

1 Individual identification from genetic marker data: developments and  
2 accuracy comparisons of methods

3

4 **J. Wang**

5 *Institute of Zoology, Zoological Society of London, London NW1 4RY, United Kingdom*

6

7 *Left running head:* J. Wang

8 *Right running head:* Marker-based individual identification

9 *Key words:* Relationship, Relatedness, Genetic Markers, Clone, Duplicates

10 *Corresponding author*

11

12 Jinliang Wang

13 Institute of Zoology

14 Regent's Park

15 London NW1 4RY

16 United Kingdom

17 Tel: 0044 20 74496620

18 Fax: 0044 20 75862870

19 Email: jinliang.wang@ioz.ac.uk

20

## Abstract

21

22 Genetic marker based identification of distinct individuals and recognition of duplicated individuals  
23 has important applications in many research areas in ecology, evolutionary biology, conservation  
24 biology and forensics. The widely applied genotype mismatch (MM) method, however, is  
25 inaccurate because it relies on a fixed and suboptimal threshold number ( $T_M$ ) of mismatches, and  
26 often yields self-inconsistent pairwise inferences. In this paper I improved MM method by  
27 calculating an optimal  $T_M$  to accommodate the number, mistyping rates, missing data and allele  
28 frequencies of the markers. I also developed a pairwise likelihood relationship (LR) method and a  
29 likelihood clustering (LC) method for individual identification, using poor-quality data that may  
30 have high and variable rates of allelic dropouts and false alleles at genotyped loci. The 3 methods  
31 together with the relatedness (RL) method were then compared in accuracy by analysing an  
32 empirical frog dataset and many simulated datasets generated under different parameter  
33 combinations. The analysis results showed that LC is generally one or two orders more accurate for  
34 individual identification than the other methods. Its accuracy is especially superior when the  
35 sampled multilocus genotypes have poor quality (i.e. teemed with genotyping errors and missing  
36 data) and highly replicated, a situation typical of noninvasive sampling used in estimating  
37 population size. Importantly, LC is the only method that guarantees to produce self-consistent  
38 results by partitioning the entire set of multilocus genotypes into distinct clusters, each cluster  
39 containing one or more genotypes that all represent the same individual. The LC and LR methods  
40 were implemented in a computer program COLONY for free download from the internet.

41

## 42 **Introduction**

43 Identification of distinct individuals and recognition of duplicated individuals from genetic marker  
44 data is important in many research areas in ecology, evolutionary biology, conservation biology and  
45 forensics. It has been used to estimate population size (or species abundance) in the traditional  
46 capture-mark-recapture (CMR) framework (Palsbøll *et al.* 1997; Schwartz *et al.* 1998; Creel *et al.*  
47 2003; Luikart *et al.* 2010), to track individuals across different life cycle stages in studying  
48 population parameters such as survivorship (Ringler *et al.* 2015) and migration, to infer colonial  
49 reproduction rates (Escaravage *et al.* 1998; Halkett *et al.* 2005), and to trace illegally killed animals  
50 or illegal trading animal products in wildlife forensics (Alacs *et al.* 2010). It can and should also be  
51 routinely used as a data cleaning tool to remove accidentally duplicated individuals before  
52 conducting various analyses of the raw genotype data. This is because, similar to close relatives but  
53 to a greater extent, duplicated individuals inadvertently included in a genetic analysis can reduce the  
54 estimates of genetic diversity, bias the estimates of fixation indices ( $F_{IS}$ ,  $F_{IT}$  and  $F_{ST}$ ), induce  
55 deviations from Hardy-Weinberg and linkage equilibrium, and ruin a population structuring  
56 inference (Anderson & Dunham 2008; Rodríguez-Ramilo & Wang 2012).

57         When marker information is ample (i.e. many polymorphic loci) and completely reliable (i.e.  
58 no mutations and no genotyping errors), individual identification is straightforward. In this ideal  
59 situation, identical multilocus genotypes (MGs) represent duplicated individuals and non-identical  
60 MGs correspond to distinct individuals. Nowadays with the wide application of highly polymorphic  
61 markers such as microsatellites and many genomic markers of SNPs, information content is no  
62 longer considered a constraint in practice. However, data quality could be a serious problem,  
63 especially in the case of noninvasive DNA samples such as hair, feathers and scats (Taberlet *et al.*  
64 1999; Pompanon *et al.* 2005). Due to the limited quantity and quality of DNA extracted from  
65 noninvasive samples, the presence of PCR inhibitors and DNA contaminations, noninvasive  
66 genotype data are characterized by high rates of missing data, false alleles and allelic dropouts  
67 (Bonin *et al.* 2004). Indeed, genotyping errors are a rule rather than an exception. Even genotypes  
68 obtained from DNA of high quality and quantity (e.g. extracted from fresh tissue or blood samples)  
69 are not exempt from mistypings (Pompanon *et al.* 2005). The more markers are genotyped, the  
70 higher the probability that an MG contains genotyping errors.

71         Unfortunately, individual identification is particularly vulnerable to genotyping errors in  
72 comparison with other genetic data analyses such as population genetic diversity or structure,  
73 because just one single error in an MG could create a false (ghost) individual. Even if genotyping  
74 errors occur at a very low rate  $e$  per locus, the probability that an MG contains one or more errors,

75  $E = 1 - (1 - e)^L$ , can be high, and increases rapidly with the number of loci  $L$ . For example, a 10-,  
76 50- and 250-locus genotype is expected to contain at least one mistyping with a probability of 1.0%,  
77 4.9% and 22.1% respectively when  $e=0.001$ , of 9.6%, 39.5% and 91.9% respectively when  $e=0.01$ ,  
78 and of 40.1%, 92.3% and 100% respectively when  $e=0.05$ . This result has prompted several  
79 researchers to suggest that individual identification should use the minimum number of loci  
80 required to attain a low probability of identity among samples from different individuals (Waits *et*  
81 *al.* 2001; Creel *et al.* 2003). This suggestion can reduce ghost individuals due to genotyping errors,  
82 but unfortunately it could also seriously limit the power of individual identification, especially in  
83 the difficult situation where many close relatives such as full siblings are present (Waits *et al.* 2001).

84 The problems of and difficulties in individual identification due to genotyping errors are  
85 made more prominent by high sample replications where many replicated samples could be  
86 collected from a single individual. Scat or hair based non-invasive samples (e.g. Creel *et al.* 2003)  
87 often exhibit massive replications with potentially tens to hundreds of replicated samples per  
88 individual. At this high level of replications, even a very small genotyping error rate could result in  
89 extreme overestimates of distinct individuals and of population size (Waits & Leberg 2000; Creel *et*  
90 *al.* 2003; McKelvey & Schwartz 2004). High sample replications coupled with genotyping errors  
91 and missing data can also result in numerous conflicts in pairwise inferences by any method  
92 (including the mismatch method) that compares pairs of samples (multilocus genotypes). For  
93 example, sample A may be inferred to be a duplicate of both sample B and sample C, but B and C  
94 may be inferred to come from distinct individuals.

95 A more robust and error-tolerant approach is to accept the presence of genotyping errors and  
96 accommodate them in recognizing individuals from MG data by the mismatch (MM) method. A  
97 common practice is that two samples having identical genotypes at all but 1 or 2 loci are accepted as  
98 being from a single individual and the mismatches are regarded as genotyping errors. This approach  
99 has been implemented in several computer programs, such as GENECAP (Wilberg & Dreher 2004).  
100 The allowance of a small threshold number,  $T_m$ , of 1 or 2 mismatches could reduce ghost  
101 individuals substantially. However, this threshold is obviously arbitrary, the optimum being  
102 dependent on factors such as the mistyping rates and the number of loci. While 1 or 2 mismatches  
103 may be sufficient to reduce ghost individuals when both mistyping rates and number of loci are low  
104 (say,  $e < 0.05$  and  $L < 20$ ), more mismatches should be allowed for when  $e$  or/and  $L$  are high. To  
105 overcome the problem, Galpern *et al.* (2012) proposed to determine  $T_m$  as the value where the  
106 number of individuals with more than one MG in a sample has a second minimum. Although their  
107  $T_m$  no longer relies on a predefined value, it depends on a similarity index defined to penalize  
108 arbitrarily missing and mismatched genotypes at a locus by  $1/(2L)$  and  $1/L$  respectively.

109 Furthermore, analyses of simulated (Galpern *et al.* 2012, Table 3) and empirical data (Ringler *et al.*  
110 2015) showed that this flexible  $T_m$  approach has a similar accuracy to the approach with a fixed  $T_m$   
111 =2.

