

1 KIMGENS: A novel method to estimate kinship in organisms with
2 mixed haploid diploid genetic systems robust to population structure

3

4 Yen-Wen Wang^{1,*}, Cécile Ané^{1,2}

5

6 ¹ Department of Botany, University of Wisconsin-Madison, Madison, WI, 53706, USA

7 ² Department of Statistics, University of Wisconsin-Madison, Madison, WI, 53706, USA

8 * Correspondence author

9

10

11 Abstract

12 Motivation:

13 Kinship estimation is necessary for evaluating violations of assumptions or testing certain
14 hypotheses in many population genomic studies. However, kinship estimators are usually
15 designed for diploid systems and cannot be used in populations with mixed haploid diploid
16 genetic systems. The only estimators for different ploidies require datasets free of population
17 structure, limiting their usage.

18

19 Results:

20 We present KIMGENS, an estimator for kinship estimation among individuals of various
21 ploidies, that is robust to population structure. This estimator is based on the popular KING-
22 robust estimator but uses diploid relatives of the individuals of interest as references of
23 heterozygosity and extends its use to haploid-diploid and haploid pairs of individuals. We
24 demonstrate that KIMGENS estimates kinship more accurately than previously developed
25 estimators in simulated panmictic, structured and admixed populations, but has lower accuracy
26 when the individual of interest is inbred. KIMGENS also outperforms other estimators in a
27 honeybee dataset. Therefore, KIMGENS is a valuable addition to a population geneticist's
28 toolbox.

29

30 Availability and Implementation:

31 KIMGENS and its association simulation tool are implemented and available open-source at
32 <https://github.com/YenWenWang/HapDipKinship>.

33

34 Contact:

35 Yen-Wen Wang

36 Email: ywang883@wisc.edu

37

38 Introduction

39 Kinship estimation is crucial to the evaluation of assumption violations (such as when
40 estimating population nucleotide diversity) or to testing various ecological or evolutionary
41 hypotheses (e.g., kin selection). However, kinship estimators for whole genome datasets are
42 mainly developed for human populations (Ramstetter *et al.*, 2017). Although these estimators
43 have been widely used in non-human systems, their applications are restricted to diploid-only
44 populations. Nonetheless, a large portion of life forms show plasticity in ploidy (Otto and
45 Gerstein, 2008), which is not accounted for in these estimators. Many plants (e.g. ferns, mosses),
46 fungi (e.g. mushrooms) and algae (e.g. sea lettuces) have complex, multistage life cycles,
47 perform alternation of generations and form haploid structures independent of their diploid
48 counterpart (Brown and Casselton, 2001; John, 1994). Furthermore, most Hymenopterans (e.g.
49 bees and ants), Thysanopterans (e.g. thrips) and some other invertebrates (e.g. some spider mites
50 and rotifers), have an arrhenotokous haplodiploidy system, where males are haploid and females
51 are diploid (Cruickshank and Thomas, 1999; Normark, 2003). Because of the widely present
52 mixed ploidy life-forms, it is crucial to develop estimators that can estimate kinship between
53 individuals with different ploidy levels.

54 Two marker-based estimators have been developed specifically to estimate relatedness
55 among individuals of different ploidy, including Huang2014, a method-of-moments (MOM)
56 estimator, and Huang2015, a maximum likelihood (ML) estimator (Huang *et al.*, 2015). These

57 estimators can thus be used for genome sequencing data directly. In addition, some classical
58 estimators can be extended to estimate relatedness between different ploidies. For example, two
59 kinship estimators, Loiselle1995 and Ritland1996 (Loiselle *et al.*, 1995; Ritland, 1996), are
60 adapted and implemented in the program PolyRelatedness to resolve inequivalent ploidy (Huang
61 *et al.*, 2015). All estimators mentioned above are capable of using multi-allelic loci, which allow
62 them to take advantage of a diversity of genetic markers (e.g., microsatellites). However, these
63 estimators require allele frequencies at each locus in the population, which may not be available
64 in some studies due to sampling strategies (Hahn, 2019). In addition, relying on allele
65 frequencies essentially assumes a population free of stratification. Therefore, the estimators do
66 not account for cryptic population structure, which can result in overestimating in kinship
67 (Manichaikul *et al.*, 2010).

68 To remove the requirement of no population structure, we built on the KING-robust
69 estimator by Manichaikul *et al.*, (2010). We extended the estimator's use to haploid-diploid pairs
70 and named this extension exKING-robust. The exKING-robust estimator uses the heterozygosity
71 of the individuals in the pair of interest as a diversity estimate for background identity-by-
72 descent (IBD). Next, we developed KIMGENS (Kinship Inference for Mixed GENetic Systems),
73 which instead uses the heterozygosity of relatives of the individuals of interest, allowing
74 estimating kinship for diploid, haploid-diploid and haploid pairs of individuals. We showed the
75 estimators are robust to population structure. KIMGENS also performs relatively well under
76 admixture, but can underestimate kinship if an individual is inbred.

77

78 Materials and Methods

79 We aim to develop a simple kinship estimator that applies to haploid, haploid-diploid and
80 diploid pairs of individuals and is robust to population structure. KING-robust's strategy is
81 useful for developing a novel kinship estimator for haploid-diploid pairs of individuals. But the
82 strategy cannot apply to haploid pairs directly because it requires the number of heterozygotic
83 sites of an individual to estimate expected heterozygosity ($2pq$) in the ancestral subpopulation of
84 an individual.

