticae* Koch (Acari:  
72 Tetranychidae), which is conferred by an amino acid substitution in chitin synthase 1 (*CHS1*; I1017F)  
73 (Van Leeuwen *et al.*, 2010). They used a bulk sample to measure the frequency of the resistant point  
74 mutation in field mite populations. To calculate the point estimate, these studies compared the relative  
75 quantity of the resistance allele with an internal reference (housekeeping gene) in the sample, known  
76 as the  $\Delta\Delta Cq$  method (Livak and Schmittgen, 2001). In the etoxazole-R diagnosis by Osakabe *et al.*  
77 (2017), glyceraldehyde-3-phosphate dehydrogenase (*GAPDH*) was used as the housekeeping gene.

78 In this study, we propose a statistical method for obtaining the interval estimate of allele frequency  
79 using  $\Delta\Delta Cq$ -based qPCR analyses for multiple bulk samples taken from a population. We first  
80 introduced the random error structure to approximate the relative abundance of the two alleles (wild  
81 type and mutant) and their ratios in the bulk DNA sample. Thereafter, we formulated how the relative  
82 amounts of the two alleles in a sample solution impacted the Cq measurements through  $\Delta\Delta Cq$ -based  
83 qPCR analysis. Finally, we combined the models and developed a maximum likelihood estimation  
84 procedure to estimate an allele frequency implemented using the R language. The package source is  
85 available on the Internet (<https://github.com/sudoms/freqpcr>).

## 86 **Model**

### 87 *Approximation of allele quantities contained in a bulk DNA sample*

88 When DNA is directly extracted from the whole body of a living organism, the DNA yield is roughly  
89 proportional to its body weight (Chen et al., 2010). For insects, the intra-population frequency  
90 distribution of body weight is often approximated using a unimodal and right-skewed continuous  
91 distribution, typically lognormal or gamma distribution (May, 1976; Rakovski et al., 2011; Knapp,  
92 2016). In fact, it has been suggested that body weights are distributed lognormally in many non-social  
93 insect species (Gouws et al., 2011).

94 In this study, we adopted a gamma, rather than lognormal, distribution to approximate the DNA  
95 amount per individual organism for two reasons. First, it is difficult to distinguish which distribution a  
96 real population obeys when the sample size is small. The two distributions are considered  
97 interchangeable (Wiens, 1999; Kundu and Manglick, 2005). Second, the sum and proportion of  
98 independent gamma distributions have closed forms under certain conditions. Using Eq. 1, let  $X$  ( $X \geq$   
99 0) be the DNA yield per single locus per individual:

$$100 \quad \text{Ga}(X|k, \theta) = \frac{1}{\Gamma(k)} \left(\frac{1}{\theta}\right)^k X^{k-1} \exp\left(-\frac{X}{\theta}\right),$$

101 *Eq. 1*

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

104 Using Eq. 1, let us consider the amounts of allelic DNA in the sample extracted from multiple  
105 individuals at once, hereafter referred to as “a bulk sample.” Table 1 lists the variables and parameters  
106 of the model structure. For simplicity, we model the case of haploidy in the main text, while Appendix  
107 S1 describes the approximated formulation for diploids. Now, we have  $n$  insects, of which  $m$  ( $m =$   
108  $0, 1, \dots, n$ ) are the genotypes resistant to an insecticide (hereafter denoted by R). The rest  $n - m$   
109 carried S, the susceptible allele. When we capture insects from a wild population, the size of  $n$  is

110 obvious, however  $m$  is usually unknown (Figure 1A). Assuming random sampling from an infinite  
111 population with the R allele at the frequency  $p$ ,  $m$  follows a binomial distribution (Eq. 2):

$$112 \quad \text{Bin}(m|n, p) = \frac{n!}{m!(n-m)!} p^m (1-p)^{n-m}.$$

113 Eq. 2

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

$$120 \quad X_R \sim \text{Ga}(mk, \theta),$$
$$121 \quad X_S \sim \text{Ga}((n-m)k, \theta).$$

122 Eq. 3

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

$$126 \quad \text{Beta}(Y_R|mk, (n-m)k) = \frac{Y_R^{mk-1} (1-Y_R)^{(n-m)k-1}}{B(mk, (n-m)k)},$$

127 Eq. 4

128 where  $B(\cdot)$  is a beta function. This error structure was originally developed to model allele frequencies  
129 measured via quantitative sequencing (Sudo et al., in press).

## 130 ***Relative quantification of DNA by real-time PCR***

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

132 The  $\Delta\Delta Cq$  (quantification cycle) method (Livak, 1997) is the most common method for relative  
133 quantification using qPCR, in which the quantities of complementary cDNA libraries are compared  
134 between samples to determine the relative expression levels of the genes of interest. Osakabe *et al.*  
135 (2017) expanded this concept and proposed the “RED- $\Delta\Delta Cq$  method” (RED, restriction enzyme  
136 digestion), a derivative method that can measure the allele frequency from a single sample solution, to  
137 diagnose the regional prevalence of an acaricide-resistant point mutation in a *T. urticae* population.

138 In the RED- $\Delta\Delta Cq$  method, the control was prepared as an intact sample containing total DNA (=  $X_R + X_S$ )  
139 on the target locus. The sample in question was the same DNA extract, however, was  
140 digested with restriction endonucleases prior to qPCR analysis (Figure 1A). The restriction site is  
141 designed to recognize the S allele on the target locus so that the operation digests the major part of S  
142 (denoted by  $1 - z$ :  $z$  is a small yet, positive variable giving the residual rate). Consequently, we  
143 obtained the template amount  $X_R + zX_S$  at the target locus after digestion. To calibrate the template  
144 DNA amounts, the samples before and after digestion were also amplified using the primer set for a  
145 housekeeping gene as an internal reference.

146 Taken together, the single bulk sample results in a quartet of Cq measurements differentiating at the  
147 target loci (resistance-associated and housekeeping genes)  $\times$  restriction enzyme digestion (undigested  
148 and digested). We can then formulate the allele frequencies by letting  $X^{HW}$  and  $X^{TW}$  represent the total  
149 amounts of template DNA at the housekeeping (H) and target (T) loci, respectively, included in the  
150 sample without digestion, the state denoted by W (Figure 1C).

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

152 *Eq. 5*

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

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

Eq. 6

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

160 As a result of qPCR, the Cq quartet,  $\tau^{HW}$ ,  $\tau^{TW}$ ,  $\tau^{HD}$ , and  $\tau^{TD}$  were obtained as:

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

Eq. 7

164 Here,  $1 + \eta$  ( $\eta > 0$ ) and  $X_\Theta$  denote the amplification efficiency per PCR cycle and its threshold,  
 165 respectively. According to Livak and Schmittgen (2001), we assume an ideal amplification, where  $X_\Theta$   
 166 is set within the early exponential amplification phase. The actual Cq data contain measurement errors  
 167 in addition to uncertainty due to experimental operations, such as sample dispensation or PCR  
 168 amplification. We express these using the common error term  $\varepsilon_c \sim N(0, \sigma_c^2)$ , following the normal  
 169 distribution of mean = 0 and variance =  $\sigma_c^2$  in the scale of raw Cq values. The validity of this error  
 170 structure is verified later.

171 The two  $\Delta Cq$  values were then defined for the undigested and digested samples, as  $\Delta\tau^W = \tau^{TW} -$   
 172  $\tau^{HW}$  and  $\Delta\tau^D = \tau^{TD} - \tau^{HD}$ , respectively. Their  $\Delta\Delta Cq$  are:



173 
$$\Delta\Delta\tau = \Delta\tau^D - \Delta\tau^W = -\frac{\ln\left(\frac{X_R + zX_S}{X_R + X_S}\right)}{\ln(1 + \eta)} + \varepsilon, \quad \varepsilon \sim N(0, 4\sigma_c^2).$$

174 Eq. 8

175 From Eq. 8, the expected value of  $(X_R + zX_S)/(X_R + X_S)$  is calculated as  $(1 + \eta)^{-\Delta\Delta\tau}$ . The  
176 coefficients  $\delta_B$  and  $\delta_T$  in Eq. 5 and Eq. 6 vanished by subtracting the Cq values and  $\Delta Cq$  values,  
177 respectively.

178 The point estimate of the resistance allele frequency,  $\hat{Y}_R$ , is defined as  $X_R/(X_R + X_S)$  for each bulk  
179 sample. When  $z$  is much smaller than  $\hat{Y}_R$ , the quantity  $(X_R + zX_S)/(X_R + X_S) = \hat{Y}_R + z(1 - \hat{Y}_R)$  itself  
180 can approximate the frequency, which will be the case with enough digestion time before qPCR.

181 However, the use of the point estimate may introduce a problem in that the size of  $\hat{Y}_R$  often exceeds 1  
182 when the R frequency is high and a larger error exists in the Cq measurement (also see the result of  
183 Experiment 2).

184 Although the value of  $1 + \eta$  may vary on the primer sets, both target and housekeeping loci share  
185 the same amplification efficiency in Eq. 7, because practical PCR protocols were designed to be  $1 +$   
186  $\eta \cong 2$ . We can also approximately cancel the effect of heterogeneous amplification efficiencies by  
187 fitting the  $\delta_T$  size of the sample sets with known allele ratios (Experiment 1).

### 188 *Measurement of $\Delta\Delta Cq$ using allele-specific primer sets*

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

193 Maeoka *et al.* (2020) demonstrated that a general  $\Delta\Delta Cq$  method without restriction enzyme  
194 treatment could be used for allele-frequency measurement if a specific primer set were to be designed  
195 to amplify only the R allele at the target locus. Similar to the RED- $\Delta\Delta Cq$  method, DNA samples with

196 unknown mixing ratios were dispensed and amplified using primer sets corresponding to T and H loci,  
197 respectively. Unlike the RED- $\Delta\Delta\text{Cq}$  method, the control sample was not taken from the test sample  
198 solution, but rather was prepared as a DNA solution containing 100% R, hereafter denoted as U (=   
199 pUre R line) (Figure 1C).

200  $X^{\text{HU}}$  and  $X^{\text{TU}}$  then denote the template DNA quantities:

$$\begin{aligned} 201 \quad X^{\text{HU}} &= X'_{\text{R}}, \\ X^{\text{TU}} &= \delta_{\text{T}} X'_{\text{R}}. \end{aligned}$$

202 Eq. 9

203 Though the definition of  $\delta_{\text{T}}$  is the same as Eq. 5, the quantity is denoted by  $X'_{\text{R}}$  instead of  $X_{\text{S}} + X_{\text{R}}$  as  
204 it no longer originates from the R portion of the test sample itself (i.e., not internal).

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

$$\begin{aligned} 207 \quad X^{\text{HV}} &= X_{\text{R}} + X_{\text{S}}, \\ X^{\text{TV}} &= \delta_{\text{T}}(X_{\text{R}} + zX_{\text{S}}). \end{aligned}$$

208 Eq. 10

209 In the PCR process of the modified  $\Delta\Delta\text{Cq}$  method, the small positive number  $z$  provides the template  
210 quantity of S, which is non-specifically amplified even with the R-specific primer set. As the primer  
211 set for the housekeeping gene was nonspecific,  $X^{\text{HV}}$  was fully amplified. Assuming that all four  
212 template DNAs are amplified with efficiency  $1 + \eta$ , we define the two  $\Delta\text{Cq}$  values as  $\Delta\tau^{\text{U}} = \tau^{\text{TU}} -$   
213  $\tau^{\text{HU}}$  and  $\Delta\tau^{\text{V}} = \tau^{\text{TV}} - \tau^{\text{HV}}$ . Finally, their  $\Delta\Delta\text{Cq}$  values are  $\Delta\Delta\tau = \Delta\tau^{\text{V}} - \Delta\tau^{\text{U}}$ , which yields a formula  
214 identical to Eq. 8.

### 215 ***Interval estimation of allele frequency and experimental parameters based on*** 216 ***qPCR over multiple bulk samples***

217 Finally, we consider the likelihood model to obtain the interval estimate of the allele frequency based  
218 on the (RED-) $\Delta\Delta\text{Cq}$  analysis over multiple bulk samples. Assume that the population has the R allele

219 at the frequency  $p$  from which  $N$  bulk samples are taken. The  $h$ th sample ( $h = 1, 2, 3, \dots, N$ ) consists of  
 220  $n_h$  haploid individuals, of which  $m_h$  are resistant mutants. As shown in Eq. 7, the Cq values (denoted  
 221 as  $\tau_h^{\text{HW}}$ ,  $\tau_h^{\text{TW}}$ ,  $\tau_h^{\text{HD}}$ , and  $\tau_h^{\text{TD}}$  for each bulk sample) are determined not only by the DNA quantities,  
 222 denoted as  $X_{h,R}$  and  $X_{h,S}$ , but also by parameters such as  $\delta_T$  or  $\sigma_c^2$  accompanying the experimental  
 223 operation. We can simultaneously estimate these if we have multiple bulk samples, for which the  
 224 likelihood function of obtaining the Cq values under the parameters is defined.