112 A more powerful approach to individual identification is via pairwise relatedness analysis.  
113 Relatedness analysis is resilient to genotyping errors (Wang 2007), and can use allele frequency as  
114 well as genotype information in identifying duplicated individuals from other competitive  
115 relationships such as full siblings (Ringler *et al.* 2015). In diploid species, two MGs are expected to  
116 have a relatedness,  $r$ , of 1 and 0.5 if they come from the same individual and from two first-class  
117 relatives (full sibs and parent-offspring), respectively. Therefore, MGs are inferred to represent  
118 duplicates of the same individual when their estimated relatedness is closer to 1 than to 0.5 (i.e.  
119 when their estimated relatedness is above an appropriate threshold  $r$  value, say  $T_r = 0.75$ ). Otherwise,  
120 they are inferred to represent distinct individuals.

121 In this study, I will improve the mismatch method by calculating and using an optimal  $T_m$   
122 that takes into account mistyping rates, missing data, and the number and allele frequencies of  
123 markers. I also propose two new likelihood approaches to efficient individual identification from  
124 genotype data of low quality. One is based on calculating the likelihood values of two MGs for their  
125 candidate relationships of clone mates (duplicates) and close competitive relationships (full siblings  
126 and parent offspring), and the other is based on partitioning (in a likelihood framework) the entire  
127 set of MGs into clusters with each cluster containing one or more genotypes that all represent the  
128 same individual. Both approaches accommodate genotyping errors and use allele frequency  
129 information, and the likelihood clustering method abandons the pairwise approach such that the  
130 inferences are guaranteed to be consistent and are especially accurate for the difficult situation of  
131 high sample replications. The accuracy of these approaches is evaluated and compared by analysing  
132 many simulated and an empirical dataset.

## 133 **Methods**

### 134 *Dyadic mismatch method (MM)*

135 The threshold value of mismatches,  $T_m$ , is critical for the mismatch method. The number of distinct  
136 individuals will be overestimated and underestimated when  $T_m$  is too small and too large,  
137 respectively. The optimal  $T_m$  that minimizes falsely detected ( $\alpha$ -error) and undetected ( $\beta$ -error)  
138 individuals depends on the rate of genotyping errors, the number of loci, the allele frequencies of  
139 each locus, and the actual genetic structure (i.e. the actual relationships) of the focal set of MGs.

140 The latter is unknown and is the target of the analysis, but the former three pieces of information are  
141 usually available and can be used to resolve an approximately optimal  $T_m$ .

142 Suppose locus  $l$  has  $K_l$  alleles with estimated frequencies  $p_{li}$ , where  $i=1, 2, \dots, K_l$  and  $l=1,$   
143  $2, \dots, L$ . I assume that a genotype  $G_l$  at locus  $l$  may be mistyped to be a phenotype  $g_l$  due to allelic  
144 dropouts (ADO) at rate  $\varepsilon_{l1}$  and false alleles (FA) at rate  $\varepsilon_{l2}$ . ADOs and FAs are the most common  
145 genotyping errors in microsatellites (Bonin *et al.* 2004; Pompanon *et al.* 2005). For ADO, I assume  
146 each of the two gene copies in a diploid genotype has the same probability of dropping out during  
147 PCR, and double dropouts (i.e. both gene copies dropping out to produce no PCR products) are rare  
148 and negligible. Under this model, ADO affects heterozygote genotypes only, and a heterozygote  $G_l$   
149  $=\{w,x\}$  ( $w \neq x$ ) is observed to be a phenotype  $g_l = \{w,x\}$ ,  $\{w,w\}$  and  $\{x,x\}$  with probabilities  $1-2e_{l1}$ ,  
150  $e_{l1}$  and  $e_{l1}$  respectively, where  $e_{l1} = \varepsilon_{l1}/(1 + \varepsilon_{l1})$ . For FA, I assume that any allele in any genotype  
151 is independently and equally probable to be mistyped to be any one of the other alleles, at a rate  $e_{l2}$   
152  $= \varepsilon_{l2}/(K_l-1)$ .

153 Given allele frequencies ( $p_{li}$ ) and mistyping rates ( $\varepsilon_{l1}$  and  $\varepsilon_{l2}$ ), and assuming genotype  
154 frequencies at Hardy-Weinberg equilibrium (HWE), I can apply the above ADO and FA models to  
155 each of the  $K_l(K_l + 1)/2$  genotypes twice to generate two phenotypes, and derive the probability  
156 that two phenotypes from the same genotype match,  $Q_l$ . The expression for  $Q_l$  is however a very  
157 complicated function of  $p_{li}$ ,  $\varepsilon_{l1}$  and  $\varepsilon_{l2}$ , and is not enlightening. For simplicity,  $Q_l$  is determined by  
158 simulations. First, a genotype is drawn at random from a population in HWE with allele frequency  
159  $p_{li}$ . Second, a phenotype is generated from the genotype, following the ADO model. Third, the  
160 phenotype is further modified according to the FA model. Fourth, steps 2 and 3 are repeated to  
161 generate another phenotype independently from the same genotype. Fifth, the two phenotypes are  
162 compared to determine whether they match or not. Steps 1-5 are repeated for a sufficiently large  
163 number of replicates, and the frequency of matching phenotypes gives a good estimate of  $Q_l$ .

164 The average number of mismatches between two phenotypes having the same underlying  
165 genotype at a set of  $L$  loci is calculated by  $\sum_{l=1}^L (1 - Q_l)$  rounded to the nearest integer. This  
166 optimal  $T_m$  value is expected to minimise both  $\alpha$  and  $\beta$  errors in individual identification by the MM  
167 method. Note that  $T_m$  is calculated for each pair of MGs in a sample such that missing data can be  
168 easily accommodated. If any or both MGs have missing data at locus  $l$ , then  $Q_l$  is set to 1 for the  
169 locus in the calculation. Therefore the calculated  $T_m$  values are dyad specific and lower for dyads  
170 with more missing data. In contrast to the widely applied fixed  $T_m = 2$ , this  $T_m$  value calculated from  
171  $Q_l$  accounts for allele frequencies, mistyping rates, number of loci, and missing data. Two MGs are

172 inferred to be from a single and two distinct individuals when their observed number of mismatches  
173 is not and is greater than their  $T_m$  value, respectively.

#### 174 *Dyadic relatedness method (RL)*

175 The genetic relatedness,  $r$ , between two MGs can be calculated by a marker-based moment or  
176 likelihood estimator (Wang 2007; 2014). Duplicated individuals and first-order relatives (e.g. full-  
177 sib or parent-offspring) are expected to have an  $r$  value of 1 or close to 1 and of 0.5 or close to 0.5,  
178 respectively, even when they have mismatches at a small fraction of loci due to genotyping errors  
179 (Wang 2007). To distinguish duplicated individuals from first-order relatives and to minimise both  
180  $\alpha$ - and  $\beta$ -error rates, I choose a threshold  $r$  value of  $T_r=0.75$ , which is the midpoint between the  
181 expected  $r$  values for duplicates and first-order relatives. Two MGs are inferred to be duplicates and  
182 distinct individuals when their  $r$  value is and is not greater than  $T_r$ , respectively. There are quite a  
183 few  $r$  estimators available (Wang 2014), among which I chose to use the one based on phenotype  
184 similarity, proposed by Lynch (1988) and improved by Li *et al.* (1993). It is chosen because it is  
185 simple to calculate and is expected to have a higher accuracy than other moment estimators when  
186 applied to close relationships such as identical twins (duplicates) and full sibs (Wang 2007).

#### 187 *Dyadic likelihood relationship method (LR)*

188 It is also possible to calculate directly the likelihoods of two MGs for the candidate relationships of  
189 duplicates (clone mates or identical twins, denoted by DP), full sibs (FS), half sibs (HS), parent  
190 offspring (PO) and unrelated (UR). If DP has the highest likelihood, then the two MGs are inferred  
191 to come from the same individual. Otherwise, they are inferred to come from distinct individuals.

192 In contrast to pairwise relatedness estimation, relationship inference is highly vulnerable to  
193 genotyping errors. A single error could exclude truly duplicated MGs from being inferred as such.  
194 The more markers one uses, the more serious the false exclusion problem will become. The  
195 likelihood functions of FS, HS and PO are available in the literature, but they do not account for  
196 genotyping errors (e.g. Goodnight & Queller 1999) or account for ADO only (e.g. Wagner *et al.*  
197 2006). Herein I show the general likelihood function applying to any pairwise relationship  
198 (including DP, FS, HS, PO and UR) and allowing for both ADO and FA occurring at rates variable  
199 across loci.

