

1 Running title (optional): **Choosing the optimal population for GWAS**

2

3 Core ideas (3-5 impact statements, 85 char max for each)

4 - Genome-wide association studies with mixture populations are expected to improve the

5 detection power of novel genes due to the increase of the sample size although the influence of

6 population structure is a concern.

7 - When a quantitative trait nucleotide (QTN) is polymorphic in a target population, a

8 combination of the target population and a population with higher diversity than the target

9 population improves the detection power of the QTN.

10 - We found that the fixation index (F_{ST}) and the expected heterozygosity (H_e) were strongly

11 related to the detection power of QTNs.

12 - Germplasm collections which have been already sequenced/genotyped are useful for improving

13 the detection power of GWAS without any addition of sequence costs by using a subset of them

14 with a target population.

15

16 **Choosing the optimal population for a genome-wide association study: a simulation using**
17 **whole-genome sequences from rice**

18 Kosuke Hamazaki, Hiromi Kajiya-Kanegae, Masanori Yamasaki, Kaworu Ebana, Shiori Yabe,
19 Hiroshi Nakagawa and Hiroyoshi Iwata*

20

21 Affiliations:

22 K. Hamazaki, H. Kajiya-Kanegae and H. Iwata, Department of Agricultural and Environmental
23 Biology, Graduate School of Agricultural and Life Sciences, The University of Tokyo, 1-1-1
24 Yayoi, Bunkyo-ku, Tokyo 113-8657, Japan; H. Kajiya-Kanegae, current address: Research
25 Center for Agricultural Information Technology, National Agriculture and Food Research
26 Organization, 3-1-1 Kannondai, Tsukuba, Ibaraki 305-8517, Japan; M. Yamasaki, Food
27 Resources Education and Research Center, Graduate School of Agricultural Science, Kobe
28 University, 1348 Uzurano, Kasai, Hyogo 675-2103, Japan; K. Ebana, Genetic Resources Center,
29 National Agriculture and Food Research Organization, 2-1-2 Kannondai, Tsukuba, Ibaraki 305-
30 8602, Japan; S. Yabe, Institute of Crop Science, National Agriculture and Food Research
31 Organization, 2-1-2 Kannondai, Tsukuba, Ibaraki 305-8518, Japan; H. Nakagawa, Institute for
32 Agro-Environmental Sciences, National Agriculture and Food Research Organization, 3-1-3
33 Kannondai, Tsukuba, Ibaraki 305-8604, Japan. *Corresponding author ([aiwata@mail.ecc.u-](mailto:aiwata@mail.ecc.u-tokyo.ac.jp)
34 [tokyo.ac.jp](mailto:aiwata@mail.ecc.u-tokyo.ac.jp))

35

36 Abbreviations:

37 AUC, area under the curve; CDR, correct detection rate; FDR, false discovery rate; FN, false
38 negative FP, false positive; F_{ST} , the fixation index; GWAS, genome-wide association study; H,
39 high; HM, higher-middle; H_e , the expected heterozygosity; L, low; LD, linkage disequilibrium;
40 LM, lower-middle; M, middle MAF, minor allele frequency QTL, quantitative trait loci; QTN,
41 quantitative trait nucleotide; ROC, receiver operating characteristic; SNP, single nucleotide
42 polymorphism; TN, true negative; TP, true positive.

43

44

ABSTRACT

45 A genome-wide association study (GWAS) needs to have a suitable population. The
46 factors that affect a GWAS, e.g. population structure, sample size, and sequence analysis and
47 field testing costs, need to be considered. Mixture populations containing subpopulations of
48 different genetic backgrounds may be suitable populations. We conducted simulation
49 experiments to see if a population with high genetic diversity, e.g., a diversity panel, should be
50 added to a target population, especially when the target population harbors small genetic
51 diversity. The target population was 112 accessions of *Oryza sativa* subsp. *japonica*, mainly
52 developed in Japan. We combined the target population with three populations that had higher
53 genetic diversities. These were 100 *indica* accessions, 100 *japonica* accessions, and 100
54 accessions with various genetic backgrounds. The results showed that the GWAS power with a
55 mixture population was generally higher than with a separate population. Also, the GWAS
56 optimal population varied depending on the fixation index F_{ST} of the quantitative trait nucleotide
57 (QTN) and its polymorphism of QTN in each population. When a QTN is polymorphic in a
58 target population, a target population combined with a higher diversity population improves the

59 QTN detection power. Investigating F_{ST} and the expected heterozygosity H_e as factors
60 influencing the detection power, we showed that SNPs with high F_{ST} or low H_e are less likely to
61 be detected by GWAS with mixture populations. Sequenced/genotyped germplasm collections
62 can improve the GWAS detection power by using a subset of them with a target population.

63

64

INTRODUCTION

65 Recently, as genome sequencing costs have continued to decrease (Metzker, 2010), the
66 whole-genome sequences of a large number of cultivars/lines have become available for major
67 crop species, such as rice (Li et al., 2014; Wang et al., 2018). A genome-wide association study
68 (GWAS) based on whole-genome sequences can more efficiently and accurately identify genes
69 that control important agronomic traits than previous methods (Koboldt et al., 2013; Ott et al.,
70 2015; Yano et al., 2016).

71 It is important to prepare an appropriate population to be analyzed when attempting to detect
72 candidate genes using GWAS techniques. For example, to avoid potential false positives caused
73 by population stratification/structure, a GWAS population should be selected that results in low
74 stratification (Begum et al., 2015; Yano et al., 2016). However, if such a population is selected as
75 an analytical population for a GWAS, the sample size may be limited and the detection power of
76 the GWAS will decrease (Korte and Farlow, 2013). Therefore, when designing an appropriate
77 GWAS population, one should be aware of the trade-off relationship between population
78 stratification and sample size.