225 We propose the joint likelihood for the two  $\Delta\text{Cq}$  values,  $\Delta\tau_h^{\text{W}} = \tau_h^{\text{TW}} - \tau_h^{\text{HW}}$  and  $\Delta\tau_h^{\text{D}} = \tau_h^{\text{TD}} - \tau_h^{\text{HD}}$ ,  
 226 for the convenience of numerical calculation:

$$227 \quad \Delta\tau_h^{\text{W}} \sim \text{N}\left(-\frac{\ln\delta_T}{\ln(1+\eta)}, 2\sigma_c^2\right),$$

$$228 \quad \Delta\tau_h^{\text{D}} \sim \text{N}\left(-\frac{\ln\delta_T + \ln\left(\frac{X_{h,R} + zX_{h,S}}{X_{h,R} + X_{h,S}}\right)}{\ln(1+\eta)}, 2\sigma_c^2\right).$$

229 *Eq. 11*

230 Although Eq. 11 is defined for the RED- $\Delta\Delta\text{Cq}$  method, it is also applicable to the  $\Delta\Delta\text{Cq}$  method by  
 231 Maeoka *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  $\Delta\tau_h^{\text{V}} = \tau_h^{\text{TV}} - \tau_h^{\text{HV}}$ ,  
 232 respectively.

### 233 *Formulation of likelihood based on gamma or beta distribution*

234 Using the relationship between  $m_h$ ,  $n_h$ , and  $p$  in Eq. 2, we proceed to the likelihood function defined  
 235 as the probability of observing the set of  $\Delta\tau_h^{\text{W}}$  and  $\Delta\tau_h^{\text{D}}$  under the given values of  $p$ ,  $n_h$ , and other  
 236 experimental parameters. In Eq. 11,  $\Delta\tau_h^{\text{W}}$  is not affected by the R : S ratio in the bulk sample; it is only  
 237 affected by the experimental parameters,  $\delta_T$ ,  $\eta$ , and  $\sigma_c^2$ . In addition, by taking the differences, there is  
 238 no need to estimate as  $X_\theta$  and  $\delta_B$  appear in Eq. 7. Moreover, cancelation of  $\delta_B$  also ensures that we  
 239 can apply the model of Eq. 11 to the general  $\Delta\Delta\text{Cq}$  method of Eq. 9 and Eq. 10.

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

246 The likelihood function for the observed  $\Delta Cq$  values on the  $h$ th bulk sample  $L_h$  is defined as  
 247 follows:

$$248 \quad L_h = P(\Delta\tau_h^W|\delta_T, \eta, \sigma_c^2) \sum_{m_h=0}^{n_h} [\text{Bin}(m_h|n_h, p)P(\Delta\tau_h^D|m_h, \delta_T, z, \eta, \sigma_c^2)],$$

$$249 \quad P(\Delta\tau_h^D|m_h, \delta_T, z, \eta, \sigma_c^2) = \begin{cases} N\left(-\frac{\ln(z\delta_T)}{\ln(1+\eta)}, 2\sigma_c^2\right) & (m_h = 0) \\ \psi_G \text{ or } \psi_B & (m_h = 1, 2, \dots, n_h - 1) \\ N\left(-\frac{\ln\delta_T}{\ln(1+\eta)}, 2\sigma_c^2\right) & (m_h = n_h) \end{cases}$$

250 Eq. 12

251 In Eq. 12,  $\psi_G$  or  $\psi_B$  denotes the probability of obtaining  $\Delta\tau_h^D$  under the template DNA quantities of  
 252  $X_{h,R|m_h} = r$  and  $X_{h,S|m_h} = s$  if we formularize the two quantities by gamma distribution, or if we  
 253 formularize their mixing ratio by the single beta distribution, respectively. We must consider not only  
 254 the possible cases of  $m_h$ , but also the entire range of the DNA amounts. If we use the gamma  
 255 distributions, for every case  $m_h = 1, 2, \dots, n_h - 1$ , we need to calculate the double integration for  $\psi_G$   
 256 under the whole region of  $X_{h,R|m_h} = r$  and  $X_{h,S|m_h} = s$  for the interval  $\{D: 0 \leq r < \infty, 0 \leq s < \infty\}$ .

$$257 \quad \psi_G = \iint_D N\left(-\frac{\ln\delta_T + \ln\left(\frac{r+zs}{r+s}\right)}{\ln(1+\eta)}, 2\sigma_c^2\right) \text{Ga}(r|m_h k, \theta) \text{Ga}(s|(n_h - m_h)k, \theta) dr ds.$$

258 Eq. 13

259 The common scale parameter of the gamma distributions,  $\theta$ , is not identifiable from the data, although  
 260 we can substitute arbitrary values  $\theta = 1$  for it because it is canceled in  $\ln[(r + zs)/(r + s)]$  in Eq. 13.

261 Since the computational burden for the double integration is large, we simplified the likelihood  
262 model with the beta distribution. By introducing  $y = r/(r + s)$ , the probability of obtaining  $\Delta\tau_h^D$  is  
263 replaced with  $\psi_B$  defined as follows:

$$264 \quad \psi_B = \int_0^1 N\left(-\frac{\ln\delta_T + \ln(z + y(1 - z))}{\ln(1 + \eta)}, 2\sigma_c^2\right) \text{Beta}(y|m_h k, (n_h - m_h)k) dy.$$

265 Eq. 14

266 We provide an R function “freqpcr()” to estimate the parameters  $p$ ,  $k$ ,  $\delta_T$ , and  $\sigma_c$  simultaneously when  
267 the set of Cq measurements ( $\tau_h^{HW}$ ,  $\tau_h^{TW}$ ,  $\tau_h^{HD}$ , and  $\tau_h^{TD}$ ) and  $n_h$  are given for each of the  $N$  bulk  
268 samples. The default is freqpcr(..., beta = TRUE), where the beta distribution model of Eq. 14 was  
269 used instead of gamma. Regardless of the algorithms, the asymptotic confidence intervals are  
270 calculated using the inverse of the Hessian matrix evaluated at the last iteration. The functions nlm()  
271 of R and cubintegrate() in the R package “cubature” (Narasimhan et al., 2019) are used for the  
272 iterative optimization and the (double) integration, respectively.

### 273 ***Identification of auxiliary parameters using DNA samples with known allele-*** 274 ***mixing ratios***

275 The likelihood introduced above ensures that we can estimate the sizes of  $p$  and  $k$  together with other  
276 experimental parameters if we have conducted a (RED-)ΔΔCq analysis on multiple bulk samples.  
277 However, the size of  $z$ , the residue rate of the S allele, is not identified and must be specified as a fixed  
278 parameter. The amplification efficiency,  $\eta$ , is estimated in theory over the iterative calculation of Eq.  
279 11, but in fact, simultaneous estimation sometimes fails when  $\eta$  is set as unknown.

280 Therefore, the experimenter should identify the sizes of these auxiliary parameters. To estimate  
281 their plausible sizes, one can conduct (RED-)ΔΔCq analysis using DNA solutions with known allele  
282 ratios; for instance, DNA can be extracted from each of the pure breeding lines of S and R and mix the  
283 solutions at multiple ratios, or make a dilution series of R by S. As the ratio of  $X_R$  to  $X_S$  is strictly  
284 fixed, Eq. 7 is directly applicable to express the relationship between DNA quantities and the four Cq

285 measurements. The R functions `knownqpcr()` and `knownqpcr_unpaired()` appearing in the package  
286 provide the maximum likelihood estimation for  $\delta_B$ ,  $\delta_T$ ,  $\sigma_c$ ,  $z$ , and  $\eta$ . These values can be used as fixed  
287 parameters in the `freqqpcr()` function. The “knownqpcr\_unpaired” function was developed to handle  
288 incomplete data (i.e., the observations of  $\tau^{HW}$ ,  $\tau^{TW}$ ,  $\tau^{HD}$ , and  $\tau^{TD}$  have different data lengths). If the  
289 four Cq measures are available for all samples, then “knownqpcr” is used.

290 Another objective of the analysis with known-ratio samples is to test the homoscedasticity of the  
291 qPCR data at the scale of Cq measures. Regarding the relationship between the etoxazole-R allele  
292 frequency in *T. urticae* and the corresponding  $2^{-\Delta\Delta Cq}$  measures (the approximate point estimate of the  
293 frequency), Osakabe *et al.* (2017) demonstrated linearity using a sample series of DNA with multiple  
294 mixing ratios on CHS1 (I1017F). In the next section, we recycled the same data to compare whether  
295 the Cq measurements in the RED- $\Delta\Delta Cq$  analysis obey the homoscedasticity in the scale of  $\Delta\Delta Cq$  or  
296  $(1 + \eta)^{-\Delta\Delta Cq}$ .

## 297 **Materials and Methods**

### 298 *Experiment 1: estimation of auxiliary parameters and verification of* 299 *homoscedasticity in Cq measurements based on mite DNA samples with* 300 *known allele-mixing ratios*

#### 301 *Experimental setup*

302 In the experiment by Osakabe *et al.* (2017), the resistant mite strain (SoOm1-etoR strain) originated  
303 from a field population collected in Omaezaki City, Shizuoka, Japan (34.7°N, 138.1°E) in January  
304 2012. The susceptible strain was obtained from Kyoyu Agri Co., Ltd. (Kanagawa, Japan) (Kyoyu-S  
305 strain). For each strain, two pairs of females and males were used separately. Each pair was allowed to  
306 mate and oviposit on a kidney bean leaf square (2 × 2 cm) for four days. The mites were then  
307 confirmed to be homozygous on the CHS1 locus using sequence analysis. Genomic DNA extracted

308 from the offspring of each pair was used for qPCR analysis. For each pair, the DNA extracts were  
309 prepared twice, each of which was a mixture from 50 adult females homogenized together, that is, four  
310 extracts (replicates) for each strain.

311 To verify the validity of the RED- $\Delta\Delta$ Cq method, qPCR analysis was performed with heterogeneous  
312 DNA solutions with ten mixing ratios of  $X_R/(X_R + X_S) = \{0, 0.001, 0.005, 0.01, 0.05, 0.1, 0.25, 0.5,$   
313  $0.75, 1\}$ . The net DNA concentration of each mixed solution was adjusted to  $1 \text{ ng } \mu\text{l}^{-1}$ , from which 15  
314 ng was dispensed into each of the two tubes. Only one was digested with the restriction enzymes  
315 before qPCR. For digestion, the samples were treated with a mixture of two enzymes, *MluC* I (10  
316 units) and *Taq<sup>q</sup>I* (20 units; New England BioLabs, Ipswich, MA, USA), at  $37^\circ\text{C}$  for 3 h, followed by  
317 incubation at  $65^\circ\text{C}$  for 3 h. This is due to the polymorphism of the *CHS1* loci; the 1017 codon of *T.*  
318 *urticae* displays ATT (Kyoyu-S strain) or TTT (SoOm1-etoR) sequences, whereas the upstream 1016  
319 codon displays a synonymous TCG or TCA independent of the strains (Van Leeuwen et al., 2012).  
320 Therefore, we need to digest both TCGATT (underline shows the restriction site of *Taq<sup>q</sup>I*) and  
321 TCAATT (*MluC* I) to diminish the entire S allele.

322 qPCR analysis using the intercalator method was performed using the LightCycler Nano System  
323 (Roche Diagnostics, Basel, Switzerland) with SYBR Fast qPCR Mix (Takara, Kusatsu, Japan) as  
324 described previously (Osakabe et al., 2017). The primer sets were tu03*CHS1* (forward: 5'-  
325 GGCAGTCTTCATCCACAAG-3' and reverse: 5'-GTGTTCCCAAGTAACAACGTTC-3') and  
326 tu25*GAPDH* (forward: 5'-GCACCAAGTGCTAAAGCATGGAG-3' and reverse: 5'-  
327 GAACTGGAACACGGAAAGCCATAC-3').