85 To resolve this issue, we propose a two-step approach: (1) We extend KING-robust to
86 obtain a haploid-diploid kinship analysis and to identify a set of diploid relatives for each
87 individual; and (2) for two individuals of interest, i and j , we use their diploid relatives from step
88 1 to estimate mean heterozygosity for this pair and to modify their kinship estimate. In the
89 following sections we will first describe a haploid-diploid kinship estimator. Next, we will
90 demonstrate the modification to KING-robust and haploid-diploid kinship estimators for using
91 related individuals k . Then, we describe the haploid kinship analysis. Lastly, we will evaluate the
92 performance of these estimators with simulations and a biological dataset (on honeybee).

93

94 Kinship estimation in haploid-diploid pairs of individuals: exKING-robust

95 The kinship coefficient ϕ_{ij} , originally termed correlation coefficient of two individuals i
96 and j , is defined as the probability that two randomly sampled alleles from two individuals are
97 identical-by-descent (IBD) (Lange, 1997; Malécot, 1948). In this section, we derive an estimator
98 for the kinship of a pair of individuals i_d and j_h , where i_d is diploid and j_h is haploid. $\phi_{i_d j_h}$ can be
99 calculated with:

100

101

$$\phi_{i_d j_h} = (1/2) \pi_{1i_d j_h}$$

102 (1)

103 where $\pi_{ni_dj_h}$ denotes the probability of individuals i_d and j_h sharing n alleles being IBD. The
104 probability of individual i_d being homozygotic and not in identical-by-state (IBS) with individual
105 j_h at a site can be calculated with:

$$\Pr(AA, a \text{ or } aa, A) = p^2q\pi_{0i_dj_h} + pq^2\pi_{0i_dj_h} = pq\pi_{0i_dj_h}$$

106
107
108 (2)

109 and the probability of individual i_d being heterozygotic and in IBS at the allele in individual j_h at
110 a site can be calculated with:

$$\Pr(Aa, a \text{ or } Aa, A) = \Pr(Aa) = 2pq.$$

111
112
113 (3)

114 Because i_d and j_h share either 0 or 1 allele by descent (j_h being haploid),

$$\pi_{0i_dj_h} + \pi_{1i_dj_h} = 1.$$

115
116
117 (4)

118 With equation (1), we derive

$$\phi_{i_dj_h} = (1/2)(1 - \pi_{0i_dj_h}).$$

119
120
121 (5)

122 We can combine equation (5) with equation (3) to get

$$\phi_{i_dj_h} = \frac{1}{2} - \frac{\Pr(AA, a \text{ or } aa, A)}{2pq}.$$

123
124

125 (6)

126 Because only individual i_d is heterozygotic, the expected genome-wide heterozygosity,
127 $\sum_m 2p_m q_m$, can be estimated with $N_{Aa}^{(i_d)} / M_{i_d j_h}$ (Manichaikul *et al.*, 2010), where $N_{Aa}^{(i_d)}$ is the
128 number of heterozygotic sites in individual i_d and $M_{i_d j_h}$ is the number of sites with non-missing
129 data in both i_d and j_h . $\sum_m \Pr(AA, a \text{ or } aa, A)_m$ can be estimated with $N_{AA, a \text{ or } aa, A} / M_{i_d j_h}$, where
130 $N_{AA, a \text{ or } aa, A}$ is the number of sites where individual i is homozygotic but not in IBS with
131 individual j_h . Therefore, kinship between individuals i_d and j_h can be estimated with:

132

$$\widehat{\Phi}_{i_d j_h} = \frac{1}{2} - \frac{N_{AA, a \text{ or } aa, A}}{N_{Aa}^{(i_d)}},$$

133

134 (7)

135 which constitute our exKING-robust estimator for a haploid-diploid pair.

136

137 Methods for using related individuals to estimate pq

138 The KING-robust extension, including KING-robust (Manichaikul *et al.*, 2010) for

139 diploid pairs and exKING-robust for haploid-diploid pairs (7), relies on $N_{Aa}^{(i_d)}$ (and $N_{Aa}^{(j_d)}$).

140 However, in haploid pairs, we do not have the luxury of using the heterozygosity of individuals

141 of interest, so we develop a different estimator, KINGENS, which uses “heterozygosity

142 references” to estimate pq . The accuracy of $\widehat{\Phi}_{i_h j_h}$ highly depends on the choice of references. To

143 accurately capture heterozygosity, references should come from the same subpopulation as

144 individuals i_h and j_h . Identification of appropriate references can be done by examining kinship

145 estimates between the individual i_h and j_h and the potential references. Since some individuals

146 may deviate from Hardy-Weinberg equilibrium (HWE) in subpopulations, choosing a single

147 reference from the relatives of either individual i_h or j_h may result in using an inbred or admixed
 148 product, biasing the estimate. So, we choose two sets of individuals $K(i_h, t)$ and $K(j_h, t)$, which are
 149 related to either one of the two individuals of interest, given a kinship threshold t . Every
 150 individual k in $K(i_h, t)$ or $K(j_h, t)$ is used as a heterozygosity reference for an intermediate kinship
 151 estimate, $\hat{\phi}_{i_h j_h}^{[k]}$. Then, we calculate two medians of intermediate kinship estimates, one from
 152 $K(i_h, t)$ and one from $K(j_h, t)$. Finally, we take the mean of these two medians as our final estimate
 153 $\hat{\phi}_{i_h j_h}$. We explain this estimation procedure below in detail.