200 The genetic relationship between two non-inbred individuals is fully specified by 3 identical  
201 by descent (IBD) coefficients  $\Delta_i$ , where  $\Delta_i$  is the probability that the two individuals share exactly  $i$   
202 ( $i=0, 1, 2$ ) pairs of gene copies IBD at a locus. Obviously,  $\Delta_0 + \Delta_1 + \Delta_2 \equiv 1$ . In diploid species,  $\Delta_0, \Delta_1$   
203 and  $\Delta_2$  have values 0, 0 and 1 for DP, 0.25, 0.5 and 0.25 for FS, 0, 0.5, 0.5 for HS, 0, 1 and 0 for PO,



204 and 0, 0, 1 for UR. The probability of observing a phenotype  $g_A=\{a,b\}$  for individual A and a  
 205 phenotype  $g_B=\{c,d\}$  for individual B at a locus with  $K$  codominant alleles, given their relationship  
 206 defined by  $\Delta_0, \Delta_1$  and  $\Delta_2$ , is (Wang 2006)

$$208 \quad \Pr[a, b; c, d|\Delta_0, \Delta_1, \Delta_2] = \sum_{u=1}^K \sum_{v=u}^K \sum_{w=1}^K \sum_{x=w}^K R[u, v; w, x|\Delta_0, \Delta_1, \Delta_2] \Pr[a, b|u, v] \Pr[c, d|w, x], \quad (1)$$

207 where

$$209 \quad R[u, v; w, x|\Delta_0, \Delta_1, \Delta_2] \\
 210 \quad = (2 - \delta_{uv})p_u p_v \left( \Delta_0(2 - \delta_{wx})p_w p_x + \frac{1}{4}\Delta_1(2 - \delta_{wx})((\delta_{uw} + \delta_{vw})p_x) \right. \\
 211 \quad \left. + (\delta_{ux} + \delta_{vx})p_w \right) + \Delta_2(\delta_{uw}\delta_{vx} + \delta_{ux}\delta_{vw} - \delta_{uw}\delta_{vx}\delta_{ux}\delta_{vw}) \quad (2)$$

212 is the probability that A and B have genotype  $\{u,v\}$  and  $\{w,x\}$  respectively conditional on their  
 213 relationship or IBD coefficients  $\Delta_0, \Delta_1, \Delta_2$ , and  $\delta_{uv}$  (and similarly for other  $\delta$  variables) is the  
 214 Kronecker delta variable with values 1 and 0 when  $u=v$  and  $u \neq v$ , respectively. In (1),  $\Pr[u, v|w, x]$   
 215 is the probability that a genotype  $\{w,x\}$  shows a phenotype  $\{u,v\}$  due to ADO and FA. It is derived  
 216 as (Wang 2004)

$$217 \quad \Pr[u, v|w, x] = \begin{cases} (1 - \varepsilon_2)^2 + e_2^2 - 2e_1 e_3^2 & (u = w, v = x) \\ e_2(1 - \varepsilon_2) + e_1 e_3^2 & (u = v = w) \text{ or } (u = v = x) \\ (2 - \delta_{u,v})e_2^2 & (u \neq w, u \neq x, v \neq w, v \neq x) \\ e_2 e_3 & (\text{otherwise}) \end{cases} \quad (3)$$

218 for a heterozygous genotype ( $w \neq x$ ) where  $e_3 = 1 - \varepsilon_2 - e_2$ , and

$$219 \quad \Pr[u, v|w, x] = \begin{cases} (1 - \varepsilon_2)^2 & (u = v = w) \\ 2e_2(1 - \varepsilon_2) & (u = w, v \neq w) \text{ or } (v = w, u \neq w) \\ (2 - \delta_{u,v})e_2^2 & (u \neq w, v \neq w) \end{cases} \quad (4)$$

220 for a homozygous genotype ( $w=x$ ) under the ADO and FA models described above.

221 Note that equations (1-4) give the likelihood of a relationship for a single locus  $l$ , and  
 222 subscript  $l$  is dropped from error rates ( $\varepsilon_{l1}, \varepsilon_{l2}, e_{l1}, e_{l2}$ ) and allele frequencies ( $p_{li}$ ) for clarity. The  
 223 multilocus likelihood is simply a product of single locus likelihood values, assuming linkage  
 224 equilibrium among loci.

### 225 *Likelihood clustering method (LC)*

226 The above 3 methods take a pairwise approach, which considers whether two MGs are duplicates or  
 227 not in isolation of others. When an individual has more than 2 replicated MGs, pairwise approaches

228 may yield conflicting results. Among 3 replicated MGs A, B and C of an individual, for example, A  
229 and B as well as A and C may be inferred as DP while B and C may be inferred as distinct  
230 individuals. This happens when, for an example, A, B and C have genotypes identical at all but a  
231 single locus at which A has missing data while B and C show different alleles. The 3 pairwise  
232 inferences are obviously in conflict. The frequency of these inconsistencies increases rapidly with  
233 an increasing level of individual replications, and decreasing data information quality and quantity.  
234 Furthermore, pairwise approaches do not use marker information fully, and thus are expected to  
235 have a lower power (accuracy) than approaches that consider the relationship among all MGs  
236 simultaneously (Wang 2004).

237 A more desirable approach is to partition the entire set of MGs into  $N$  (unknown) individual  
238 clusters, with each cluster containing one or more MGs that all represent the same individual. To  
239 reduce both  $\alpha$  and  $\beta$  errors, the clustering should be better made by considering several competitive  
240 relationships such as DP, FS and HS which could generate similar patterns of MGs. The algorithm  
241 used for sibship inference (Wang 2004) can be modified to identify individuals, as shown below.

242 First, assuming each MG corresponds to a distinct individual, a sibship analysis is conducted  
243 to partition the entire set of individuals into full-sib clusters. The analysis could adopt the simple  
244 monogamy model (i.e. no inference of half sibs), or the sophisticated polygamy model (i.e.  
245 inference of half sibs). The monogamy model is preferred because it runs much faster than, but has  
246 the same or very similar accuracy to, the polygamy model for individual identification. This is  
247 because DP is much closer to FS in relatedness than to HS and is thus much less likely to confuse  
248 with HS than FS. Second, each inferred FS cluster is further partitioned by a likelihood approach  
249 into a number of individual clusters, with each cluster containing one or more MGs that all  
250 represent the same individual. The first step has been described before (Wang 2004), and the second  
251 step is detailed below.

252 Suppose an inferred FS cluster contains  $M$  ( $\geq 1$ ) MGs. If  $M=1$ , then no further analysis is  
253 needed. Otherwise, the MGs can be divided into one of a number of  $B_M$  possible partitions (or  
254 configurations), where  $B_M$  is the Bell number. A partition contains a number of  $m$  (where  $m \geq 1$  and  
255  $m \leq M$ ) individual clusters, with each cluster containing one or more MGs that all represent the same  
256 individual. Three MGs ( $M=3$ ) of A, B and C, for example, have  $B_3=5$  different partitions, which are  
257  $\{(A), (B), (C)\}$ ,  $\{(A, B), (C)\}$ ,  $\{(A, C), (B)\}$ ,  $\{(B, C), (A)\}$ ,  $\{(A, B, C)\}$  where all MGs in a pair of  
258 parentheses come from the same individual and constitute an individual cluster. Partition  $\{(A, B),$   
259  $(C)\}$ , for example, has two individual clusters which are (A, B) and (C), meaning that A and B  
260 come from one individual and C comes from another individual. Each partition is evaluated for its

261 likelihood which is equal to the probability of the genotype data given the partition, and the one  
 262 with the maximum likelihood is returned as the best estimate. The challenge is to construct, and  
 263 calculate the likelihood values of, the  $B_M$  partitions, where  $B_M$  increases explosively with  $M$ . Even  
 264 for small  $M$  value of 5, 10 and 15, for example, the corresponding  $B_M$  values are 52, 115975 and  
 265 1382958545, respectively.