### 328 *Statistical analysis*

329 The maximum likelihood of  $\delta_B$ ,  $\delta_T$ ,  $\sigma_c$ ,  $z$ , and  $\eta$  was conducted with the “knownqpcr\_unpaired”  
330 function of the freqpcr package (version 0.3.2). The raw Cq data are available as ESM 1 along with a  
331 step-by-step guide for statistical analyses (ESM 2). Due to the limitation of the handling capacity of  
332 the thermal cycler, qPCR analysis was not conducted on undigested samples of the nine mixing ratios

333 other than  $X_R/(X_R + X_S) = 1$  (i.e., pure R solution). Thus, in each replicate, Osakabe *et al.* (2017)  
334 used the observed  $\Delta\tau^W$  value when the ratio = 1 for other ratios to calculate the conventional  $\Delta\Delta Cq$   
335 indices. As we have shown in Eq. 7, this operation does not affect the point estimates of  $p$ , although  
336 the size of the Cq measurement error ( $\sigma_c$ ) will be underestimated if we recycle the observed Cq value  
337 multiple times.

338 Regarding the relationship between the true mixing ratio and the RED- $\Delta\Delta Cq$  measures in the  
339 sample, the linearity was analyzed using a linear model via the function “lm” running on R version  
340 3.6.1 (R Core Team, 2019), where the response variables were put into the model at the scale of Cq or  
341  $(1 + \eta)^{-\Delta\Delta Cq}$ . Based on the linear models, we tested heteroscedasticity using the Breusch-Pagan test  
342 via the `bptest()` function of the R library “lmtest” (Hothorn et al., 2019).

## 343 ***Experiment 2: evaluation of the simultaneous estimation method with*** 344 ***randomly generated data***

345 Since the experiment by Osakabe *et al.* (2017) used a sample series with strict mixing ratios, the effect  
346 of individual differences in DNA yield was not evaluated. Instead, we conducted a numerical  
347 experiment to verify the accuracy of the simultaneous parameter estimation under uncertainty in the  
348 individual DNA yield. The frequency of the R allele in the population,  $p$ , was set to 0.01, 0.05, 0.1,  
349 0.25, 0.5, or 0.75.

350 For the sampling strategy,  $N$  bulk samples (the parameter ‘ntrap’ in the R source code), each  
351 comprising  $n$  individuals ( $n$  was fixed among the samples: the parameter ‘npertrap’ in the code), were  
352 generated by random sampling from a wild population of a haploid organism. To assess how the  
353 estimation interval responds to the sample sizes, we evaluated the combination of  $N = \{2, 4, 8, 16, 32,$   
354  $64\}$  and  $n = \{4, 8, 16, 32, 64\}$ , though the combinations with  $Nn > 128$  were excluded ( $Nn$   
355 corresponds to ‘ntotal’ in the code). The DNA quantities ( $X_R$  and  $X_S$ ) contained in each bulk sample  
356 were generated as random numbers that followed the gamma distributions of Eq. 3. To cover a



357 plausible variability range of the DNA yield, the gamma shape parameter was varied as  $k = \{1, 3, 9,$   
358  $27\}$ . Depending on the size of  $k$ , the gamma scale parameter was set at  $\theta = 1 \times 10^{-6}/k$  to fix the  
359 mean of the individual DNA yield to  $1 \times 10^{-6}$ . The termination threshold for qPCR,  $X_{\theta}$  was fixed at  
360 1.

361 We fixed the other parameters due to limitations of the computing resources. From the results of  
362 Experiment 1,  $\delta_T = 1.2$ ,  $\delta_B = 0.24$ ,  $z = 0.0016$ , and  $\eta = 0.97$  were presupposed. As for the random  
363 errors in the PCR amplification process and/or the Cq measurement,  $\sigma_c = 0.2$  was assumed regardless  
364 of the initial template quantity. For each of the 624 parameter regions, the dummy datasets comprising  
365  $N$  bulk samples were generated 1,000 times independently with different random number seeds (i.e.,  
366 1,000 replicates), for which the parameter estimation with `freqpcr(..., beta = TRUE)` of the `freqpcr`  
367 package version 0.3.1 was run on the R 3.6.1 environment. The simulation code is available in ESM 3.

368 As we also implemented the gamma distribution model as `freqpcr(..., beta = FALSE)`, a numerical  
369 experiment with the gamma model was also conducted for the first 250 replicates, and the estimation  
370 accuracy was compared between the two assumptions. Furthermore, we also fitted the function with  
371 the settings `freqpcr(..., K = 1)`, that is, assuming the gamma shape parameter was fixed at 1 (a.k.a.  
372 exponential distribution), in addition to the default simulation with all parameters ( $p$ ,  $k$ ,  $\delta_T$ , and  $\sigma_c$ )  
373 unknown. Further, the easiest way to estimate  $p$  derived from Eq. 8, we averaged the observed  $\Delta\Delta Cq$   
374 values for  $N$  bulk samples and transformed them as  $\hat{p} = (1 + \eta)^{-\overline{\Delta\Delta Cq}}$ .

## 375 **Results**

### 376 *Estimation of auxiliary parameters and verification of homoscedasticity*

377 Based on the Cq measures, the auxiliary parameters were estimated based on the RED- $\Delta\Delta Cq$  analysis  
378 of the I1017F mutation of *T. urticae*. As for the initial quantity of template DNA (the parameter  
379 “meanDNA” on the R code; defined as  $X/X_{\theta}$ ), the maximum-likelihood estimate was  $1.256 \times 10^{-6}$   
380 (95% confidence interval:  $7.722 \times 10^{-7}$  to  $2.041 \times 10^{-6}$ ). The relative quantity of the target gene to the

381 housekeeping gene  $\delta_T$  (targetScale) was estimated to be 1.170 (95% CI: 1.069–1.280). The locus-  
382 independent change rate in the template quantity accompanying the restriction enzyme treatment  $\delta_B$   
383 (baseChange) was 0.2361 (95% CI: 0.2040 to 0.2731). The measurement error in the scale of Cq  $\sigma_c$   
384 (SD) was 0.2376 (95% CI: 0.2050 to 0.2755). The residue rate of the S allele after digestion  $z$   
385 (zeroAmount) was 0.001564 (95% CI: 0.001197–0.002044). The efficiency of amplification per PCR  
386 cycle  $\eta$  (EPCR) was 0.9712 (95% CI: 0.9231–1.022).

387 In the RED- $\Delta\Delta Cq$  analysis of the etoxazole resistance of *T. urticae*, the relationship between the  
388 true R allele frequency ( $Y_R = X_R/(X_R + X_S)$  in the sample) and the corresponding Cq measures  
389 exhibited higher homoscedasticity in the scale of the measured  $\Delta\Delta Cq$  values rather than in  
390  $(1 + \eta)^{-\Delta\Delta Cq}$ , the transformation to  $\hat{Y}_R$  (Figure 2). The linear regression of the  $\Delta\Delta Cq$  values on  
391  $-\ln[0.001564 \times (1 - Y_R) + Y_R]/\ln(1 + 0.971)$  showed high linearity (intercept =  $-0.07694$ ,  
392 coefficient = 1.025, adjusted  $R^2 = 0.9936$ ). The homoscedasticity of the coefficient of determination  
393 was not rejected at the 5% level of significance (Breusch-Pagan test: BP = 3.1577,  $df = 1$ ,  $p = 0.07557$ )  
394 (Figure 2A). Conversely, the linear regression of  $1.971^{-\Delta\Delta Cq}$  on  $[0.001564 \times (1 - Y_R) + Y_R]$  showed a  
395 slightly lower linearity (intercept =  $-0.008625$ , coefficient = 1.092, adjusted  $R^2 = 0.9709$ ). The  
396 Breusch-Pagan test was highly significant (BP = 13.978,  $df = 1$ ,  $p = 0.0001849$ ), rejecting  
397 homoscedasticity (Figure 2B). These results suggest that it is easier to model the error structure of the  
398 RED- $\Delta\Delta Cq$  method on the scale of Cq values (logarithm) rather than frequency (linear scale).

### 399 ***Evaluation of the simultaneous estimation method with randomly generated*** 400 ***data***

401 Among the 624 parameter regions of the numerical simulation with 1,000 replicates (250 for the  
402 gamma model), the total success rate of the interval estimation  $p$  using `frequpcr(..., beta = TRUE)` was  
403 70.6% and 94.5% when all parameters were unknown, and when the gamma shape parameter was  
404 fixed as  $k = 1$ . The “success rate” here indicates the probability when the function returns certain  
405 values other than NA (i.e., the diagonal of the Hessian was not negative): no guarantee that the

406 estimated confidence interval was accurate. The estimation success for the  $C_q$  measurement error,  $\sigma_c$ ,  
407 was 69.6% and 97.6% in the beta-distribution model with unknown  $k$  and  $k = 1$ , respectively. The  
408 relative quantity of the target gene,  $\delta_T$ , was 68.1% and 96.1%, respectively. However, the estimated  
409 success of  $k$  was 59.9% with the beta distribution model, showing a lower performance than the other  
410 parameters, implying that the likelihood is insensitive to the size of  $k$ . Conversely, the estimation of  $p$   
411 is robust to the size of  $k$ , as we show later in this section.

412 The estimation success of `freqpqr()` largely depended on the total sample size ( $Nn$  corresponding to  
413 the facet 'ntotal' in the figures), as well as the level of  $p$  (Figure S1 and S2 for the beta and gamma  
414 models, with all parameters unknown). In each parameter region, the quantity  $\text{Bin}(0|Nn, p)$  generally  
415 gives the probability that the whole sample contains no R individuals. When  $Nn$  is larger enough,  
416  $Nn > 3/p$  is approximately the requirement for the total sample size to contain at least one R  
417 individual with 95% confidence, called the “rule of three” (Eypasch et al., 1995). The gray  
418 backgrounds in the facets of Figures 3–4 and S1–S7 signify the regions where the total sample sizes  
419 are smaller than the thresholds (e.g., 60 haploid individuals are required when  $p = 0.05$ ). As shown in  
420 Figures S1 and S2, the parameter estimation often failed when  $Nn$  did not meet the rule of three. Once  
421 we exclude the parameter regions of  $Nn \leq 3/p$ , the estimation success rate of  $p$  with `freqpqr(..., beta`  
422 `= TRUE)` improved to 84.3% and 99.9% with all parameters unknown and assuming  $k = 1$ ,  
423 respectively.

424 As for the estimation accuracy of  $p$ , the `freqpqr()` function assuming beta distribution provides an  
425 unbiased estimator. Figures 3 and S3 show the estimated sizes of  $p$  using the beta model with all  
426 parameters unknown and assuming  $k = 1$ , respectively. Both settings demonstrated that the estimator  
427 converged to the true R frequency; the upper/lower bounds of the estimated 95% confidence intervals  
428 (yellow/blue boxes in each plot) became narrower as we increased the total sample sizes ( $Nn$ ) or  
429 included more bulk DNA samples ( $N$ ). Fixing the size of the gamma shape parameter to  $k = 1$  scarcely  
430 affected the point estimates and intervals of  $p$ , as long as  $Nn > 3/p$  is satisfied (Figure S3). However,

431 if every individual was analyzed separately, the interval estimation was only possible when  $k$  was  
432 fixed (see the regions of “sample division = ntotal” cases in Figure 3).

433 When we used the gamma distribution model, the interval estimation of  $p$  was also possible and  
434 unbiased (Figure S4). However, when we defined the point estimator of  $p$  as a simple average, that is,  
435  $\hat{p} = (1 + \eta)^{-\overline{\Delta\Delta\tau}}$ , it was strongly underestimated as the samples were more divided ( $N/Nn$  was  
436 large) (Figure 4). The upper limit of 95% CI often violated 1, suggesting that the “simple average of  
437  $\Delta\Delta Cq$ ”  $\pm 1.96$  SE is inadequate for the interval estimation based on the RED- $\Delta\Delta Cq$  method.

438 Although the `freqpcr()` function with the gamma and beta distributions both showed an unbiased  
439 estimation of  $p$ , the gamma model was disadvantageous regarding calculation time and the number of  
440 iterations before convergence. The time varied largely in the model settings and sample sizes (Figures  
441 S5–S7). Among the settings we tried, beta model with fixed  $k$  was the fastest and converged within a  
442 few seconds in most parameter regions (median and 75 percentile: 0.32 and 0.69 s: Figure S6). It was  
443 three and >10 times faster than the beta (0.91 and 2.4 seconds: Figure S5) and gamma (3.0 and 15 s:  
444 Figure S7) model, respectively with all parameters unknown. The calculation time increased as the  
445 dataset size increased -  $Nn$  and the sample was more divided (larger  $N/Nn$ ) in the beta distribution  
446 model, because the marginal likelihood was calculated for each bulk sample (Figures S5 and S6).  
447 Conversely, the gamma distribution model (Figure S7) requires increased calculation time as the size  
448 of each bulk sample becomes larger (larger  $n_h$ ). This was considered because the combination of  
449  $\text{Bin}(m_h|n_h, p)$  exploded when  $n_h$  was large.