154 For more generality, we introduce this procedure for diploid individuals to modify the
 155 exKING-robust estimators as well. For two diploids i_d and j_d and for a reference individual k , we
 156 define the intermediate kinship estimate $\hat{\phi}_{i_d j_d}$ as

$$158 \quad \hat{\phi}_{i_d j_d}^{[k]} = \frac{1}{2} - \frac{1}{4} \frac{4N_{AA,aa \text{ or } aa,AA}^{(i_d, j_d)} - 2N_{Aa, Aa}^{(i_d, j_d)} + N_{Aa}^{(i_d)} + N_{Aa}^{(j_d)}}{N_{Aa}^{(k)}}.$$

159 (8)

160 Next, for an individual x , we consider its references to be the set $K(x, t)$ of diploid individuals that
 161 share kinship with x greater than a given threshold t (including x itself if x is diploid), based on
 162 the exKING-robust kinship estimate. Finally, we define the KIMGENS estimate as follows:

$$164 \quad \hat{\phi}_{i_d j_d} = \frac{1}{2} \left(\text{Median}_{k \in K(i_d, t)} \{ \hat{\phi}_{i_d j_d}^{[k]} \} + \text{Median}_{l \in K(j_d, t)} \{ \hat{\phi}_{i_d j_d}^{[l]} \} \right).$$

165 (9)

166 Parenthesized superscripts denote the individuals with which sequences are compared to derive
 167 the number of sites with a particular pattern, and bracketed superscripts denote the individuals
 168 used for intermediate kinship estimates. Note, (8) corresponds to equation (11) in Manichaikul *et*

169 *al.*, (2010) for k taken to be either i_d or j_d , whichever has the smallest $N_{Aa}^{(k)}$. The innovation here
 170 is to consider the median of kinship estimates, and to use close relatives (not just i_d or j_d) to
 171 approximate heterozygosity at the denominator.

172 Using the same idea, we use (7) to define the intermediate kinship estimate between a
 173 pair of diploid and haploid individuals, given a reference individual k as:

174

$$\hat{\phi}_{i_d j_h}^{[k]} = \frac{1}{2} - \frac{N_{AA,a \text{ or } aa,A}^{(i_d, j_h)}}{N_{Aa}^{(k)}} \quad (10)$$

177 and for a haploid-diploid pair we define the KIMGENS estimate as:

$$\hat{\phi}_{i_d j_h} = \frac{1}{2} \left(\text{Median}_{k \in K(i_d, t)} \{ \hat{\phi}_{i_d j_h}^{[k]} \} + \text{Median}_{l \in K(j_h, t)} \{ \hat{\phi}_{i_d j_h}^{[l]} \} \right). \quad (11)$$

180

181 When calculating a $\hat{\phi}_{i_d j_d}^{[k]}$ (or $\hat{\phi}_{i_d j_h}^{[k]}$), there are three individuals involved: i_d, j_d (or j_h) and
 182 k . The amount of missing data are not the same in these three individuals. So, we only consider
 183 the sites that are non-missing in all three individuals for each $\hat{\phi}_{i_d j_d}^{[k]}$ (or $\hat{\phi}_{i_d j_h}^{[k]}$).

184

185 Kinship estimation in haploid pairs of individuals

186 Under the same definition for kinship, in haploid pairs, the kinship coefficient $\phi_{i_h j_h}$ can
 187 be calculated with:

188

$$189 \quad \phi_{i_h j_h} = \pi_{1i_h j_h}.$$

190 (12)

191 The probability of individuals i_h and j_h not in IBS at a site can be calculated with:

192

$$\Pr(A, a \text{ or } a, A) = 2pq\pi_{0i_hj_h}. \quad (13)$$

195 Because

$$\pi_{0i_hj_h} + \pi_{1i_hj_h} = 1,$$

197

$$\phi_{i_hj_h} = 1 - \frac{\Pr(A, a \text{ or } a, A)}{2pq}. \quad (14)$$

200 Using the same strategy described above, an intermediate kinship for haploid pairs of
201 individuals can be estimated using a reference diploid individual k with:

202

$$\hat{\phi}_{i_hj_h}^{[k]} = 1 - \frac{N_{A,a \text{ or } a,A}^{(i_hj_h)}}{N_{Aa}^{(k)}} \quad (15)$$

205 and the KIMGENS estimate for a haploid-haploid pair is defined as

206

$$\hat{\phi}_{i_hj_h} = \frac{1}{2} \left(\text{Median}_{k \in K(i_h,t)} \{ \hat{\phi}_{i_hj_h}^{[k]} \} + \text{Median}_{l \in K(j_h,t)} \{ \hat{\phi}_{i_hj_h}^{[l]} \} \right). \quad (16)$$

208

210 Simulations

211 To assess the performance of these estimators, we simulated panmictic, structured and
212 admixed populations of species with haplodiploid or diploid genetic system. For panmictic
213 populations, the allele frequency of each site was simulated from a uniform distribution between
214 0.1 and 0.9, $U(0.1,0.9)$. The genotypes for starter individuals (those without known parents) in
215 pedigree simulations were drawn from the allele frequency. For structured and admixture
216 populations, the allele frequencies of three subpopulations were simulated following the Balding-
217 Nichols model from a panmictic ancestral population (allele frequency drawn from $U(0.1,0.9)$).
218 The Wright's $F_{st} (\theta_k)$ of the subpopulations was set to 0.05, 0.15 and 0.25. In structured
219 populations, each family was drawn from a random subpopulation. To simulate admixture,
220 Conomos's strategy was used (Conomos *et al.*, 2016). In pedigree simulations, the ancestry
221 proportions of the founders were drawn independently from either of two Dirichlet distributions:
222 $Dir(6, 2, 0.3)$ and $Dir(2, 6, 0.3)$, and the genotypes of the founders were drawn from the ancestry
223 and allele frequencies of the three subpopulations.