266 Instead of using the simulated annealing approach in sibship analysis (Wang 2004), I take a  
 267 systematic approach to individual identification. The approach is deterministic and fast, because a  
 268 FS cluster is usually small. For a FS cluster with  $M$  MGs, the algorithm starts with an initial  
 269 configuration,  $C_0$ , of  $M$  individual clusters, each containing one MG. Round 1 searching works on  
 270  $C_0$ . Each of the  $M(M-1)/2$  possible configurations is constructed by merging two of the  $M$   
 271 clusters, and is evaluated for likelihood. The best of these configurations,  $C_1$ , with the maximum  
 272 likelihood value is then compared with  $C_0$ . If the former has a smaller likelihood, then  $C_0$  is  
 273 returned as the best estimate and the searching process terminates. Otherwise,  $C_0$  is abandoned and  
 274  $C_1$  is accepted, and round 2 searching is initiated to improve on  $C_1$ . Following exactly the same  
 275 procedure in constructing new configurations as in round 1, round 2 returns the best configuration  
 276 with  $M-2$  clusters,  $C_2$ . If  $C_2$  has a lower likelihood than  $C_1$ , then the latter is reported as the best  
 277 estimate and the searching process terminates. Otherwise,  $C_1$  is replaced by  $C_2$ , and round 3  
 278 searching is initiated to work on  $C_2$ , following the same process as in previous rounds. The whole  
 279 searching process stops when, at round  $m$ , the best of the  $(M-m+1)(M-m)/2$  reconfigurations,  
 280  $C_m$ , has a lower likelihood than that of the previous round,  $C_{m-1}$ , which is returned as the best  
 281 estimate.

282 Now consider the likelihood of a configuration with  $m$  ( $=1\sim M$ ) individual clusters, with  
 283 cluster  $i$  ( $=1, 2, \dots, m$ ) containing  $n_i$  genotypes  $g_{ij}$  ( $j=1, 2, \dots, n_i$ ) at a locus with  $K$  alleles. All  
 284 genotypes within a cluster are duplicates of the same individual, and genotypes from different  
 285 clusters represent different individuals. Obviously, we have  $\sum_{i=1}^m n_i \equiv M$ . The likelihood function is

$$286 \sum_{u=1}^K p_u \sum_{v=1}^K p_v \sum_{w=1}^K p_w \sum_{x=1}^K p_x \prod_{i=1}^m \frac{1}{4} \left( \sum_{a=u,v} \sum_{b=w,x} \prod_{j=1}^{n_i} \Pr[g_{ij}|a,b] \right), \quad (5)$$

287 where the probability of observing a phenotype  $g_{ij}$  given its underlying genotype  $G_{ij}=\{a,b\}$ ,  
 288  $\Pr[g_{ij}|a,b]$ , is calculated by (3-4). The computational cost of (5) can be much reduced by pooling  
 289 all unobserved alleles in the FS cluster into a single ‘‘allele’’ and by pooling identical parental  
 290 genotypes (e.g.  $\{u,v\}$  and  $\{v,u\}$ ) and parental genotype combinations (e.g.  $\{\{u,v\}, \{w,x\}\}$  and  
 291  $\{\{w,x\}, \{u,v\}\}$ ), as in sibship likelihood calculations (Wang 2004). For multiple loci in linkage

292 equilibrium, the likelihood is simply the product of single locus values calculated by (5). When one  
293 or both parents of the FS family are assigned to candidate adults with genotype data, the likelihood  
294 function is slightly more complicated and is not shown herein.

### 295 *Simulations*

296 Simulated data were generated and analysed comparatively by the above described 4 methods to  
297 evaluate their accuracies. A number of factors are expected to affect individual identifications, and  
298 are thus considered in the simulations.

299 First, the simulations considered the actual relatedness structures of the sampled individuals.  
300 Presence of close relatives, such as full sibs, makes individual identification more difficult by  
301 increasing  $\beta$  errors. I considered 3 sibship structures to reflect low, medium and high relatedness.  
302 These are denoted by 40(1, 1), 16(1, 1, 3) and 4(1, 2, 3, 4, 10), where the value before the brackets  
303 gives the number of replicate half-sib families and the values within the brackets are the sizes of  
304 full-sib families that are nested within a half-sib family. For example, 16(1, 1, 3) means there are 16  
305 half-sib families, and each family has a single father mated with 3 mothers who give 1, 1, and 3 full  
306 siblings. Each of the 3 sibship structures yields 80 distinct individuals (genotypes) in a sample.  
307 Other close relatives such as parent-offspring may also be present in a practical sample. However,  
308 these relationships have much smaller effect on individual identification than full sibs, because the  
309 latter are more likely to generate identical or nearly identical MGs. Therefore, relatives other than  
310 full sibs are not considered in the simulations.

311 Second, the simulations allowed for different extents of individual replications. The number  
312 of individual genotype replications is assumed to follow a Poisson distribution with parameter  $\lambda$ ,  
313 taking values between 0 and 5. For each of the 80 distinct individuals in a sample, a random number  
314  $R \sim \text{Poisson}[\lambda]$  is generated and the MG of the individual is replicated by  $R$  times.

315 Third, the simulations considered different numbers and polymorphisms of markers. For  
316 given numbers of loci ( $L$ ) and alleles ( $K_l$ ) per locus, allele frequencies were drawn from a uniform  
317 distribution at each locus, and the 80 MGs in a given sibship structure were generated by assuming  
318 Hardy-Weinberg and linkage equilibrium. These MGs were faithfully replicated according to  
319 Poisson[ $\lambda$ ] as described above. When considering the impact of  $K_l$ , I vary  $K_l$  and  $L$  simultaneously  
320 such that the total number of independent alleles across loci,  $\sum_{l=1}^L (K_l - 1)$ , is fixed at 160.

321 Fourth, the simulations allowed for different rates of ADO, FA and missing data at each  
322 locus. After replications, each MG is modified independently at each locus for ADOs, FAs, and  
323 missing data to generate the corresponding multilocus phenotype. Suppose ADO, FA and missing

324 data occur at rates  $\varepsilon_{l1}$ ,  $\varepsilon_{l2}$  and  $\varepsilon_{l3}$  at locus  $l$ , respectively. A maximum of 3 steps are required to  
325 generate the phenotype at this locus from its genotype. In step 1, a random number  $R$  uniformly  
326 distributed in the range  $[0,1]$  is drawn. If  $R \leq \varepsilon_{l3}$ , then the phenotype becomes  $\{0,0\}$  to indicate  
327 missing data. Otherwise, the genotype is subject to ADO in step 2. Another random number  $R$  is  
328 drawn. If the genotype is a heterozygote and  $R \leq \varepsilon_{l1}/(1+\varepsilon_{l1})$ , then the phenotype is returned as a  
329 homozygote for an allele drawn at random from the genotype. Otherwise, the genotype has no  
330 changes in step 2. In both cases, the genotype is subject to FA in step 3. For each allele in the  
331 genotype, a random number  $R$  is drawn. If  $R \leq \varepsilon_{l2}$ , then the allele is changed to another allele  
332 randomly drawn from the  $K-1$  alleles. Otherwise, no change is made to the allele.

333 Fifth, all methods except for RL use ADO and FA rates at each locus. In practice, these  
334 mistyping rates are usually unknown, but are estimated from duplicated genotyping or pedigree  
335 based analysis (Creel *et al.* 2003; Pompanon *et al.* 2005). It is important to know how robust these  
336 methods are to mis-specified mistyping rates. For this purpose, I simulated data with a true  
337 mistyping rate of  $\varepsilon_{l1}=\varepsilon_{l2}=0.1$  for each locus  $l$ , but analysed the data assuming values of  $\varepsilon_{l1}=\varepsilon_{l2}$  in  
338 the range of 0 to 0.2.

#### 339 *Accuracy assessment*

340 Accuracy is assessed by the proportion of MG dyads in a dataset that are from a single individual  
341 but are incorrectly identified as from distinct individuals ( $\alpha$  errors, falsely identified individuals),  
342 and that are from distinct individuals but are incorrectly identified as from a single individual ( $\beta$   
343 errors, unidentified individuals). The overall accuracy including both types of errors is measured by  
344 the proportion of MG dyads in a dataset that are incorrectly inferred to be non-duplicates or  
345 duplicates,  $\gamma$ . These  $\alpha$ -,  $\beta$ - and total-error rates were calculated for each dataset and averaged across  
346 100 replicate datasets for a given parameter combination. Because most applications are affected by  
347 both  $\alpha$ - and  $\beta$ -errors, I report the total error rate,  $\gamma$ , to indicate overall accuracy in this paper to save  
348 space.