450 Furthermore, the estimation accuracy of the shape parameter,  $k$ , it was underestimated as the real  
451 size of the parameter increased (e.g.,  $k = 27$ ) when the gamma distribution model was applied (Figure  
452 S8B). Since the iterative fitting of the parameter in `freqpcr()` always starts internally from  $k = 1$  (this  
453 was determined due to the calculation stability), this bias suggests the likelihood function of  $\psi_G$  (Eq.  
454 13), with little information on the size of  $k$  compared with  $p$ . Then,  $k$  tends to stay at its initial value,  
455 suggesting that the gamma model is not suitable for the simultaneous estimation of  $p$  and  $k$ . Unlike the  
456 gamma version, the fitting of  $k$  with `freqpcr(beta = TRUE)` was satisfactory when we divided the total

457 samples into more bulk samples (larger  $N/Nn$ ), although the initial value dependence was still  
458 observed, especially when  $p$  or  $N$  was small (Figure S8A). This may be because the estimation of  $k$  via  
459  $\text{Beta}(m_h k, (n_h - m_h)k)$  in Eq. 14 is comparable with measuring the overdispersion of  $Y_{h,R|m_h}$ , which  
460 is only possible when multiple bulk samples contain both R and S alleles.

## 461 Discussion

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

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

480 The combination of mass trapping and bulk qPCR analysis solves the latter challenge by collecting  
481 more individuals and pooling them. This can result in higher precision with less work than individual  
482 PCR. For instance, we sampled 16 individuals from the population with an allele frequency of  $p = 0.05$   
483 and analyzed two individuals once in the numerical experiment (Figure 3: facet of  $n_{\text{total}} = 16$ , sample  
484 division = 8). The lower and upper limits of the 95% confidence interval  $p$  were estimated to be  
485 0.0087 and 0.34, respectively, using `freqpqr(..., beta = TRUE)` (as the medians of the 1,000  
486 independent trials). We also simulated the case of  $n_{\text{total}} = 64$  and sample division = 4 (i.e., analyzed  
487 16 individuals together) and found the upper and lower limits to be 0.015 and 0.15, respectively. Thus,  
488 we improved the precision of the interval estimate with half the handling effort.

489 Also, in non-quantitative PCR, sample pooling is considered as a tool for the detection of rare  
490 (c)DNA in the population with practical labor requirements, and has been used as high throughput pre-  
491 screening system for many samples e.g. in clinical examinations (Taylor et al., 2010; Yelin et al.,  
492 2020). However, in some fields, such as plant quarantine, it is important to guarantee that a product is  
493 not contaminated with pests or unapproved genetically modified seeds at a certain consumer risk. As  
494 the assumed frequency range is low ( $p \approx 0.001$ ), frequency estimation is not realistic (3,000 seeds are  
495 needed to meet the “rule of three” when  $p = 0.001$ ) and is not required for the current inspection  
496 routine. Thus, group testing based on non-quantitative PCR has been conducted in these fields  
497 (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  
499 “positive.”

500 Overall, there has been a gap in methodology between the frequency estimation based on the  
501 individual PCR and the non- or semi-quantitative PCR based on the non-quantitative bulk PCR.  
502 Although it provides the highest estimation precision following binomial distribution, the former is  
503 only available at a higher  $p$ ; it becomes labor-intensive once we try to quantify rare alleles. The latter  
504 can be applied to a lower range of  $p$ , but the precision is generally low or even non-quantitative.  
505 Bridging the gap, our qPCR-based procedure offers an allele frequency estimation in the mid-low

506 range ( $p = 0.01$  to  $0.25$ ), which is considered a critical range for decision making in some fields like  
507 pesticide resistance management (Takahashi et al., 2017; Sudo et al., 2018).

508 Although this study focused on resistance genes, the likelihood model in Eq. 11 can also be applied  
509 for other qPCR protocols based on  $\Delta\Delta Cq$ . If both the specific and nonspecific primer sets are available  
510 to amplify the “mutant” and “wild type + mutant” alleles at the target locus, they can be used for the  
511 test and control samples equivalent to  $X^{TV}$  in Eq. 10 and  $X^{TU}$  in Eq. 9, respectively. However, there is a  
512 caveat in determining which allele should be amplified with a specific primer set and which affects the  
513 estimation accuracy due to the intrinsic nature of  $(1 + \eta)^{-\Delta\Delta\tau}$ . As shown, the 95% confidence  
514 intervals were broader when  $p = 0.75$  than when  $p = 0.25$  (Figure 3), the accuracy was not symmetric  
515 around 0.5, but more accurate when the frequency was low. That is, one should design a specific  
516 primer set to amplify the allele that would be rare in the population to improve the signal-to-noise  
517 ratio.

518 The maximum likelihood estimation with `freqpqr()` relies on the assumption that the quantities of  
519 the S and R alleles in each bulk sample independently follow gamma distribution and that their  
520 quotient is expressed using beta distribution. Fixing the size of the gamma shape parameter  $k$  further  
521 accelerated the optimization, which was owing to the robustness of  $p$  to the size of  $k$ . However, once  
522 the size of  $k$  was fixed much larger than the actual size of the gamma shape parameter (i.e., the  
523 individual DNA yield was regarded as almost a fixed value), the iterative optimization using the `nlm()`  
524 function sometimes returned an error. Therefore, one should start with a smaller shape parameter e.g.,  
525  $k = 1$  (the exponential distribution: Figure S3), which is currently the default setting of the `freqpqr`  
526 package.

527 In qPCR applications for diagnostic use,  $\Delta\Delta Cq$  is often used with calibration. One of the popular  
528 methods is the involvement of technical replicates; each sample is dispensed and analyzed using qPCR  
529 multiple times, which negates the  $Cq$  measurement error. The measurement error obeys a  
530 homoscedastic normal distribution in the  $Cq$  scale, as shown in Experiment 1. Thus, a simple solution  
531 is to average the  $Cq$  values measured for each bulk sample before the estimation with `freqpqr()`,

532 although the estimated size of  $\sigma_c$  changes from its original definition in Eq. 7. However, it is trivial if  
533 the number of technical replicates is unified between bulk samples.

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

## 541 **Acknowledgments**

542 We appreciate Dr. Kohji Yamamura and Dr. Takehiko Yamanaka for earlier discussion on the gamma  
543 assumption of the individual DNA yield. The work was supported by a grant from the Ministry of  
544 Agriculture, Forestry, and Fisheries of Japan (Genomics-based Technology for Agricultural  
545 Improvement): PRM05 to M.O. and PRM07 to M.S.

## 546 **References**

- 547 Andow, D. A., & Alstad, D. N. (1998). F2 screen for rare resistance alleles. *Journal of Economic*  
548 *Entomology*, 91(3), 572–578.
- 549 Chen, H., Rangasamy, M., Tan, S. Y., Wang, H., & Siegfried, B. D. (2010). Evaluation of five  
550 methods for total DNA extraction from western corn rootworm beetles. *PLoS One*, 5(8), e11963.
- 551 Donnelly, M. J., Isaacs, A. T., & Weetman, D. (2016). Identification, validation, and application of  
552 molecular diagnostics for insecticide resistance in malaria vectors. *Trends in Parasitology*, 32(3),  
553 197–206.
- 554 Eypasch, E., Lefering, R., Kum, C. K., & Troidl, H. (1995). Probability of adverse events that have not  
555 yet occurred: A statistical reminder. *BMJ*, 311(7005), 619–620.
- 556 Falconer, D. S. (1960). *Introduction to quantitative genetics*. Oliver And Boyd; Edinburgh; London.



- 557 ffrench-Constant, R. H. (2013). The molecular genetics of insecticide resistance. *Genetics*, 194(4),  
558 807–815.
- 559 Germer, S., Holland, M. J., & Higuchi, R. (2000). High-throughput SNP allele-frequency  
560 determination in pooled DNA samples by kinetic PCR. *Genome Research*, 10(2), 258–266.
- 561 Gould, F., Anderson, A., Jones, A., Sumerford, D., Heckel, D. G., Lopez, J., Micinski, S., Leonard, R.,  
562 & Laster, M. (1997). Initial frequency of alleles for resistance to *Bacillus thuringiensis* toxins in  
563 field populations of *Heliothis virescens*. *Proceedings of the National Academy of Sciences*, 94(8),  
564 3519–3523.
- 565 Gouws, E. J., Gaston, K. J., & Chown, S. L. (2011). Intraspecific body size frequency distributions of  
566 insects. *PLoS One*, 6(3), e16606.
- 567 Hothorn, T., Zeileis, A., Farebrother (pan.f), R. W., Cummins (pan.f), C., Millo, G., & Mitchell, D.  
568 (2019). lmtest: Testing Linear Regression Models (0.9-37) [Computer software]. [https://CRAN.R-](https://CRAN.R-project.org/package=lmtest)  
569 [project.org/package=lmtest](https://CRAN.R-project.org/package=lmtest)
- 570 Kim, S. Y., Lohmueller, K. E., Albrechtsen, A., Li, Y., Korneliussen, T., Tian, G., Grarup, N., Jiang,  
571 T., Andersen, G., & Witte, D. (2011). Estimation of allele frequency and association mapping using  
572 next-generation sequencing data. *BMC Bioinformatics*, 12(1), 231.
- 573 Knapp, M. (2016). Relative importance of sex, pre-starvation body mass and structural body size in  
574 the determination of exceptional starvation resistance of *Anchomenus dorsalis* (Coleoptera:  
575 Carabidae). *PloS One*, 11(3), e0151459.
- 576 Kundu, D., & Manglick, A. (2005). Discriminating between the log-normal and gamma distributions.  
577 *Journal of the Applied Statistical Sciences*, 14, 175–187.
- 578 Kwon, D. H., Yoon, K. S., Strycharz, J. P., Clark, J. M., & Lee, S. H. (2008). Determination of  
579 permethrin resistance allele frequency of human head louse populations by quantitative sequencing.  
580 *Journal of Medical Entomology*, 45(5), 912–920.
- 581 Li, G., Reising, D., Miao, J., Gould, F., Huang, F., & Feng, H. (2016). Frequency of Cry1F non-  
582 recessive resistance alleles in North Carolina field populations of *Spodoptera frugiperda*  
583 (Lepidoptera: Noctuidae). *PloS One*, 11(4), e0154492.
- 584 Livak, K. J. (1997). Comparative Ct method. In *ABI Prism 7700 sequence detection system User*  
585 *Bulletin* #2, P/N 4303859 (Vol. 2, pp. 11–15). Applied Biosystems.  
586 [http://tools.thermofisher.com/content/sfs/manuals/cms\\_040980.pdf](http://tools.thermofisher.com/content/sfs/manuals/cms_040980.pdf)
- 587 Livak, Kenneth J., & Schmittgen, T. D. (2001). Analysis of relative gene expression data using real-  
588 time quantitative PCR and the 2-  $\Delta\Delta$ CT method. *Methods*, 25(4), 402–408.
- 589 Maeoka, A., Yuan, L., Itoh, Y., Saito, C., Doi, M., Imamura, T., Yamaguchi, T., Imura, T., &  
590 Osakabe, M. (2020). Diagnostic prediction of acaricide resistance gene frequency using quantitative  
591 real-time PCR with resistance allele-specific primers in the two-spotted spider mite *Tetranychus*  
592 *urticae* population (Acari: Tetranychidae). *Applied Entomology and Zoology*, 55, 329–335.  
593 <https://doi.org/10.1007/s13355-020-00686-7>