224 While simulating pedigrees, nine different scenarios were simulated 1000 times each.
225 The scenarios differed by four factors: (1) the number of independent SNP sites: 20k or 100k, (2)
226 genetic system: arrhenotokous haplodiploidy or diploidy, (3) population structure: panmictic,
227 structured or admixture and (4) pedigrees (Supplementary Table 1). Overall, 100k SNP sites
228 were simulated unless when the estimators being compared included those implemented in
229 PolyRelatedness, in which case 20k SNPs were simulated. All simulations are under
230 haplodiploidy unless otherwise noted. First, to evaluate the performance of exKING-robust and
231 KIMGENS, we simulated a single large family (Supplementary Figure 1) from a panmictic
232 population (scenario 1). To compare the performance with that of previously published
233 estimators, we simulated 11 families (Supplementary Figure 2) from a panmictic or structured

234 population (scenarios 2 and 3). To explore the performance of the estimators under admixture,
235 we simulated the single large family (Supplementary Figure 1) or 11 families (Supplementary
236 Figure 2) from an admixed population with 100k or 20k sites (scenarios 4 and 5). To understand
237 how different estimators perform on inbred products, we simulated five inbreeding families
238 (Supplementary Figure 3) and ten unrelated individuals (so PolyRelatedness can estimate allele
239 frequency more accurately) from a panmictic diploid or haplodiploid population (scenarios 6 and
240 7). Lastly, to explore the use of different thresholds (t), we simulated a new family
241 (Supplementary Figure 4) with twenty unrelated diploid individuals in a structured or admixed
242 population (scenarios 8 and 9).

243 We estimated pairwise kinships for all individuals using different estimators and
244 extracted the estimates of the pairs of interest. To convert relatedness (calculated by the
245 estimators implemented in PolyRelatedness) to kinship, the relatedness estimates for diploid
246 pairs were divided by two and those for haploid-haploid and haploid-diploid pairs were divided
247 by one. To summarize the estimation, we calculated the bias ($\sum(\hat{\phi}_i - \phi_{true})/n$) and root-mean-
248 square error (RMSE; $\sqrt{\sum(\hat{\phi}_i - \phi_{true})^2/n}$) of each estimator. For KIMGENS, the threshold t
249 was arbitrarily set to 0.1, except when exploring different thresholds, in which case t was set to
250 either 0.1 or 0. Inbreeding coefficients were calculated from pedigrees with the R package
251 kinship2 (Sinnwell *et al.*, 2014).

252

253 Biological data

254 In addition to simulations, we used a honeybee dataset which was originally collected for
255 estimating crossover rate (Liu *et al.*, 2015). This dataset includes three monogynous colonies
256 (one queen per colony). One queen (diploid) and multiple drones (haploids) were sampled from

257 all three colonies. Six additional workers (diploids) were sampled from one of the colonies. Also,
258 three drones were sequenced twice. We therefore expect that from a single colony, (1) the drones
259 and the queen share a kinship of 0.5, (2) the workers and the queen share a kinship of 0.25, (3)
260 the drones share a kinship of 0.5 with each other, (3) the workers share a kinship of 0.375 (full-
261 siblings) or 0.125 (half-siblings) with each other, (4) the drones and workers share a kinship of
262 0.25, and (5) the two sequences from the same drone share a kinship of 1 with each other.

263 The genomic raw reads were downloaded from NCBI and mapped to reference genome
264 (GCF_000002195.4) with BWA mem ver. 0.7.17. Duplicated reads were filtered with samtools
265 ver. 1.9 and SNPs were called with bcftools ver. 1.9. To avoid identifying SNPs due to indels,
266 we applied four filters: (1) the repetitive regions identified by RepeatMasker, (2) sites with read
267 depth higher than 1.3X mean depth or lower than 0.75X mean depth, (3) sites with minor allele
268 frequency lower than 0.01 and (4) sites that are called heterozygous in any haploid individuals
269 (drones). All filtered SNPs (N=1,008,683) were used to estimate kinship without LD correction.
270 As the previous section, for KIMGENS, the threshold t was arbitrarily set to 0.1. To compare
271 KIMGENS with other published estimators, we sampled one every twenty SNPs to avoid
272 segmentation faults for these other estimators.

273

274 Results and discussion

275 Evaluation of the methods under a panmictic population

276 We simulated a single large family in haplodiploidy with 100k sites in a panmictic
277 population (Supplementary Table 1 and Supplementary Figure 1) and compared the performance
278 of exKING-robust and KIMGENS, for each ploidy level of the individuals of interests (diploid,
279 haploid-diploid or haploid). For diploid pairs of individuals, the estimates from both exKING-

280 robust and KIMGENS are accurate with no bias and small RMSE (Figure 1, Supplementary
281 Table 2). The same is observed for haploid-diploid and haploid pairs (Figure 1, Supplementary
282 Table 2). The variance of estimates is usually higher in haploid pairs and lower in diploid pairs,
283 likely due to the fact that the amount of allelic data is halved in haploid compared to diploid
284 individuals, causing a precision decrease.

285

286 Comparison with previous methods in panmictic and structured populations

287 We compared the performance of KIMGENS with other relatedness estimators
288 implemented in the package PolyRelatedness, including Huang2014 (MOM), Huang2015
289 (MLE), Ritland1996 and Loiselle1995 (Huang *et al.*, 2015; Loiselle *et al.*, 1995; Ritland, 1996).
290 We simulated 11 families from a panmictic or structured population (Supplementary Table 1 and
291 Supplementary Figure 2).