#### 349 *Empirical data*

350 Ringler *et al.* (2015) showed that microsatellites can be used to reliably mark amphibian larvae and  
351 to re-identify them after metamorphosis. They genotyped 1800 tadpoles of the dendrobatid frog  
352 (*Allobates femoralis*) at 14 highly polymorphic microsatellite loci before releasing them on a 5-ha  
353 river island which was previously uninhabited by this species. They surveyed the island and  
354 sampled 42 juvenile individuals six months after the release, and sampled 36 males and 31 females  
355 one year after the release. The sampled juveniles and adults were released to their capture sites after

356 taking DNA samples, which were genotyped at the same set of 14 loci as the tadpoles. Based on  
357 their unique ventral patterns, 20 of the 67 adults were identified to correspond to one of the 42  
358 juveniles. These 20 individuals sampled as both juveniles and adults were mostly confirmed by  
359 relatedness analysis of marker data. Individual identification between tadpoles and juveniles or  
360 between tadpoles and adults was based on the mismatch and relatedness methods. In the present  
361 study, the genotype data are comparatively analysed by the 4 individual identification methods.

## 362 **Results**

363 Simulations under the three sibship structures yield qualitatively similar results, and thus only the  
364 analysis results for sibship structure 4(1, 2, 3, 4, 10) are reported below.

### 365 *Effect of the number of markers*

366 The optimal  $T_m$  determined by the simulation procedure gives an unbiased estimate of the average  
367 number of mismatches between duplicated MGs for different numbers of loci  $L$  (Fig. 1) and for  
368 different mistyping and data missing rates (not shown). For a given  $L$ , calculated  $T_m$  values vary  
369 because different MG dyads may have different numbers of loci at which genotype data are missing,  
370 and because different loci may have different  $Q_l$  values. However, the variation of  $T_m$  values is  
371 much smaller than the variation of the observed numbers of mismatches, and the difference  
372 increases with  $L$ . Part of the reason that the mismatch method is less accurate than other methods  
373 (see below) is the high variation of the observed number of mismatches around  $T_m$ , which results in  
374 high rates of both  $\alpha$ - and  $\beta$ -errors.

375 With an increasing number of markers, the accuracy of mismatch (MM) method is almost  
376 constant, while that of relatedness (RL), likelihood relationship (LR), and likelihood clustering (LC)  
377 methods increases rapidly (Fig. 1). This means RL, LR and LC are statistically consistent, but MM  
378 is not, even when an optimal  $T_m$  value was calculated and used in the analysis. MM makes  
379 decreasing  $\beta$ -errors (undetected individuals) but increasing  $\alpha$ -errors (falsely detected individuals)  
380 with an increasing  $L$ , as expected. As a result, the overall error rate  $\gamma$  is almost constant with an  
381 increasing  $L$  (Fig. 1). If a fixed value of  $T_m = 2$  were used, MM method would perform much worse  
382 with a much higher  $\gamma$  due to excessive  $\beta$ -errors when  $L < 10$  or excessive  $\alpha$ -errors when  $L > 10$ .

383 LC is the most accurate method for different numbers of markers, followed by LR. These  
384 two methods become more and more accurate than RL method with an increasing number of loci.  
385 When  $L=80$ , perfect inference ( $\alpha = \beta = 0$ ) is obtained by both LC and LR methods.

### 386 *Effect of the number of alleles*

387 For different numbers of alleles per locus and thus different numbers of loci when the total number  
388 of independent alleles is fixed at 160, LC method always has the lowest  $\alpha$  error rate and the second  
389 lowest  $\beta$  error rate (Fig. 2). LR has an  $\alpha$  error rate only slightly larger than LC, but has the highest  
390  $\beta$  error rate. MM has an  $\alpha$  error rate much larger and a  $\beta$  error rate much smaller than the other three  
391 methods. Overall, LC is the most accurate, making much fewer  $\alpha$  and  $\beta$  errors than the other  
392 methods.

393 At a fixed total number of 160 independent alleles, the overall accuracy of the 4 methods  
394 first increases and then decreases with an increasing number of alleles per locus,  $K$  (Fig. 2). The  
395 maximal accuracy is achieved when  $K=5$  for all methods except for the mismatch (MM) method.  
396 The RL and LR methods have an indistinguishable overall accuracy, which is higher than that of  
397 MM but much lower than that of LC for different numbers of alleles per locus. The accuracy  
398 differences among methods increases with a decreasing number of alleles per locus and  
399 correspondingly an increasing number of loci.

#### 400 *Effect of the extent of individual replication*

401 Contrasting behaviours of different methods are observed for different levels of individual  
402 replications,  $\lambda$  (Fig. 3). With an increasing  $\lambda$ , the accuracy of LR is almost constant, that of MM and  
403 RL decreases, while that of LC increases. When a sample contains no replicated individuals (i.e.  
404  $\lambda=0$ ), MM has the lowest overall error rate  $\gamma$  because it has no chance to falsely identify individuals  
405 ( $\alpha$  errors) to which the method is particularly vulnerable. However, MM quickly becomes the least  
406 accurate method at a low value of  $\lambda=0.3$ , when roughly each of 30% individuals is replicated only  
407 once. The clustering method LC always outperforms the 3 pairwise approaches when there exist  
408 replicated individuals in a sample, and this advantage increases steadily with the replication level  $\lambda$ .

#### 409 *Effect of mistyping and missing data rates*

410 Genotyping errors and missing data decrease marker information and increase noises. As a result,  
411 all 4 methods show a decreasing accuracy with an increasing mistyping and missing data rate (Fig.  
412 4). The mismatch method is especially susceptible to mistyping and missing data. Its accuracy  
413 quickly reduces to the lowest when  $\varepsilon_{l1}=\varepsilon_{l2}=\varepsilon_{l3}$  raises to a low value of 0.01 for each of 20 loci. For  
414 the entire range of mistyping and missing data rates from 0 to 0.16, LC has the highest accuracy,  
415 followed by LR.

#### 416 *Robustness to mis-specified mistyping rates*

417 The relatedness method does not use (account for) mistyping rates and thus its accuracy is  
418 unaffected by the assumed mistyping rate  $\hat{\epsilon}$  (Fig. 5). The behaviour of MM is perplexing, as its  
419 accuracy increases slowly with an increasing  $\hat{\epsilon}$  when it is actually larger than the true simulated  
420 mistyping rate  $\epsilon$ . This is because the dominating errors made by MM when marker information is  
421 not small are falsely identified individuals ( $\alpha$  errors), which can be reduced by the use of an  
422 overestimated mistyping rate. The two likelihood methods, LR and LC, have the highest accuracy  
423 when  $\hat{\epsilon}$  is roughly equal to  $\epsilon$ . Their accuracy decreases as  $\hat{\epsilon}$  deviates from  $\epsilon$ . Relatively, LR is much  
424 more vulnerable than LC to mis-specified mistyping rates, and becomes the least accurate method  
425 when roughly  $\hat{\epsilon} > 1.25\epsilon$ . Although LC is also affected by mis-specified  $\hat{\epsilon}$ , it is always the most  
426 accurate method in the range between  $\hat{\epsilon}=0$  and  $\hat{\epsilon} = 2\epsilon$ .

### 427 *Results of empirical data analysis*

428 The 1909 MGs (1800 tadpoles, 42 juveniles, 67 adults) were partitioned by LC into 1766 individual  
429 clusters, each corresponding to an inferred distinct individual. Among these clusters, 1651, 92 and  
430 23 are singletons, dyads, and trios, each containing 1, 2 and 3 MGs, respectively. Among the 23  
431 trios, each of 20 contains a morphologically identified juvenile-adult dyad and a tadpole, one  
432 contains 2 tadpoles and a juvenile, one contains 2 tadpoles and an adult, and one contains 3 tadpoles.  
433 The first 20 trios confirm morphological observations and are highly likely to be correct, while the  
434 last 3 trios are probably incorrect if no tadpoles are actually replicated in the sample. The last 3 trios  
435 have similar numbers of missing and mismatched genotypes to the first 20 trios.

436 Because juveniles and adults are subsamples of tadpoles, we expect each juvenile or adult  
437 should have a corresponding tadpole. Indeed, each of all 67 adults and each of 38 juveniles was  
438 inferred to match a tadpole, and each of the 4 remaining juveniles was inferred to match no tadpoles.  
439 This means the  $\alpha$  error (falsely identified individuals) rate of LC for this dataset is low, only about  
440 3.6% (4 out of 109). It is also possible to calculate  $\beta$  error (unidentified individuals) rate of LC for  
441 this dataset, if no individuals within a life stage (tadpoles, juveniles, adults) are actually replicated.  
442 Among the 1821186 possible dyads, only 41 dyads within a life stage were identified by LC as  
443 single individuals, giving a  $\beta$  error rate of 0.0000225. It turns out that all of the 41 dyads are  
444 tadpoles, and no adults and no juveniles were found duplicated. This is not surprising because  
445 tadpoles are much more numerous than juveniles and adults, and many tadpoles were inferred to  
446 come from large full sib families (data not shown).