- 594 May, R. M. (1976). Models for single populations. In *Theoretical ecology* (pp. 5–29). Blackwell.
- 595 Narasimhan, B., Koller, M., Johnson, S. G., Hahn, T., Bouvier, A., Kiêu, K., & Gaure, S. (2019).  
596 cubature: Adaptive Multivariate Integration over Hypercubes (2.0.4) [Computer software].  
597 <https://CRAN.R-project.org/package=cubature>
- 598 Osakabe, M., Imamura, T., Nakano, R., Kamikawa, S., Tadatsu, M., Kunimoto, Y., & Doi, M. (2017).  
599 Combination of restriction endonuclease digestion with the  $\Delta\Delta\text{Ct}$  method in real-time PCR to  
600 monitor etoxazole resistance allele frequency in the two-spotted spider mite. *Pesticide Biochemistry*  
601 *and Physiology*, 139, 1–8.
- 602 R Core Team. (2019). R version 3.6.1. <https://www.r-project.org/>
- 603 Rakovski, C., Weisenberger, D. J., Marjoram, P., Laird, P. W., & Siegmund, K. D. (2011). Modeling  
604 measurement error in tumor characterization studies. *BMC Bioinformatics*, 12(1), 284.
- 605 Samayoa, L. F., Malvar, R. A., Olukolu, B. A., Holland, J. B., & Butrón, A. (2015). Genome-wide  
606 association study reveals a set of genes associated with resistance to the Mediterranean corn borer  
607 (*Sesamia nonagrioides* L.) in a maize diversity panel. *BMC Plant Biology*, 15(1), 35.
- 608 Snoeck, S., Kurlavs, A. H., Bajda, S., Feyereisen, R., Greenhalgh, R., Villacis-Perez, E., Kosterlitz,  
609 O., Dermauw, W., Clark, R. M., & Van Leeuwen, T. (2019). High-resolution QTL mapping in  
610 *Tetranychus urticae* reveals acaricide-specific responses and common target-site resistance after  
611 selection by different METI-I acaricides. *Insect Biochemistry and Molecular Biology*, 110, 19–33.
- 612 Sonoda, S., Inukai, K., Kitabayashi, S., Kuwazaki, S., & Jouraku, A. (2017). Molecular evaluation of  
613 diamide resistance in diamondback moth (Lepidoptera: Yponomeutidae) populations using  
614 quantitative sequencing. *Applied Entomology and Zoology*, 52(2), 353–357.
- 615 Sudo, M., Takahashi, D., Andow, D. A., Suzuki, Y., & Yamanaka, T. (2018). Optimal management  
616 strategy of insecticide resistance under various insect life histories: Heterogeneous timing of  
617 selection and interpatch dispersal. *Evolutionary Applications*, 11(2), 271–283.
- 618 Sudo, M., Yamamura, K., Sonoda, S., & Yamanaka, T. (in press). Estimating the proportion of  
619 resistance alleles from bulk Sanger sequencing, circumventing the variability of individual DNA.  
620 *Journal of Pesticide Science*.
- 621 Sugimoto, N., Takahashi, A., Ihara, R., Itoh, Y., Jouraku, A., Van Leeuwen, T., & Osakabe, M.  
622 (2020). QTL mapping using microsatellite linkage reveals target-site mutations associated with high  
623 levels of resistance against three mitochondrial complex II inhibitors in *Tetranychus urticae*. *Insect*  
624 *Biochemistry and Molecular Biology*, 103410.
- 625 Tabashnik, B. E., Patin, A. L., Dennehy, T. J., Liu, Y.-B., Carriere, Y., Sims, M. A., & Antilla, L.  
626 (2000). Frequency of resistance to *Bacillus thuringiensis* in field populations of pink bollworm.  
627 *Proceedings of the National Academy of Sciences*, 97(24), 12980–12984.
- 628 Takahashi, D., Yamanaka, T., Sudo, M., & Andow, D. A. (2017). Is a larger refuge always better?  
629 Dispersal and dose in pesticide resistance evolution. *Evolution*, 71(6), 1494–1503.

- 630 Taylor, S. M., Juliano, J. J., Trottman, P. A., Griffin, J. B., Landis, S. H., Kitsa, P., Tshetu, A. K., &  
631 Meshnick, S. R. (2010). High-throughput pooling and real-time PCR-based strategy for malaria  
632 detection. *Journal of Clinical Microbiology*, 48(2), 512–519.
- 633 Toda, S., Hirata, K., Yamamoto, A., & Matsuura, A. (2017). Molecular diagnostics of the R81T  
634 mutation on the D-loop region of the  $\beta 1$  subunit of the nicotinic acetylcholine receptor gene  
635 conferring resistance to neonicotinoids in the cotton aphid, *Aphis gossypii* (Hemiptera: Aphididae).  
636 *Applied Entomology and Zoology*, 52(1), 147–151.
- 637 Uesugi, R., Hinomoto, N., & Goto, C. (2016). Estimated time frame for successful PCR analysis of  
638 diamondback moths, *Plutella xylostella* (Lepidoptera: Plutellidae), collected from sticky traps in  
639 field conditions. *Applied Entomology and Zoology*, 51(3), 505–510.
- 640 Van Leeuwen, T., Demaeght, P., Osborne, E. J., Dermauw, W., Gohlke, S., Nauen, R., Grbić, M.,  
641 Tirry, L., Merzendorfer, H., & Clark, R. M. (2012). Population bulk segregant mapping uncovers  
642 resistance mutations and the mode of action of a chitin synthesis inhibitor in arthropods.  
643 *Proceedings of the National Academy of Sciences*, 109(12), 4407–4412.
- 644 Van Leeuwen, T., Vanholme, B., Van Pottelberge, S., Van Nieuwenhuyse, P., Nauen, R., Tirry, L., &  
645 Denholm, I. (2008). Mitochondrial heteroplasmy and the evolution of insecticide resistance: Non-  
646 Mendelian inheritance in action. *Proceedings of the National Academy of Sciences*, 105(16), 5980–  
647 5985.
- 648 Van Leeuwen, T., Vontas, J., Tsagkarakou, A., Dermauw, W., & Tirry, L. (2010). Acaricide resistance  
649 mechanisms in the two-spotted spider mite *Tetranychus urticae* and other important Acari: A  
650 review. *Insect Biochemistry and Molecular Biology*, 40(8), 563–572.
- 651 Vandesompele, J., De Preter, K., Pattyn, F., Poppe, B., Van Roy, N., De Paepe, A., & Speleman, F.  
652 (2002). Accurate normalization of real-time quantitative RT-PCR data by geometric averaging of  
653 multiple internal control genes. *Genome Biology*, 3(7), research0034-1.
- 654 Wiens, B. L. (1999). When log-normal and gamma models give different results: A case study. *The*  
655 *American Statistician*, 53(2), 89–93.
- 656 Yamamura, K., & Hino, A. (2007). Estimation of the proportion of defective units by using group  
657 testing under the existence of a threshold of detection. *Communications in Statistics—Simulation*  
658 *and Computation*, 36(5), 949–957.
- 659 Yamamura, K., Mano, J., & Shibaike, H. (2019). Optimal definition of the limit of detection (LOD) in  
660 detecting genetically modified grains from heterogeneous grain lots. *Quality Technology &*  
661 *Quantitative Management*, 16(1), 36–53.
- 662 Yelin, I., Aharony, N., Shaer-Tamar, E., Argoetti, A., Messer, E., Berenbaum, D., Shafran, E., Kuzli,  
663 A., Gandali, N., & Hashimshony, T. (2020). Evaluation of COVID-19 RT-qPCR test in multi-  
664 sample pools. *Clinical Infectious Diseases*, 71, 2073–2078.

665  
666

## 667 **Data accessibility**

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

672 ESM 1

673 RED- $\Delta\Delta Cq$  dataset from Osakabe *et al.* (2017).

674 ESM 2

675 R source code for Experiment 1 (Figure 2), including a brief guide to the “freqpcr” package.

676 ESM 3

677 R source code for the numerical simulation (Experiment 2) and the codes for Figures 3 and after.

678 Appendix S1

679 Formularization in the case of diploidy.

## 680 **Author contributions**

681 M.S. designed research, made statistical models and R package, and analyzed data. M.O. conducted  
682 laboratory work. Both authors wrote the final manuscript.

683

684 **Tables**

685 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 $X_R/(X_R + X_S)$	$0 \leq Y_R \leq 1$	—
$k, \theta$	Shape and scale parameters of the gamma distribution $Ga(k, \theta)$	$k > 0, \theta > 0$	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			
$\eta$	Per-cycle efficiency in the PCR amplification (as $1 + \eta$ )	$\eta > 0$	EPCR
$X_\theta$	The termination threshold of the amplification in the real-time PCR process	$X_\theta > 0$	Fixed 1 in the package
$\tau$	Cq value: the number of PCR amplification cycles before termination	$\tau \in \mathbb{R}$	$\tau_h^{TW}$ : target0, $\tau_h^{TD}$ : target1, $\tau_h^{HW}$ : housek0, $\tau_h^{HD}$ : housek1
$\delta_T$	Relative content of the target gene to the internal reference (housekeeping gene)	$\delta_T > 0$	targetScale
$\delta_B$	(In RED- $\Delta\Delta Cq$ method) the locus-independent change rate of the template DNA quantity accompanying the restriction enzyme treatment.	$\delta_T > 0$	baseChange
$z$	(In RED- $\Delta\Delta Cq$ method) residual rate of restriction enzyme digestion, or (in general $\Delta\Delta Cq$ analyses) portion of the off-target allele amplified in the PCR	$z > 0$	zeroAmount
$\varepsilon_c$	Cq measurement error (standard deviation)	$\varepsilon_c > 0$	sdMeasure

686

## 687 Figures

688 Figure 1 Scheme of population allele frequency estimation based on qPCR analyses. A: Insect  
689 sampling and subsequent qPCR analysis using the restriction enzyme digestion (RED)- $\Delta\Delta Cq$  method.  
690 B: Probability distributions of the DNA amounts of the resistant (R) and susceptible (S) alleles and  
691 their ratio in the bulk sample. C: Template DNA involved in the RED- $\Delta\Delta Cq$  analysis and a general  
692  $\Delta\Delta Cq$  analysis using an R-specific primer set. In either method, the frequency of  $X_R$  in a test sample is  
693 quantified as  $X_R + zX_S (\cong X_R)$  measured on the target gene, divided by  $X_R + X_S$  measured on a  
694 housekeeping gene in the sample. As the copy numbers may differ between genes, the relative content  
695  $\delta_T$  is also quantified using a control sample.

696

697 Figure 2 Relationship between the allele frequency in the sample and A: the RED- $\Delta\Delta Cq$  measures, B:  
698 the observed frequency calculated as  $(1 + \eta)^{-\Delta\Delta Cq}$ , showing the results of etoxazole resistance in  
699 the two-spotted spider mites. The lines are not the regression on the actual Cq measurement (shown as  
700 points), but the theoretical relationship between true frequency of the R allele and the quantity defined  
701 as A:  $-\ln(z + Y_R(1 - z))/\ln(1 + \eta)$  or B:  $z + Y_R(1 - z)$ , where  $Y_R = X_R/(X_R + X_S)$ . Parameters are  
702  $z = 0.00156$  and  $\eta = 0.971$ .

703

704 Figure 3 Estimation accuracy of the population allele frequency,  $p$ , with `freqpcr()` when the beta  
705 distribution was assumed, and all estimable parameters (P, K, targetScale, and sdMeasure) were set as  
706 unknown. The result of numerical experiments based on 1,000 dummy datasets per parameter region.  
707 The x-axes corresponds to  $N$ , or the “ntrap” parameter, the extent to which the collected individuals  
708 (ntotal) were divided to the bulk samples. The three box plots (white thin, blue, and yellow wide) in  
709 each region show the maximum likelihood estimates (MLE), lower bound of the 95% CI, and the  
710 upper bound, respectively. In each boxplot, the horizontal line signifies the median of the simulations,  
711 hinges of the box show 25 and 75 percentiles, and the upper/lower whiskers correspond to the  $1.5 \times$   
712 interquartile ranges. The shaded facets show that the total sample sizes (ntotal) are smaller than  $3/p$ .

713

714 Figure 4 Estimation accuracy of the population allele frequency by simple averaging of  $\Delta\Delta Cq$   
715 measures. The frequency was underestimated than its true value (horizontal broken line in each facet)  
716 as the samples were more divided. The Cq dataset was derived from the numerical experiment of “beta  
717 distribution, all parameters unknown.”