292 In a panmictic population, the performance of KIMGENS outcompetes all other
293 estimators in terms of the overall RMSE and bias (Figure 2A, Supplementary Table 3).
294 KIMGENS performs slightly worse than Huang2015 only when the true kinship is zero. In a
295 structured population, KIMGENS again outperforms all other methods when the true kinship is
296 not zero (Figure 2B and Supplementary Table 4). However, KIMGENS has the highest RMSE
297 and absolute bias when true kinship is zero, and Huang2015 has the lowest. In the structured
298 population simulation, there is 1/3 chance that two unrelated individuals are from different
299 subpopulations. The fixed variants in the subpopulations increase the homozygotic differences
300 between two individuals and hence lower the kinship estimates between unrelated samples using
301 KING-robust-based strategies (Manichaikul *et al.*, 2010). Also, note that Huang2015 performs
302 the best when the true kinship equals zero in both conditions (Supplementary Table 4). This is

303 likely because Huang2015 uses a maximum likelihood strategy searching for IBD on a parameter
304 space, where the lower bound of the parameter space is zero (Huang *et al.*, 2015). If negative
305 kinship is a concern, one can enforce a lower bound of zero for all estimators. In our simulations,
306 this would vastly improve the bias and RMSE of KIMGENS when the true kinship is zero,
307 without affecting the performance when the true kinship is positive.

308

309 Estimates in an admixed population

310 Although estimating kinship in an admixed population is not the goal of this project, we
311 explored the robustness of KIMGENS in admixed population. First, we simulated the single-
312 family pedigree in an admixed population to evaluate the performance of exKING-robust and
313 KIMGENS (Supplementary Table 1 and Supplementary Figure 1). Like previous reports on
314 KING-robust (Conomos *et al.*, 2016), the accuracy of both estimators drops compared to the
315 estimates in a panmictic population because the individuals from a single family may have
316 different ancestries (Supplementary Table 5 and Supplementary Figure 5). Similarly to the
317 panmictic population simulation, the estimates in haploid and diploid pairs of individuals have
318 slightly higher and lower RMSE, respectively. KIMGENS also performs slightly better than
319 exKING-robust in terms of RMSE and bias.

320 We further compared KIMGENS with aforementioned estimators using the 11-family-
321 pedigree (Supplementary Table 1 and Supplementary Figure 2). In this admixture simulation,
322 KIMGENS has lower absolute bias than all other estimators but a RMSE slightly higher than
323 Huang 2015 (Figure 2C and Supplementary Table 6). The relatively high RMSE is also driven
324 by the lower kinship estimates on unrelated individuals due to the same reasons discussed in the
325 last section.

326

327 Estimates on inbred individuals

328 Like KING-robust, KIMGENS is not designed to calculate kinship in inbred populations,
329 but we explored its performance for inbred individuals by simulating five families and ten
330 additional unrelated individuals (half male and half female) in a panmictic population in diploid
331 or haplodiploid genetic system (Supplementary Table 1 and Supplementary Figure 3). The
332 unrelated individuals were included because all of the methods being compared require
333 population allele frequency, which is estimated with the sampled individuals in this study. In a
334 diploid genetic system, KIMGENS performs slightly better than other estimators overall in terms
335 of both RMSE and bias (Figure 3A and Supplementary Table 7). However, the RMSE and
336 absolute bias increase when the individual inbreeding coefficients of the two individuals
337 increase, and the increasing rate is faster than other kinship estimators, such as Huang2015 and
338 Loiselle1995.

339 In a haplodiploid genetic system, KIMGENS has a relatively high overall RMSE and
340 absolute bias (Figure 3B and Supplementary Table 8), so we broke down the results by the
341 ploidy of pairs and individual inbreeding coefficients. For diploid pairs, the behavior of
342 KIMGENS is very similar to that in the diploid simulation (Supplementary Table 8 and
343 Supplementary Figure 6A). KIMGENS outperforms all other estimators overall, but the accuracy
344 drops when the individual inbreeding coefficient increases. For haploid pairs, all individuals
345 have zero inbreeding coefficients and KIMGENS also performs better than other estimators
346 (Supplementary Table 8 and Supplementary Figure 6C). However, for haploid-diploid pairs,
347 KIMGENS performs worse than other estimators overall except for exKING-robust and also
348 when individual inbreeding coefficients are higher than zero (Supplementary Table 8 and

349 Supplementary Figure 6B). Like diploid pairs, the kinship estimates decrease under a higher
350 degree of inbreeding (Supplementary Table 8 and Supplementary Figure 6B). This correlation is
351 essentially the same underestimation as when the individuals in the pair of interest are from two
352 different subpopulations.

353

354 Performance on biological data

355 The kinship estimates on the honeybee dataset using KIMGENS are close to expectations
356 in all within-colony relationships except for workers-workers (Figure 4; Supplementary Table 9).
357 Kinship estimates between workers can be clustered into three groups: 0.125, 0.25 and 0.375.
358 While estimates at 0.125 and 0.375 between workers are expected for full-siblings and half-
359 siblings, estimates at 0.25 are unexpected, but is likely the result of paternal relatedness. For
360 example, the kinship between two workers whose fathers are siblings equals 0.25. In addition,
361 we found that individuals between different colonies share a considerably high degree of kinship
362 (mean= 0.08) (Figure 4). We hypothesize that the high degree of kinship is derived from true
363 background relatedness due to breeding management of the bee farm. The background
364 relatedness may also contribute to the positive biases of kinship estimates between workers
365 within a single family—that is, the putative half-siblings may have distantly related fathers
366 (Supplementary Table 9).

367 The performance of KIMGENS on the subsampled dataset is similar to that on the whole
368 dataset, while all other estimators underestimate kinships on the subsampled dataset when the
369 true kinships are lower than 1 (Supplementary Figure 7 and Supplementary Table 10). This
370 observation supports the usage of KIMGENS on biological datasets.