447 The distributions of the numbers of loci with missing data and mismatches between a pair of  
448 MGs for various classes of dyads are shown in Fig. 6, and explain the low power and accuracy of  
449 the mismatch method. As expected, there is essentially no difference in missing data for dyads of



450 various relationship classes. The average number of loci with missing data for a dyad is 2, no matter  
451 the dyadic MGs come from a single individual, two full siblings, or two non-full siblings. However,  
452 the distributions of mismatches differ among dyads of different classes. A dyad coming from a  
453 single individual most often has 0, 1 or 2 mismatches, but can occasionally have a maximal number  
454 of 7 mismatches. A full sib dyad on average has 8 mismatches, but can have a minimal number of  
455 only 2 mismatches. A non-full sib dyad on average has 11 mismatches, with the minimal number of  
456 mismatches being 6. Using a threshold value of mismatches  $T_m = 6$  or 7, the mismatch method can  
457 confidently identify duplicated MGs (Fig. 6, E and F) and unrelated individuals (Fig. 6H) with a  
458 small  $\alpha$  and  $\beta$  error rates. However, it has tremendous difficulties to distinguish duplicated MGs  
459 from full siblings (Fig 6G). Using the optimum  $T_m$  value of 4 or 5, it still could result in substantial  
460  $\alpha$  and  $\beta$  error rates. The analysis shown in Fig. 6 also demonstrates that the optimal  $T_m$  value is not  
461 only marker property (e.g. number, polymorphisms, genotyping error rates, data missing rates) but  
462 also sample genetic structure (i.e. distributions of relatedness among MGs) dependent. The optimal  
463  $T_m$  value would be 6~7 and 4~5 if full siblings occur at a very low rate and at a substantial rate,  
464 respectively. It should decrease with an increasing rate of full siblings and also a decreasing rate of  
465 duplicates to minimize both  $\alpha$  and  $\beta$  errors. Unfortunately, however, sample genetic structure is  
466 usually unknown, and is the focus of an individual identification study.

467 Results from pairwise approaches are much less accurate, as expected from the simulation  
468 results and from the fact that this dataset has a large number of individuals and contains very large  
469 full sib families. Take the LR method as an example. Among the 1821186 possible dyads, 153  
470 dyads within a life stage were identified as single individuals, yielding a  $\beta$  error rate 3.73 times  
471 larger than that of LC. A serious problem with the pairwise approach is self-conflicted inferences.  
472 Fig. 7 shows the pairwise relationships among 5 MGs inferred by LR. Obviously, these pairwise  
473 inferences are incompatible. The higher the level of individual replications, the more severe will be  
474 the problem of pairwise approaches.

475

## 476 **Discussion**

477 Although the mismatch method is the simplest and the most widely applied method for marker-  
478 based individual identification in molecular ecology, it has unfortunately several weaknesses and as  
479 a result is the least accurate method. First, the fixed threshold, typically  $T_m = 1$  or 2, is arbitrary. It is  
480 too small when the number of loci or/and the mistyping rate is high, resulting in too many ghost  
481 individuals. It is too large when the number of loci and mistyping rate are very low, or/and close  
482 relatives are frequent. It is also too rigid and inappropriate for pairs of MGs having missing data at

483 different numbers of loci. These properties of MM have been well recognized, and have led to the  
484 suggestion that the fewest possible number of markers that have sufficient power for individual  
485 identification should be used to avoid excessive mismatches and exclusions (Waits *et al.* 2001;  
486 Creel *et al.* 2003). In reality, the markers used in individual identification can be highly variable in  
487 polymorphisms and mistyping rates, and the background relationship (e.g. sibship and parentage)  
488 structure of a sample can also be highly variable. It is difficult for any fixed value of  $T_m$  to cater for  
489 all scenarios. Second, the mismatch method fails to use the mismatch information efficiently. Two  
490 single locus genotypes are regarded matched when they are identical, and mismatched when they  
491 have either one or both alleles different. Obviously, mismatched genotypes give more evidence of  
492 distinct individuals when they have both alleles rather than a single allele different. This kind of  
493 information is however unused by the mismatch method. Third, the mismatch method treats all loci  
494 equally, while they can be highly heterogeneous in information (polymorphism) and noise  
495 (mistyping) contents. The method simply counts the number of mismatches, regardless of the loci at  
496 which the mismatches occur. Obviously, mismatched MGs give more support for distinct  
497 individuals when the mismatches occur at loci with lower mistyping rates or/and higher  
498 polymorphisms.

499 I showed in this study that an optimal  $T_m$  value can be calculated by simulations,  
500 accommodating the number of loci, the mistyping and missing data rates and the allele frequencies  
501 at each locus. The optimal  $T_m$  gives an unbiased estimate of the average number of mismatches  
502 between truly duplicated MGs (Fig. 1). Applying the optimal  $T_m$  value determined by simulations,  
503 the mismatch method has almost a constant accuracy independent of the number of loci ( $L$ , Fig. 1).  
504 If the fixed  $T_m=2$  were applied, the accuracy would have decreased rapidly with  $L$  when it is larger  
505 than 20 because of the excessive  $\alpha$  errors. Compared with other methods, however, the mismatch  
506 method using the optimal  $T_m$  value is still the least accurate for various parameter combinations  
507 considered in the simulations (Figures 1-5). It is impossible for the mismatch method to use as  
508 much marker information (e.g. mistyping rates, allele frequencies) and thus to have a comparable  
509 accuracy as the other methods.

510 Relatedness method has rarely been used in individual identifications. However, recently  
511 Ringler *et al.* (2015) showed that it is much more accurate than mismatch method for analysing  
512 their frog data. Relatedness method has several advantages over mismatch method. First, it uses  
513 allele frequency information. For example, two matched genotypes lend more support for a single  
514 individual if they are rare (i.e. containing rare alleles) than if they are common. Second, relatedness  
515 calculation is robust to the presence of mistypings. The relatedness estimates between close  
516 relatives (such as duplicates and full sibs) are reduced only slightly by assuming perfect data when

517 they are actually not (Wang 2007). My simulations conducted for different parameter combinations  
518 confirm Ringler *et al.*'s conclusion that relatedness method is more accurate than mismatch method.  
519 Importantly, relatedness method is statistically consistent. With an increasing number of markers,  
520 even though they suffer from genotyping errors, the method always becomes increasingly more  
521 accurate (Fig. 1).

522 Like the mismatch method, the relatedness method requires a threshold value,  $T_r$ , to  
523 determine the relationship between two MGs. The dyad is concluded to be a single and two distinct  
524 individuals when their relatedness is greater and not greater than  $T_r$ , respectively. Ideally, the  
525 optimal  $T_r$  value that minimises both  $\alpha$ - and  $\beta$ -errors should be obtained by considering the  
526 frequencies of DPs and the most close relationship (e.g. FS) in the sample. These frequencies are  
527 usually unknown, and the close relatives are most often full siblings and parent offspring, both  
528 having an expected relatedness of 0.5. Using the average relatedness of first degree relatives (0.5)  
529 and DPs (1.0) as threshold, I obtained  $T_r=0.75$  and used it in simulated data analysis. This value is  
530 slightly smaller than the value obtained by Ringler *et al.* (2015), 0.8, in their frog data analysis.  
531 They derived this value from the estimated relatedness of the 20 juvenile-adult pairs identified as  
532 identical from morphology. In practice, whenever a sufficient number of known duplicated  
533 individuals are available, Ringler *et al.*'s approach should be followed to determine a dataset  
534 specific  $T_r$ . Otherwise, a generic  $T_r=0.75$  can be used in individual identification, bearing in mind  
535 that the optimal value depends on the relative frequencies of DPs and the most close relationships as  
536 well as genotyping error rates and other factors (e.g. number and polymorphisms of markers).  
537 Further study (via simulation or meta-analysis) is needed to investigate the optimal  $T_r$  and the  
538 factors affecting it.