# 1 freqPCR: estimation of population allele frequency using qPCR

## 2 $\Delta\Delta Cq$ measures from bulk samples

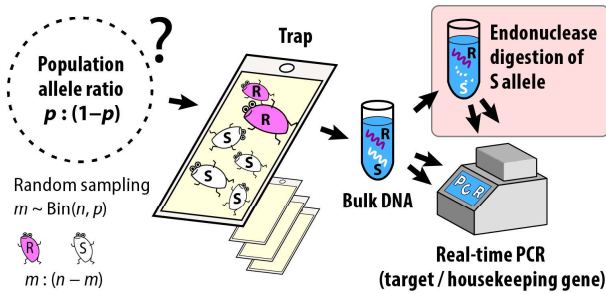
3 Masaaki Sudo, Masahiro Osakabe

4

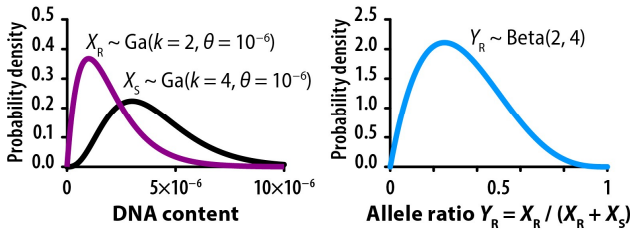
### 5 Figures

6

#### A. Allele frequency estimation by RED- $\Delta\Delta Cq$ method



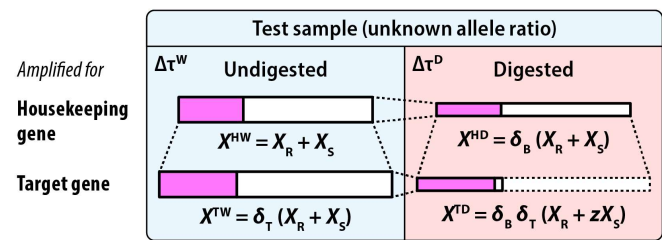
#### B. DNA content distribution in bulk DNA solution



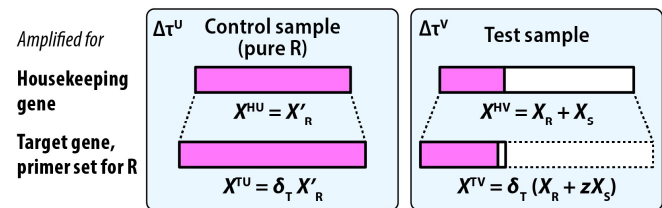
7

#### C. Template DNA quantities in $\Delta\Delta Cq$ methods

(in RED- $\Delta\Delta Cq$  method:  $\Delta\Delta\tau = \Delta\tau^D - \Delta\tau^W$ )



(in general  $\Delta\Delta Cq$  method:  $\Delta\Delta\tau = \Delta\tau^V - \Delta\tau^U$ )



8 Figure 1 Scheme of population allele frequency estimation based on qPCR analyses. A: Insect sampling and

9 subsequent qPCR analysis using the restriction enzyme digestion (RED)- $\Delta\Delta Cq$  method. B: Probability

10 distributions of the DNA amounts of the resistant (R) and susceptible (S) alleles and their ratio in the bulk

11 sample. C: Template DNA involved in the RED- $\Delta\Delta Cq$  analysis and a general  $\Delta\Delta Cq$  analysis using an R-

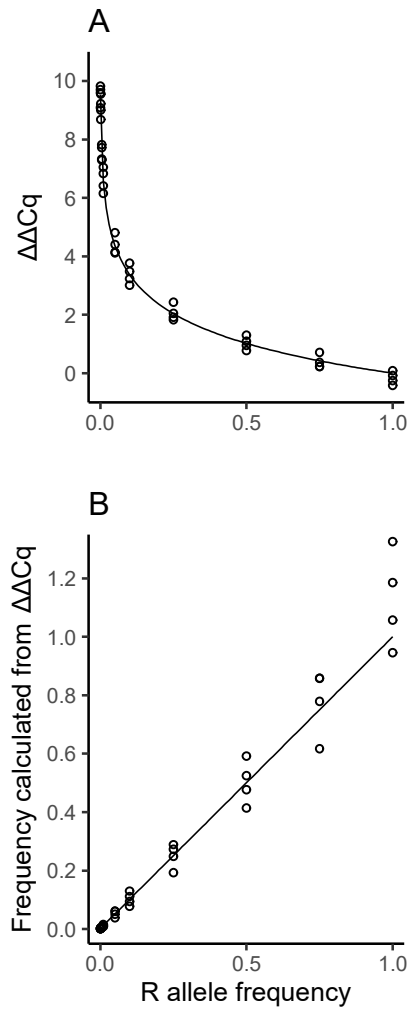
12 specific primer set. In either method, the frequency of  $X_R$  in a test sample is quantified as  $X_R + zX_S (\cong X_R)$

13 measured on the target gene, divided by  $X_R + X_S$  measured on a housekeeping gene in the sample. As the

14 copy numbers may differ between genes, the relative content  $\delta_T$  is also quantified using a control sample.

15



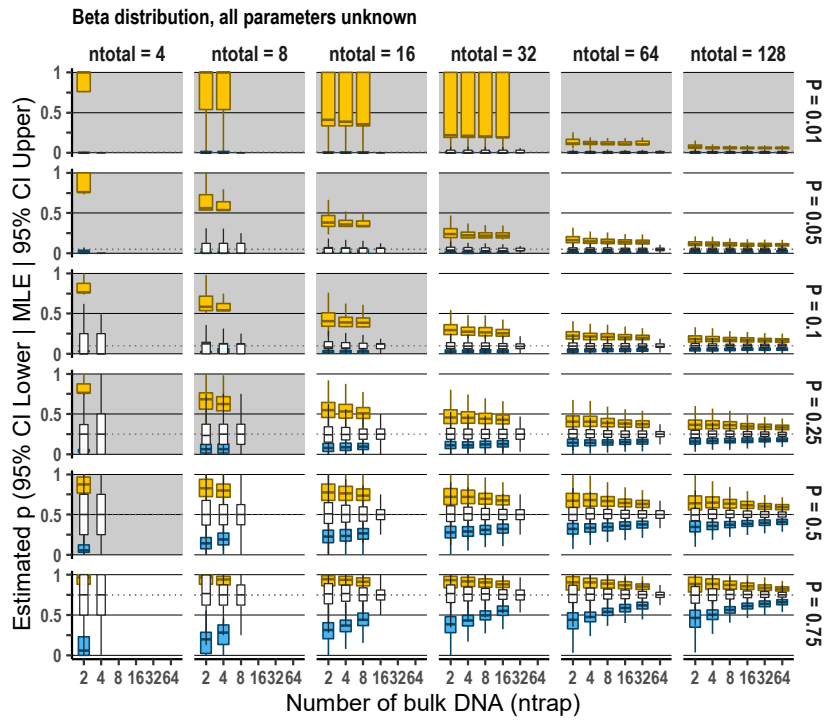


16

17 Figure 2 Relationship between the allele frequency in the sample and A: the RED- $\Delta\Delta Cq$  measures, B: the  
18 observed frequency calculated as  $(1 + \eta)^{-\Delta\Delta Cq}$ , showing the results of etoxazole resistance in the two-  
19 spotted spider mites. The lines are not the regression on the actual Cq measurement (shown as points), but  
20 the theoretical relationship between true frequency of the R allele and the quantity defined as A:  
21  $-\ln(z + Y_R(1 - z))/\ln(1 + \eta)$  or B:  $z + Y_R(1 - z)$ , where  $Y_R = X_R/(X_R + X_S)$ . Parameters are  $z = 0.00156$   
22 and  $\eta = 0.971$ .

23

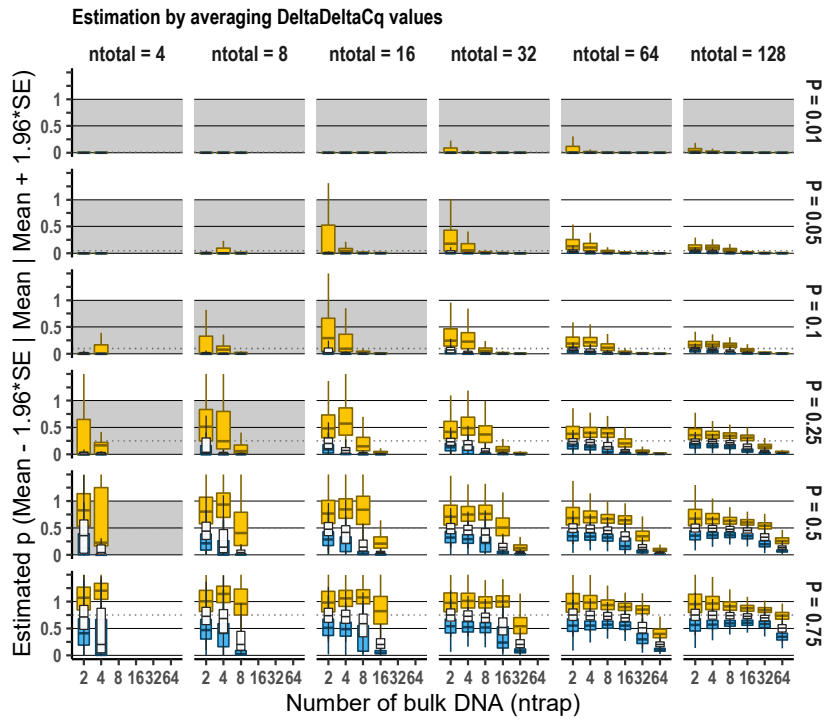




24

25 Figure 3 Estimation accuracy of the population allele frequency,  $p$ , with `freqpcr()` when the beta distribution  
 26 was assumed, and all estimable parameters ( $P$ ,  $K$ , `targetScale`, and `sdMeasure`) were set as unknown. The  
 27 result of numerical experiments based on 1,000 dummy datasets per parameter region. The x-axes  
 28 corresponds to  $N$ , or the “`ntrap`” parameter, the extent to which the collected individuals (`ntotal`) were  
 29 divided to the bulk samples. The three box plots (white thin, blue, and yellow wide) in each region show the  
 30 maximum likelihood estimates (MLE), lower bound of the 95% CI, and the upper bound, respectively. In  
 31 each boxplot, the horizontal line signifies the median of the simulations, hinges of the box show 25 and 75  
 32 percentiles, and the upper/lower whiskers correspond to the  $1.5 \times$  interquartile ranges. The shaded facets  
 33 show that the total sample sizes (`ntotal`) are smaller than  $3/p$ .