371

372 Choice of the kinship threshold t

373 The only parameter in KIMGENS is the kinship threshold t used to define the set of
374 relatives for heterozygosity referencing. Without inbreeding in a panmictic or structured
375 population, the threshold t should not affect the accuracy as long as it is positive. However,
376 admixture may elevate unrelated individuals kinship (Figure 2C), and inbreeding can lower the
377 heterozygosity in some individuals, so the choice of t needs to be taken into consideration.

378 The threshold t can affect two factors: the accuracy of heterozygosity referencing and the
379 number of reference individuals. In order to identify diploid individuals that can represent the
380 heterozygosity of the individuals of interest (in a same population), one should consider a higher
381 t , but a higher t may reduce the number of heterozygosity references. A lower number of
382 references should not directly affect the accuracy of kinship estimation. For example, although
383 there are 46 drones in the honeybee dataset and only nine females, the female honeybees are
384 closely related to the drones and hence can represent the heterozygosity of the drones well.
385 However, if a dataset includes numerous inbreeding events, using a high t may result in only
386 referencing the inbred individuals, and the kinship estimation will be inaccurate, so a lower t
387 should be considered. In some extreme cases, it may be that KIMGENS cannot estimate kinship,
388 if zero diploid relatives are available for a pair of haploid individuals. In this case, one may
389 choose zero for threshold t , which allows referencing unrelated individuals from the same
390 subpopulation for heterozygosity.

391 To explore the effect of using unrelated individuals as heterozygosity references on
392 kinship estimation for haploid pairs of individuals, we simulated one family and twenty unrelated
393 diploids in a structured or admixed population (Supplementary Table 1 and Supplementary
394 Figure 4). We first estimated kinships regularly using a threshold t at 0.1 with KIMGENS. Then,

395 we removed all diploid individuals within the family, leaving the unrelated diploids in the dataset
396 only and estimated kinship using a threshold t at 0. In a structured population, both the RMSE
397 and absolute bias are higher when using non-relatives compared to using relatives, but the
398 estimates are still accurate (Supplementary Table 11 and Supplementary Figure 8A). In admixed
399 populations, the RMSE and bias are also higher when using non-relatives; however, the RMSE is
400 noticeably higher when the true kinship is over zero (Supplementary Table 11 and
401 Supplementary Figure 8B). This is likely due to the complex ancestry of the individuals in
402 admixed populations, resulting in non-relatives providing a significantly worse heterozygosity
403 reference than relatives. Of note, in 1% of the pairs of interest, there were no diploid individuals
404 with a kinship above 0 (but less than 0.1) with either haploid individual in the pair, so no kinship
405 was estimated for these pairs. Using a negative threshold t can resolve the issue, but these
406 estimates should be interpreted with extra caution.

407

408 Conclusions

409 Here we present new kinship estimators for mixed haploid-diploid populations that are
410 robust to population structure. We demonstrate the accuracy of KIMGENS in panmictic,
411 structured and admixed simulated populations as well as in a biological dataset. Simulations and
412 biological datasets indicate that KIMGENS performs better than previously developed kinship
413 estimators, but one may choose to use previously developed kinship estimators when the dataset
414 contains many multiallelic loci or individuals of interest with high degree of inbreeding
415 coefficient. The methods are implemented in an R package available on github
416 (<https://github.com/YenWenWang/HapDipKinship>) for researchers studying population
417 genomics in mixed ploidy systems.

418

419 Data Availability Statement

420 The data underlying this article are available in Open Science Repository at
421 <https://dx.doi.org/10.17605/OSF.IO/EP6MF>. The datasets were derived from NCBI sequence
422 read archive, accession SRP043350.

423

424 Acknowledgements

425 The authors thank Dr. Anne Pringle and Dr. Jacob Golan for reviewing the manuscript,
426 and Dr. Cameron Currie and Dr. Sean Schoville for advice on biological datasets.

427 Funding

428 Y.-W.W supported by an E. K. and O. N. Allen fellowship, a Tulipa et Paeonia RA
429 support and a Taylor-Vinje research award provided by the department of Botany, and a MSA
430 graduate fellowship provided by the Mycological Society of America. C.A was supported in part
431 by a H. I. Romnes faculty fellowship provided by the University of Wisconsin-Madison Office
432 of the Vice Chancellor for Research and Graduate Education with funding from the Wisconsin
433 Alumni Research Foundation.