539 Individual identification from a pairwise likelihood relationship (LR) analysis does not  
540 require a threshold. We calculate the probability of two MGs conditional on each of a number of  
541 candidate relationships, and the probability is the likelihood of the relationship. We then simply  
542 select the relationship that has the maximal likelihood as the best estimate. Similar to the  
543 considerations in relatedness analysis, we choose FS, HS, PO as well as DP as the candidate  
544 relationships. Unlike relatedness analysis, however, relationship inference is highly susceptible to  
545 mistypings, and a relationship (such as PO and DP) can be erroneously excluded because of  
546 genotyping errors. For this reason, I used the error models of Wang (2004) to account for false  
547 alleles (FA) and allelic dropouts (ADO) separately. Overall, LR method performs slightly better  
548 than, but is more susceptible to mis-specified FA and ADO rates (Fig. 5) than relatedness (RL)  
549 method. Recently, researchers have recognized the ubiquitous presence of mistypings and its large  
550 impact on many downstream analyses (Bonin *et al.* 2004; Pompanon *et al.* 2005), and increasingly

551 quantified and reported mistyping rates. Therefore, the application of LR method should  
552 increasingly less limited by the lack of mistyping information.

553 A common problem of the above three methods is that they consider each pair of MGs in  
554 isolation of others. These pairwise approaches waste marker information and thus have low  
555 accuracy. For an example, let's consider  $n+1$  MGs which are identical except for a single locus at  
556 which there are  $n$  heterozygous genotypes {A,B} and 1 single homozygous genotype {A,A}. These  
557  $n+1$  MGs would support the hypothesis that they come from a single individual rather than two  
558 distinct individuals when ADO or FA rate is not very small at the locus showing different  
559 genotypes and when  $n$  is large. The larger the value of  $n$ , the greater is the support. However, this  
560 support is much reduced when only 2 genotypes are considered as in the pairwise approach.  
561 Confirming the reasoning, Fig. 3 shows contrasting behaviours between LC and the 3 pairwise  
562 approaches. As the replication level increases, LC becomes more accurate, while pairwise  
563 approaches either remain the same accuracy or become less accurate. As a result, the difference in  
564 accuracy between LC and pairwise approaches increases with an increasing level of individual  
565 replication.

566 Another common problem of the above three methods is that they frequently yield self-  
567 incompatible inferences, as shown in a real example (Fig. 7). In practice, what one needs is usually  
568 the MG clusters, each corresponding to a single individual. This means one has to go through these  
569 pairwise inferences and assemble them into individual clusters. The process is not only tedious  
570 because of so many pairwise inferences, in the order of  $N(N-1)/2$  where  $N$  is the number of MGs,  
571 but may fail to produce valid clusters.

572 Although the simulated data contain half sibs, they were analysed by LC by assuming  
573 monogamy for both sexes such that half sibs were not inferred. This is because half sibs are not of  
574 our interest and also have much smaller effect on individual identification than full sibs.  
575 Abandoning half sib inferences can however speed up the computation substantially and is thus  
576 especially favourable for a simulation study. In analysis of real data, it is also safe to ignore half  
577 sibs when individual identification is the purpose of analysis.

578 Highly polymorphic microsatellites from noninvasive samples have been used in identifying  
579 individuals and estimating population size (Waits & Leberg 2000; Creel *et al.* 2003; McKelvey &  
580 Schwartz 2004). It is anticipated that SNPs would become more and more widely used in the near  
581 future because of their low cost and high automation in genotyping. Although much less  
582 informative (usually biallelic) individually than microsatellites, SNPs can be genotyped at a much  
583 larger number of loci at ease and collectively they can be much more informative. My simulations

584 (Fig. 2) showed that all four methods can use markers of widely different polymorphisms in  
585 individual identification. However, the performance of the mismatch method, even when improved  
586 by using an optimal  $T_m$ , deteriorates rapidly with a decreasing marker polymorphism because of the  
587 excessive false identifications of individuals ( $\alpha$  errors). The problem is much more severe if a fixed  
588  $T_m$  value is used. In contrast, the LC method is especially more accurate than other methods with  
589 many markers of low polymorphisms. Using a number of 160 SNPs, each having 2 alleles and a  
590 mistyping rate of 0.05, LC has an overall accuracy several orders higher than other methods.

591         Except for the mismatch method that uses a fixed  $T_m$  value, allele frequencies are needed in  
592 inferring duplicates. Usually these frequencies are unavailable in practice, but can be estimated  
593 from the genotype data under the assumption that all homologous genes (within and between  
594 individuals) at a locus are non-identical by descent. The assumption is obviously violated when  
595 some sampled individuals are duplicated or otherwise related. However, violation of the assumption  
596 does not seem to cause a serious problem for all 4 methods investigated in this study, even when  
597 individual replication level is high (Fig. 3). The LC method implemented in Colony program does  
598 have the ability to account for the inferred genetic structure in refining allele frequency estimates,  
599 and has been proved to be effective in improving pedigree reconstruction when the families  
600 included in a sample are highly unbalanced in sizes (Wang 2004; Wang & Santure 2009).

601         My simulations assumed an outbred species without inbreeding. However, inbreeding or  
602 population structure could have some effects on the inference of duplicates. While it is not  
603 immediately apparent how to extend the MM, RL and LR methods to account for inbreeding, the  
604 LC method in Colony can actually accommodate inbreeding, including selfing, in relationship  
605 inference (Wang & Santure 2009). It can estimate inbreeding and relationship jointly. However,  
606 how much improvement in individual duplicate inference can be gained by allowing for inbreeding  
607 is yet to be investigated in a further study.

608         The simulation results for less related family structures, 40(1, 1) and 16(1, 1, 3), are similar  
609 to those shown in Figures 1-5. All methods become slightly more accurate, because full sib  
610 frequency is smaller and thus the chance of  $\alpha$  errors is reduced. Overall across all simulated datasets  
611 and the empirical dataset, the LC method performs substantially better than the pairwise approaches,  
612 and is highly recommended for use in practice.

613         The LC and LR methods are implemented and added to the computer program COLONY  
614 version 2.0.5.3, which was used in analysing the data shown in this paper. The program is  
615 downloadable from the website <http://www.zsl.org/science/software/colony>.

616

617 **Acknowledgements**

618 I am grateful to Eva Ringler for her stimulating paper and for sending me her frog data analysed by  
619 this MS. I thank her and four anonymous referees for valuable comments which have helped  
620 improving the MS.

621

622 **References**

- 623 Alacs EA, Georges A, FitzSimmons NN, Robertson J (2010) DNA detective: a review of molecular  
624 approaches to wildlife forensics. *Forensic Science, Medicine, and Pathology*, **6**, 180-194.
- 625 Anderson EC, Dunham KK (2008) The influence of family groups on inferences made with the  
626 program Structure. *Molecular Ecology Resources*, **8**, 1219-1229.
- 627 Bonin A, Bellemain E, Bronken Eidesen P, Pompanon F, Brochmann C, Taberlet P (2004) How to  
628 track and assess genotyping errors in population genetics studies. *Molecular Ecology*, **13**, 3261-  
629 3273.
- 630 Creel S, Spong G, Sands JL, Rotella J, Zeigle J, Joe L *et al.* (2003) Population size estimation in  
631 Yellowstone wolves with error-prone noninvasive microsatellite genotypes. *Molecular*  
632 *Ecology*, **12**, 2003–2009.
- 633 Escaravage N, Questiau S, Pornon A, Doche B, Taberlet P (1998) Clonal diversity in a  
634 *Rhododendron ferrugineum* L. (Ericaceae) population inferred from AFLP markers. *Molecular*  
635 *Ecology*, **7**, 975–982.
- 636 Galpern P, Manseau M, Hettinga P, Smith K, Wilson P (2012) ALLELEMATCH: an R package for  
637 identifying unique multilocus genotypes where genotyping error and missing data may be  
638 present. *Molecular Ecology Resources*, **12**, 771-778.
- 639 Goodnight KF, Queller DC (1999) Computer software for performing likelihood tests of pedigree  
640 relationship using genetic markers. *Molecular Ecology*, **8**, 1231-1234.
- 641 Halkett F, Simon JC, Balloux F (2005) Tackling the population genetics of clonal and partially  
642 clonal organisms. *Trends in Ecology and Evolution*, **20**, 194-201.
- 643 Li CC, Weeks DE, Chakravarti A (1993) Similarity of DNA fingerprints due to chance and  
644 relatedness. *Human Heredity*, **43**, 45–52.
- 645 Lynch M (1988) Estimation of relatedness by DNA fingerprinting. *Molecular Biology and*  
646 *Evolution*, **5**, 584–599.
- 647 Luikart G, Ryman N, Tallmon DA, Schwartz MK, Allendorf FW (2010) Estimation of census and  
648 effective population sizes: the increasing usefulness of DNA-based approaches. *Conservation*  
649 *Genetics*, **11**, 355-373.