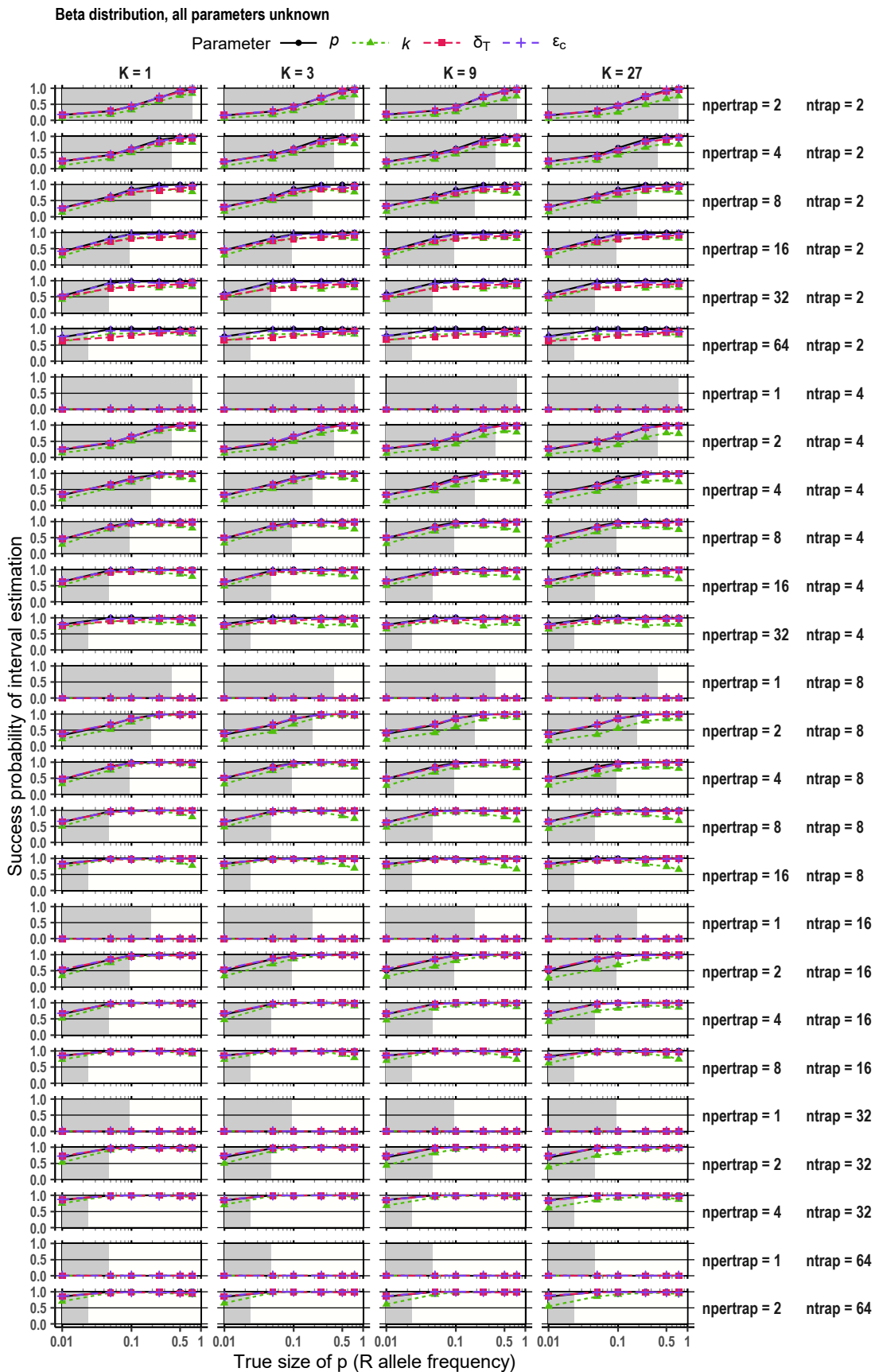
34



35

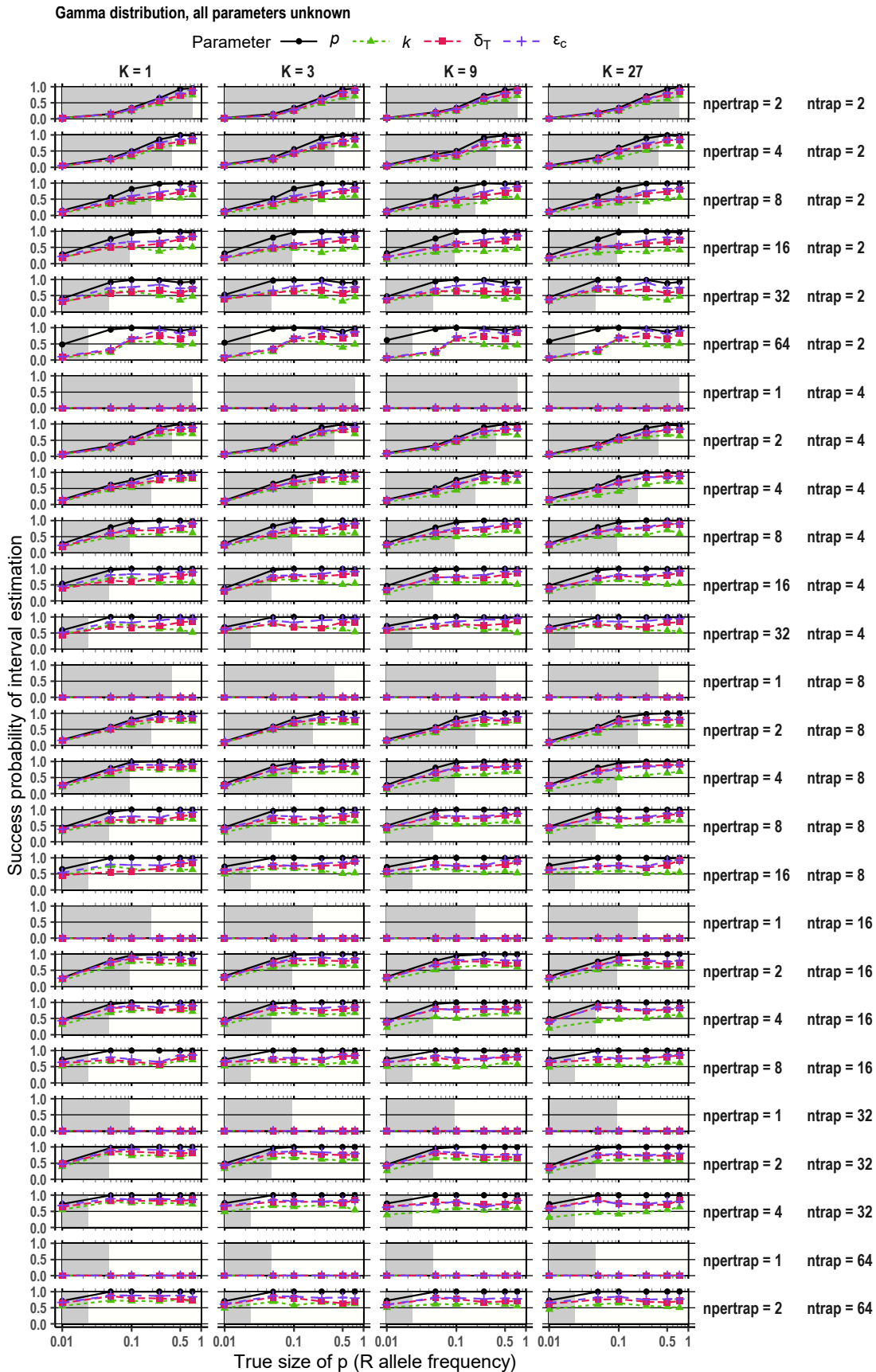
36 Figure 4 Estimation accuracy of the population allele frequency by simple averaging of  $\Delta\Delta Cq$  measures. The  
37 frequency was underestimated than its true value (horizontal broken line in each facet) as the samples were  
38 more divided. The Cq dataset was derived from the numerical experiment of “beta distribution, all  
39 parameters unknown.”

40



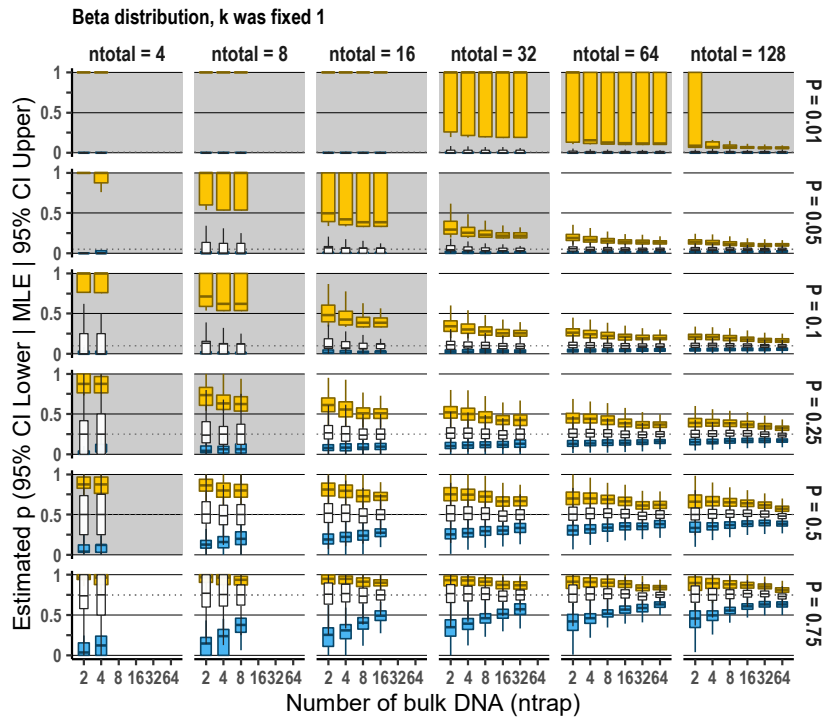
41

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



45

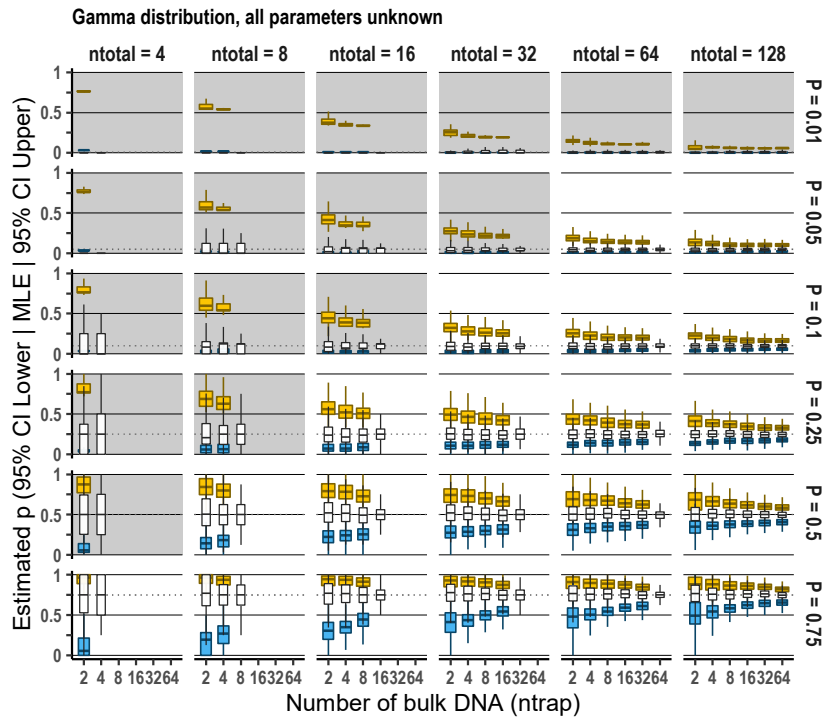
46 Figure S2 Probability of estimation success with freqpcr(). The gamma distributions were assumed, and all  
 47 estimable parameters were set as unknown. The function often failed to calculate the CIs for  $k$  when  $npertrap$   
 48 (individuals in each bulk sample) were larger.



49

50 Figure S3 Estimation accuracy of the population allele frequency,  $p$ , with `freqpcr()` when the beta distribution  
51 was assumed, considering  $K = 1$ .

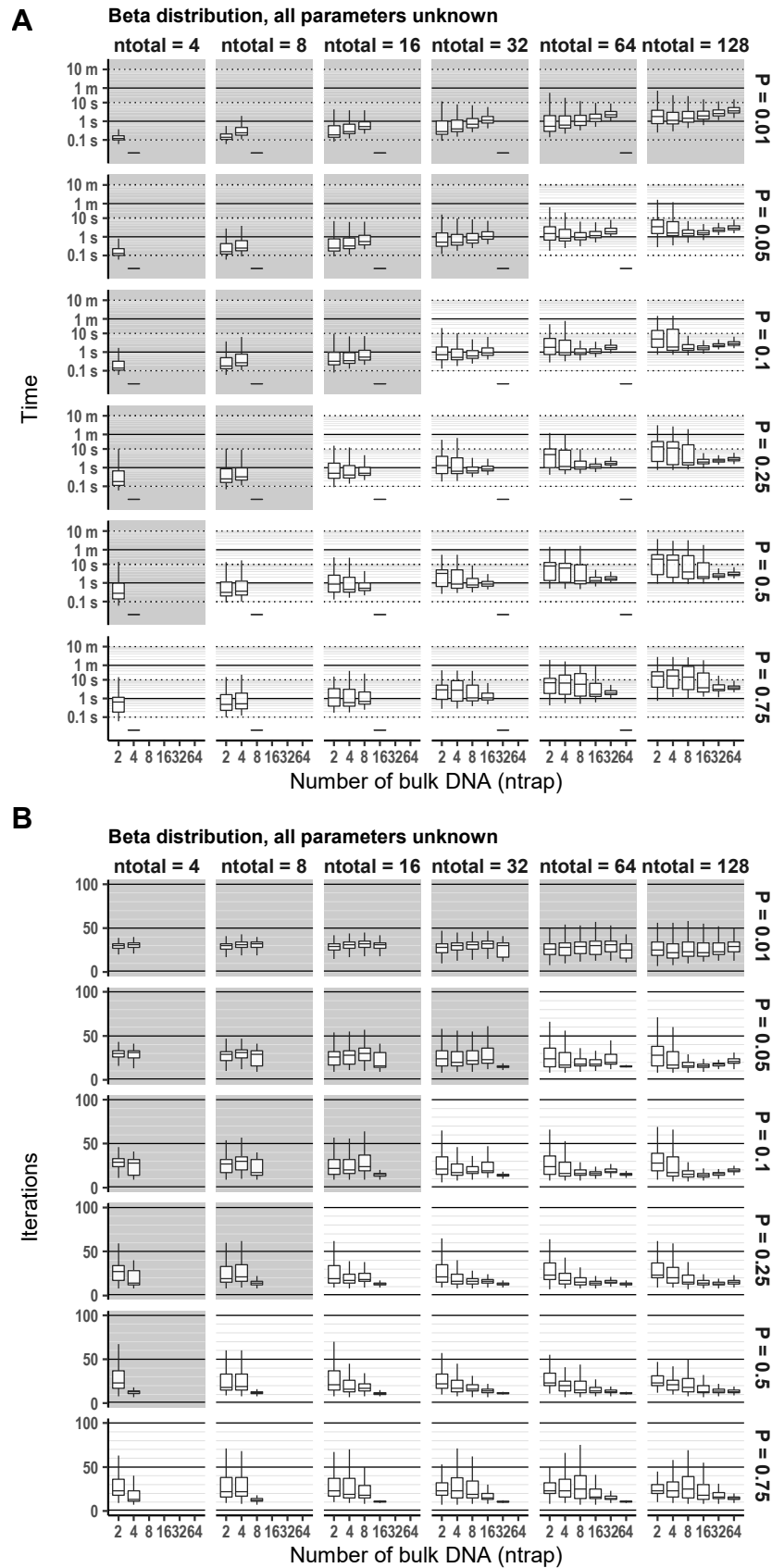
52



53

54 Figure S4 Estimation accuracy of  $p$  with `freqpcr()` when gamma distributions were assumed and all estimable  
55 parameters were set as unknown.