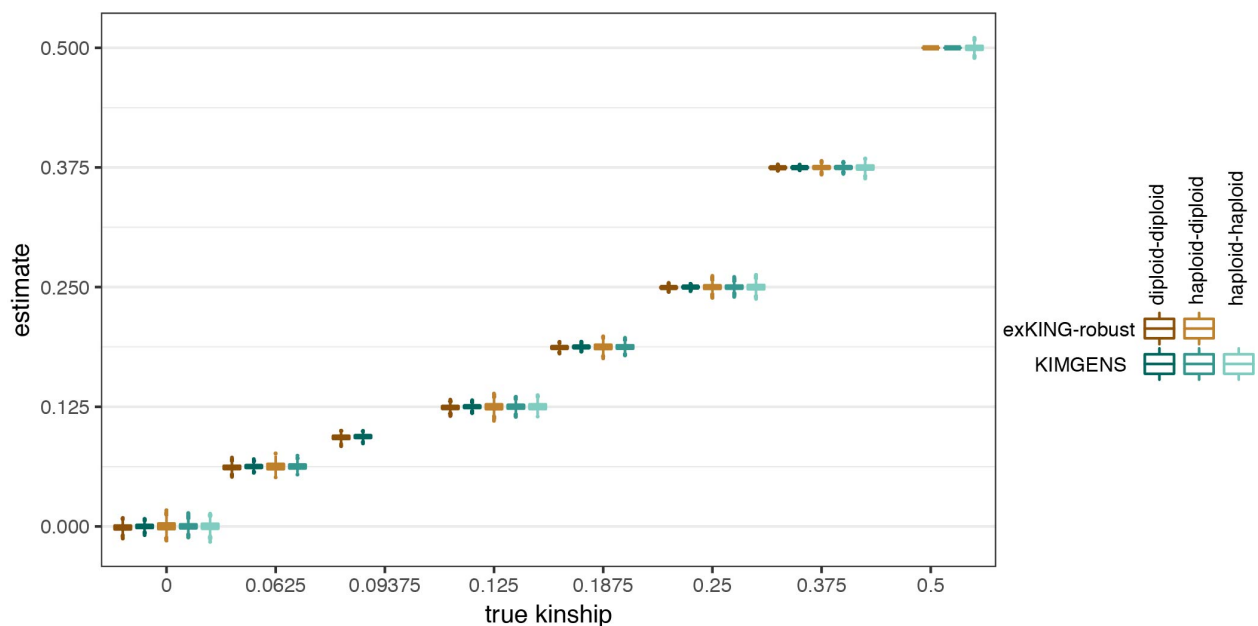
434

435 References

- 436 Brown,A.J. and Casselton,L.A. (2001) Mating in mushrooms: Increasing the chances but
437 prolonging the affair. *Trends Genet.*, **17**, 393–400.
- 438 Conomos,M.P. *et al.* (2016) Model-free estimation of recent genetic relatedness. *Am. J. Hum.*
439 *Genet.*, **98**, 127–148.
- 440 Cruickshank,R.H. and Thomas,R.H. (1999) Evolution of haplodiploidy in dermanyssine mites
441 (Acari: Mesostigmata). *Evolution (N. Y.)*, **53**, 1796–1803.
- 442 Hahn,M. (2019) Experimental design. In, *Molecular population genetics*. Sinauer Associates,
443 Inc., Sunderland, Massachusetts, pp. 25–42.
- 444 Huang,K. *et al.* (2015) Estimating pairwise relatedness between individuals with different levels
445 of ploidy. *Mol. Ecol. Resour.*, **15**, 772–784.
- 446 John,D.M. (1994) Alternation of generations in algae: Its complexity, maintenance and
447 evolution. *Biol. Rev. Camb. Philos. Soc.*, **69**, 275–291.
- 448 Lange,K. (1997) Mathematical and statistical methods for genetic analysis. Springer, New York,
449 NY, USA.
- 450 Liu,H. *et al.* (2015) Causes and consequences of crossing-over evidenced via a high-resolution
451 recombinational landscape of the honey bee. *Genome Biol.*, **16**, 15.
- 452 Loiselle,B.A. *et al.* (1995) Spatial genetic structure of a tropical understory shrub, *Psychotria*
453 *officinalis* (Rubiaceae). *Am. J. Bot.*, **82**, 1420–1425.
- 454 Malécot,G. (1948) Les mathématiques de l’hérédité. Masson et Cie, Paris, France.
- 455 Manichaikul,A. *et al.* (2010) Robust relationship inference in genome-wide association studies.
456 *Bioinformatics*, **26**, 2867–2873.
- 457 Normark,B.B. (2003) The evolution of alternative genetic systems in insects. *Annu. Rev.*

- 458 *Entomol.*, **8**, 397–423.
- 459 Otto,S.P. and Gerstein,A.C. (2008) The evolution of haploidy and diploidy. *Curr. Biol.*, **18**,
- 460 R1121–R1124.
- 461 Ramstetter,M.D. *et al.* (2017) Benchmarking relatedness inference methods with genome-wide
- 462 data from thousands of relatives. *Genetics*, **207**, 75–82.
- 463 Ritland,K. (1996) Estimators for pairwise relatedness and individual inbreeding coefficients.
- 464 *Genet. Res.*, **67**, 175–185.
- 465 Sinnwell,J.P. *et al.* (2014) The kinship2 R package for pedigree data. *Hum. Hered.*, **78**, 91–93.
- 466

467 Figures

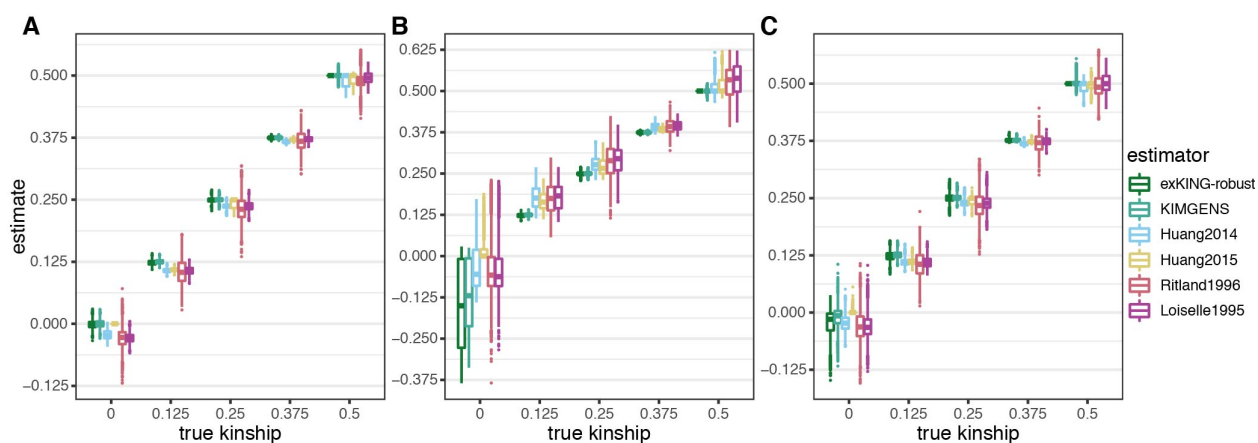


468

469 Figure 1. Distribution of kinship estimates of exKING-robust and KIMGENS in a panmictic

470 population. Boxplots show the median, first and third quartiles, and range of each distribution.