650 McKelvey KS, Schwartz MK (2004) Genetic errors associated with population estimation using  
651 non-invasive molecular tagging: problems and new solutions. *Journal of Wildlife Management*,  
652 **68**, 439-448.

653 Palsbøll P, Allen J, Bérubé M, Clapham PJ, Feddersen TP, Hammond PS, Hudson RR *et al.* (1997)  
654 Genetic tagging of humpback whales. *Nature*, **388**, 767–769.

655 Pompanon F, Bonin A, Bellemain E, Taberlet P (2005) Genotyping errors: causes, consequences  
656 and solutions. *Nature Review Genetics*, **6**, 847–850.

657 Ringler E, Mangione R, Ringler M (2015) Where have all the tadpoles gone? Individual genetic  
658 tracking of amphibian larvae until adulthood. *Molecular Ecology Resources*. (In press)

659 Rodríguez-Ramilo ST, Wang J (2012) The effect of close relatives on unsupervised Bayesian  
660 clustering algorithms in population genetic structure analysis. *Molecular Ecology Resources*, **12**,  
661 873-884.

662 Schwartz MK, Tallmon DA, Luikart G (1998) Review of DNA-based census and effective  
663 population size estimators. *Animal Conservation*, **1**, 293–299.

664 Taberlet P, Waits L, Luikart G (1999) Non-invasive genetic sampling: look before you leap. *Trends*  
665 *in Ecology and Evolution*, **14**, 323–327.

666 Wagner AP, Creel S, Kalinowski ST (2006) Estimating relatedness and relationships using  
667 microsatellite loci with null alleles. *Heredity*, **97**, 336-345.

668 Waits L, Leberg PL (2000) Biases associated with population estimation using molecular tagging.  
669 *Animal Conservation*, **3**, 191-199.

670 Waits L, Luikart G, Taberlet P (2001) Estimating the probability of identity among genotypes in  
671 natural populations: cautions and guidelines. *Molecular Ecology*, **10**, 249–256.

672 Wang J (2004) Sibship reconstruction from genetic data with typing errors. *Genetics*, **166**, 1963-  
673 1979.

674 Wang J (2006) Informativeness of genetic markers for pairwise relationship and relatedness  
675 inference. *Theoretical Population Biology*, **70**, 300-321.

676 Wang J (2007) Triadic IBD coefficients and applications to estimating pairwise relatedness.  
677 *Genetics Research*, **89**, 135–153.

678 Wang J (2014) Marker-based estimates of relatedness and inbreeding coefficients: an assessment of  
679 current methods. *Journal of Evolutionary Biology*, **27**, 518-530.

680 Wang J, Santure AW (2009) Parentage and sibship inference from multilocus genotype data under  
681 polygamy. *Genetics*, **181**, 1579–1594.

682 Wilberg MJ, Dreher BP (2004) genecap: a program for analysis of multilocus genotype data for  
683 non-invasive sampling and capture-recapture population estimation. *Molecular Ecology Notes*, **4**,  
684 783-785.

685

686

---

687 J. Wang is interested in developing population genetics models and methods of analysis of  
688 empirical data to address issues in evolutionary and conservation biology.  
689

---

690

691 **Data accessibility**

692 The simulated genotype datasets can be found on Dryad: Dryad doi: [10.5061/dryad.2q3qh](https://doi.org/10.5061/dryad.2q3qh)

693 The frog dataset of Ringler et al. (2015) can be found on Dryad:

694 <http://dx.doi.org/10.5061/dryad.db800>

695 The computer program used in the simulated and empirical data analysis, Colony, is available

696 <http://www.zsl.org/science/software/colony>.

697



## Figure Captions

699 **Fig. 1** Effect of the number of markers. The upper graph plots the observed ( $x$  axis) and threshold  
 700  $T_m$  ( $y$  axis) numbers of mismatches of each simulated duplicated MG dyad for different number of  
 701 markers ( $L$ ). The lower graph plots the error rate ( $\gamma$ ) of 4 individual identification methods as a  
 702 function of the number of markers ( $L$ ). The four methods are mismatch (MM), relatedness (RL),  
 703 likelihood relationship (LR), and likelihood clustering (LC). For both graphs, the parameters used in  
 704 the simulations are family structure 4(1, 2, 3, 4, 10),  $K_l=10$  and  $\varepsilon_{l1}=\varepsilon_{l2}=\varepsilon_{l3}=0.05$  for each locus  $l$   
 705 ( $=1, 2, \dots, L$ ),  $\lambda=0.5$ .

706 **Fig. 2**  $\alpha$ -,  $\beta$ - and total-error rates of 4 individual identification methods as a function of the number  
 707 of alleles per marker ( $K$ ). The four methods are mismatch (MM), relatedness (RL), likelihood  
 708 relationship (LR), and likelihood clustering (LC). The parameters used in the simulations are family  
 709 structure 4(1, 2, 3, 4, 10),  $L=160, 80, 40, 20, 10, 5$  when  $K=2, 3, 5, 9, 17$  and 33 respectively,  
 710  $\varepsilon_{l1}=\varepsilon_{l2}=\varepsilon_{l3}=0.05$  for each locus, and  $\lambda=0.5$ .

711 **Fig. 3** Error rate ( $\gamma$ ) of 4 individual identification methods as a function of the extent of individual  
 712 replication ( $\lambda$ ). The four methods are mismatch (MM), relatedness (RL), likelihood relationship  
 713 (LR), and likelihood clustering (LC). The parameters used in the simulations are family structure  
 714 4(1, 2, 3, 4, 10),  $L=10$ ,  $K=10$ ,  $\varepsilon_{l1}=\varepsilon_{l2}=\varepsilon_{l3}=0.05$  for each locus,  $\lambda$  ( $x$  axis) varies between 0 (no  
 715 replication) to 3.2 (an individual is on average replicated by 3.2 times).

716 **Fig. 4** Error rate ( $\gamma$ ) of 4 individual identification methods as a function of the rate of mistyping and  
 717 missing data at a locus ( $\varepsilon$ ). The four methods are mismatch (MM), relatedness (RL), likelihood  
 718 relationship (LR), and likelihood clustering (LC). The parameters used in the simulations are family  
 719 structure 4(1, 2, 3, 4, 10),  $L=20$ ,  $K=10$ ,  $\lambda=0.5$ ,  $\varepsilon_{l1} \equiv \varepsilon_{l2} \equiv \varepsilon_{l3}$  ( $x$  axis) varies between 0 (perfect data  
 720 with no mistyping and no missing data) to 0.16 at each locus  $l$ .

721 **Fig. 5** Error rate ( $\gamma$ ) of 4 individual identification methods as a function of the assumed rate of  
 722 mistyping at a locus ( $\hat{\varepsilon}$ ). The four methods are mismatch (MM), relatedness (RL), likelihood  
 723 relationship (LR), and likelihood clustering (LC). The parameters used in the simulations are family  
 724 structure 4(1, 2, 3, 4, 10),  $L=20$ ,  $K=10$ ,  $\lambda=0.5$ ,  $\varepsilon_{l1} \equiv \varepsilon_{l2} = 0.1$ ,  $\varepsilon_{l3} = 0.05$ . The analysis was  
 725 conducted assuming a mistyping rate ( $x$  axis) of  $\hat{\varepsilon}_{l1} \equiv \hat{\varepsilon}_{l2}$  between 0 (perfect data with no mistyping)  
 726 to 0.2 at each locus  $l$ .

727 **Fig. 6** Distributions of the numbers of loci with missing data (A-D) and mismatches (E-H) between  
 728 two MGs in the frog dataset. Row 1 (A and E) is for the 60 dyads in the 20 inferred trios that  
 729 contain morphologically identified juvenile-adult pairs, row 2 (B and F) is for the 106 other dyads

730 inferred to be duplicates, row 3 (C and G) is for the inferred 16620 full sib dyads, and row 4 (D and  
731 H) is for the inferred 1804400 non-full-sib dyads.

732 **Fig. 7** The relationships among 5 MGs inferred by LR for the frog dataset. In the 5 MG names, “it”,  
733 “m” and “ij” indicate tadpoles, male adults, and juveniles respectively. Two MGs are inferred by  
734 LR to come from a single individual if they are linked by a line, and from distinct individuals if they  
735 are not linked by a line.

736