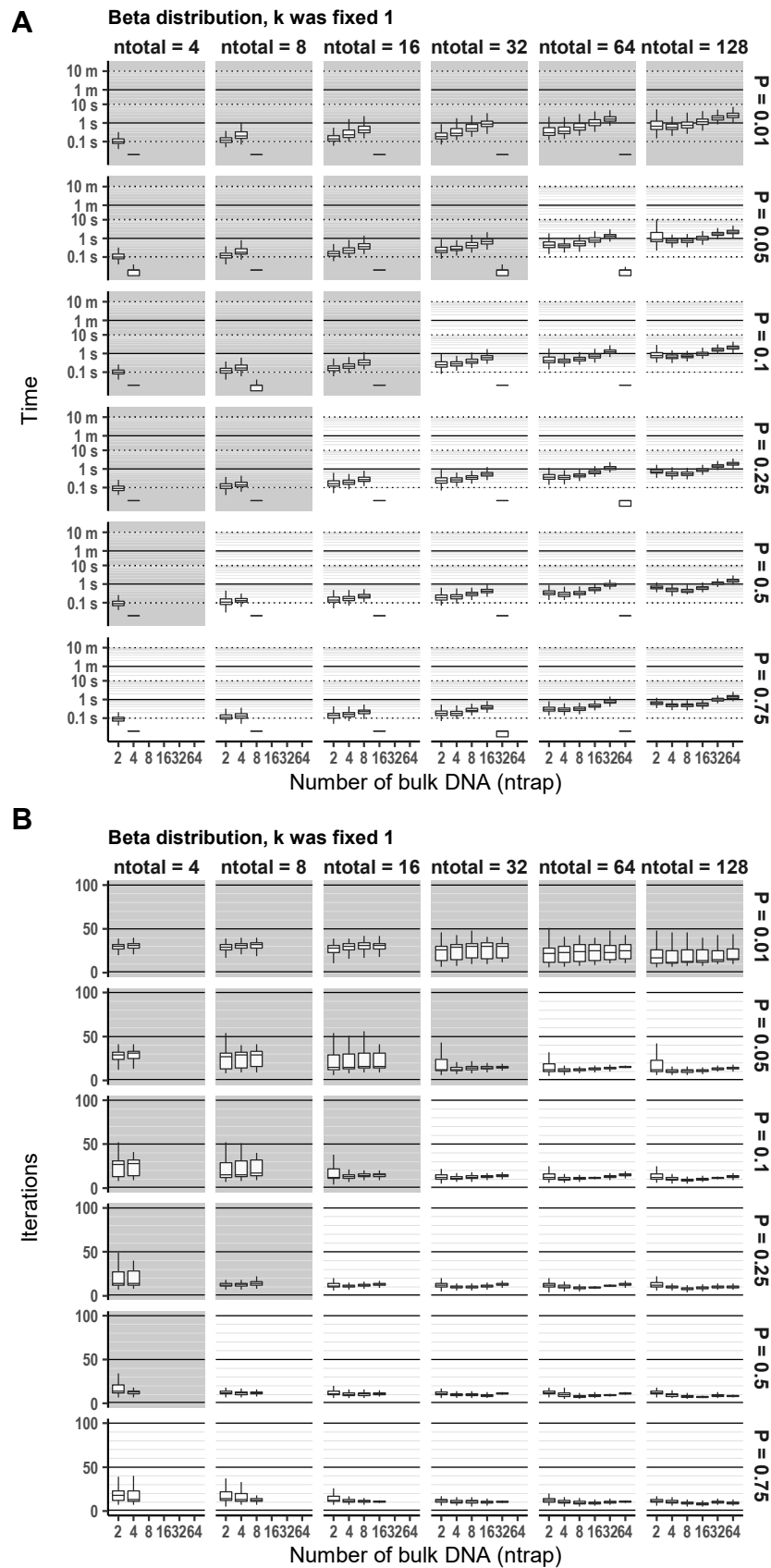
56



57

58 Figure S5 Calculation time (A) and number of iterations (B) until the freqpccr() function converges. The beta  
59 distribution was assumed, and all estimable parameters were set as unknown.

60

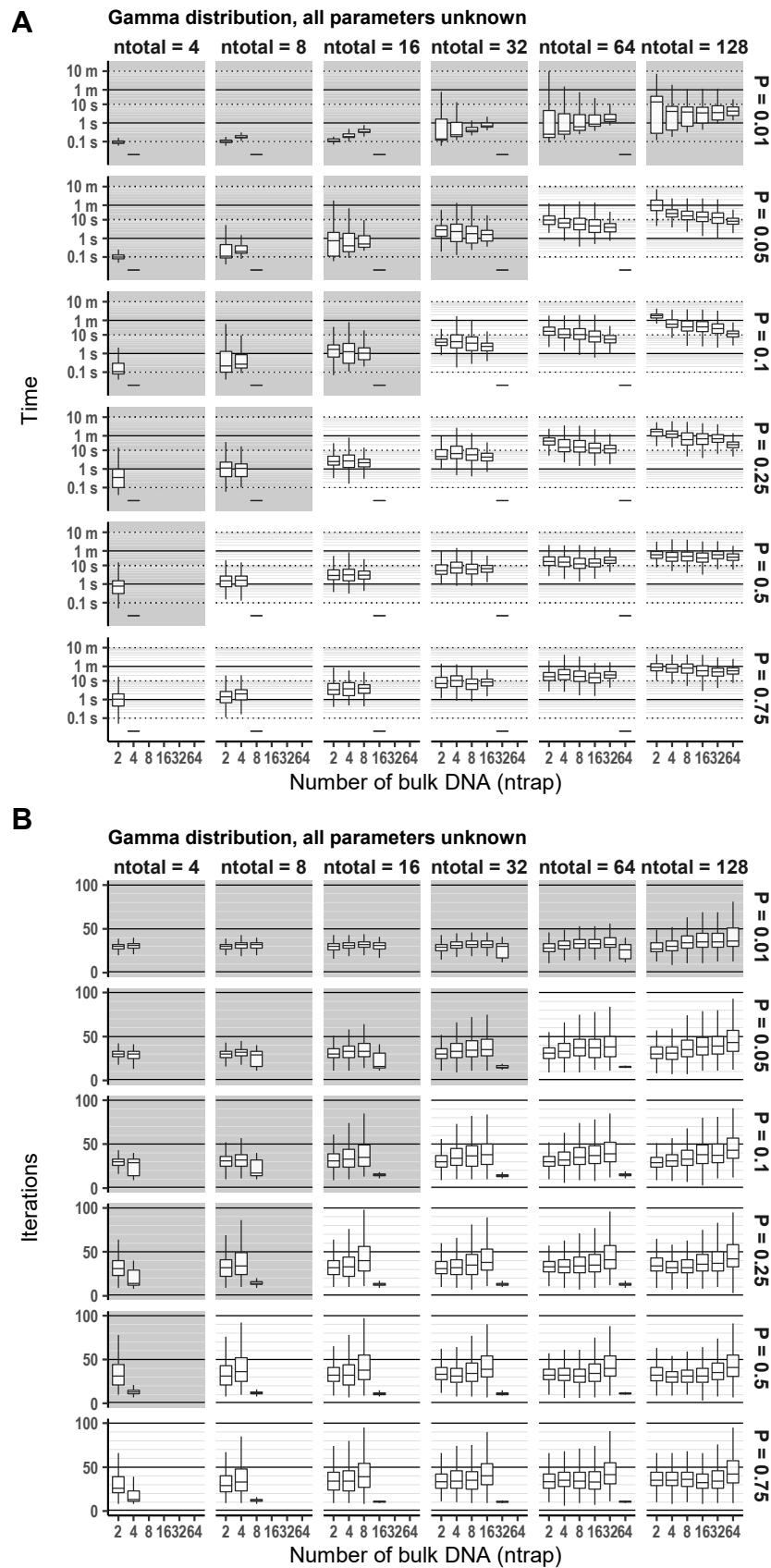


61

62 Figure S6 Calculation time (A) and number of iterations (B) until the freqpccr() function converges. The beta  
63 distribution was assumed, fixing the gamma shape parameter  $K = 1$ .

64

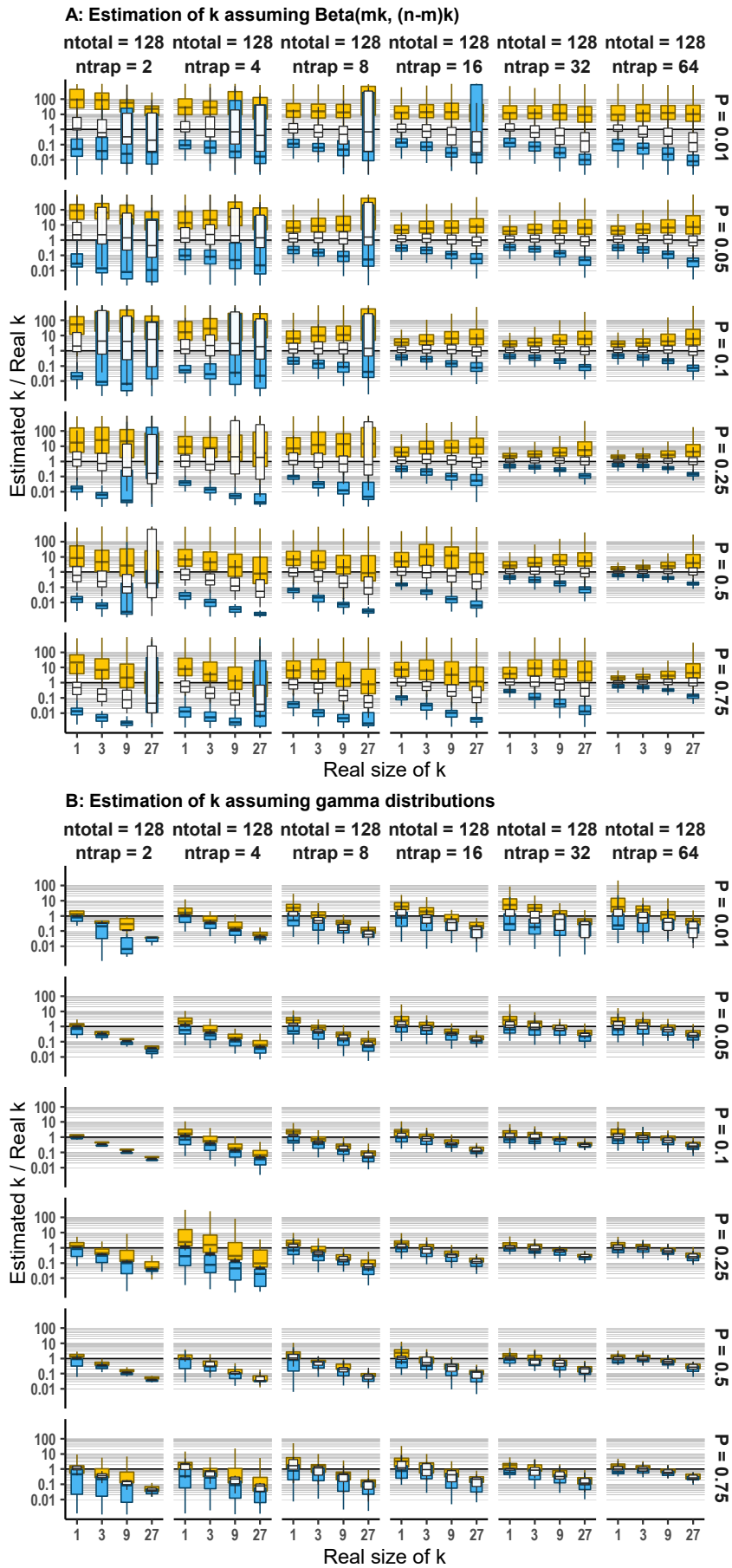




65

66 Figure S7 Calculation time (A) and number of iterations (B) until the freqpccr() function converges, assuming  
67 gamma distributions. All estimable parameters were set as unknown.

68



69

70 Figure S8 Estimation accuracy of  $k$  (the gamma shape parameter) in the simulation, showing the maximum  
71 likelihood estimate by `freqpccr()` divided by the actual parameter size.