471

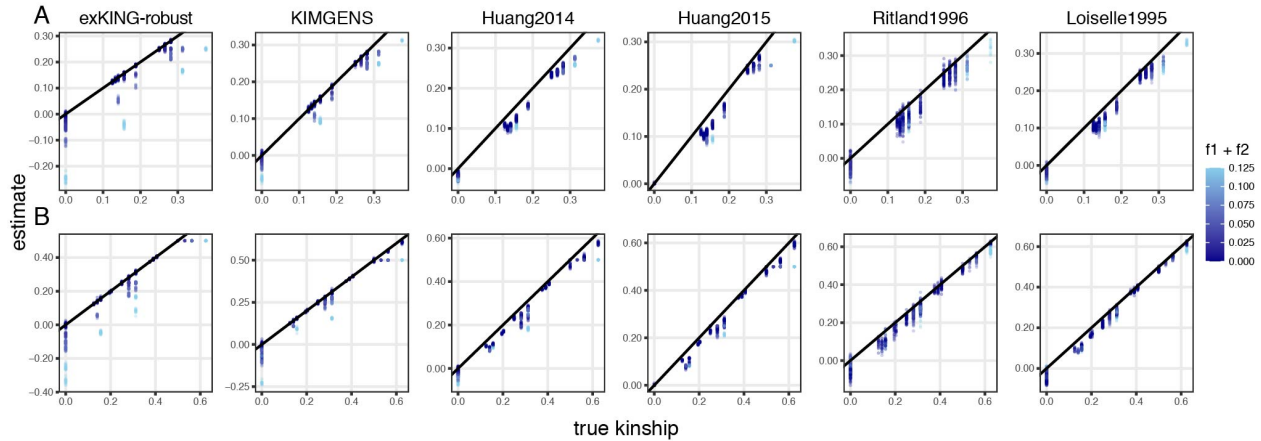


472

473 Figure 2. Kinship estimates of different estimators in panmictic (A), structured (B) and admixed

474 (C) populations.

475



476

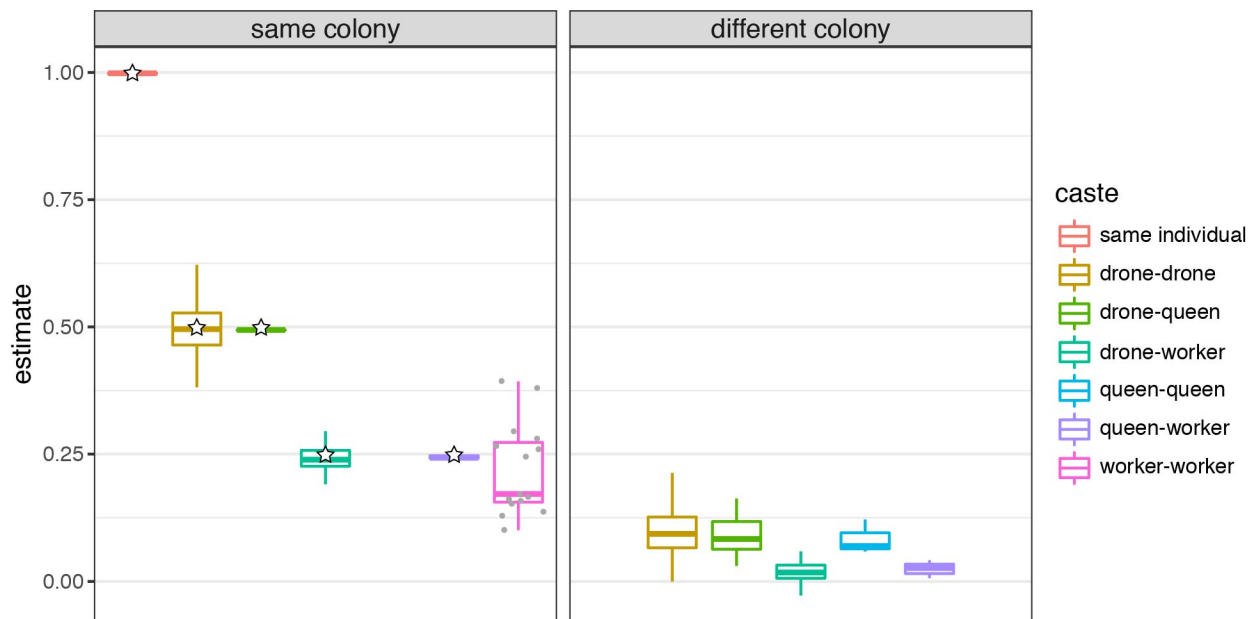
477 Figure 3. Performance of different estimators on inbred diploid (A) and haplodiploid (B)

478 populations. A thousand points were chosen randomly to be presented on each plots. Diagonal

479 line: estimated kinship = true kinship. $f_1 + f_2$: the sum of inbreeding coefficients of the two

480 individuals in a pair of interest.

481



482

483 Figure 4. Kinship estimates from KIMGENS on three honeybee colonies. Gray dots indicate

484 each kinship estimate between workers. The “same individual” category consists of pairs of

485 sequence data sets from the drones that were sequenced twice. Stars indicate expected kinship.