79 When preparing the population to be analyzed, the factors that directly affect the GWAS
80 results, such as population structure, sample size, and the sequence analysis and cultivation

81 testing costs, need to be considered. In recent years, the whole-genome sequences of a large
82 number of cultivars/lines have become publicly available due to highly efficient sequencing
83 analyses and database enrichment. The publically available whole-genome sequence data will
84 improve GWASs and could enable researchers to avoid the costs of sequencing analyses. For
85 example, in rice, "The 3,000 Rice Genomes Project" (Li et al., 2014; Wang et al., 2018) by the
86 International Rice Research Institute (IRRI) is a well-known whole-genome sequence dataset
87 that is available in the "Rice SNP-Seek Database" (Alexandrov et al., 2015; Mansueto et al.,
88 2016; 2017). Therefore, an appropriate GWAS population could potentially utilize existing
89 public sequence data.

90 A mixture population obtained by mixing subpopulations with different genetic backgrounds
91 could also potentially be used in a GWAS. An advantage of using such a mixture population is
92 that it should improve the detection of causal variants by increasing the sample size. Conversely,
93 a GWAS with a mixture population may suffer from large numbers of false positives caused by
94 the population structure. Although a mixed effect model that suppresses the influence of the
95 population structure has been proposed (Yu et al., 2006), such a mixture population has rarely
96 been analyzed by a GWAS.

97 An actual data analysis of rice using whole-genome sequences showed that the detection
98 power of a GWAS improved when *Oryza sativa* subsp. *japonica* and *Oryza sativa* subsp. *indica*
99 populations were combined (Misra et al., 2017). Furthermore, the identification of new rice
100 genes using a GWAS and populations with extremely high genetic diversities has also been
101 previously reported (Zhao et al., 2011). Conversely, it has been reported that the genetic
102 differentiation between subpopulations in a population with high genetic diversity could cause a
103 reduction in the power of a GWAS (Huang et al., 2012). Therefore, real data studies have been

104 inconsistent about whether mixture populations or populations with high genetic diversities
105 should be used in a GWAS. However, these previous studies mostly analyzed actual data, and
106 there have been no theoretical simulation studies that have considered the possibility of using a
107 mixture population in a GWAS. Furthermore, no previous studies have discussed which kinds of
108 populations should be mixed to improve the GWAS detection power or which kinds of
109 populations are most appropriate for a GWAS. Therefore, in this study, we conducted simulation
110 experiments to see whether adding a population with a high genetic diversity compared to a
111 target population (e.g., adding a diversity panel to a target population) is appropriate, especially
112 when the genetic diversity of the target population is small.

113

114

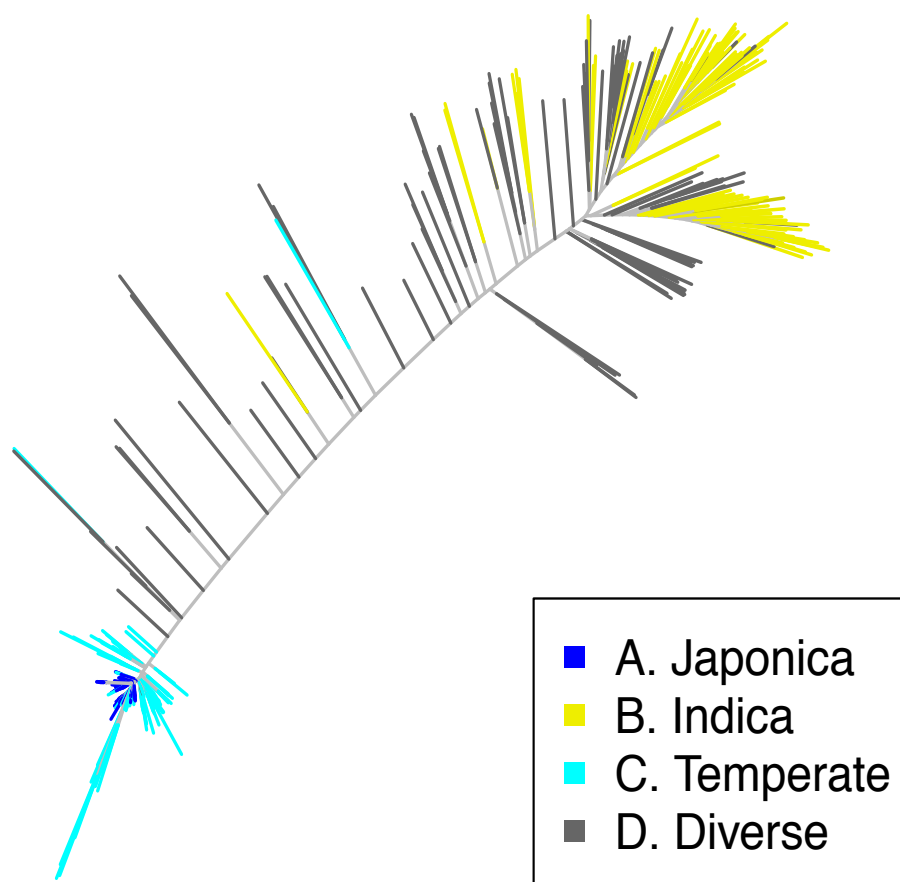
MATERIALS AND METHODS

115

Materials (populations used in the GWAS)

116 In this study, 112 accessions of *Oryza sativa* subsp. *japonica* (referred to as “A”), which
117 were accessions that had mainly been developed in Japan, were used as a target population with
118 low genetic diversity (Yabe et al., 2016). We used the following three populations selected from
119 “The 3,000 Rice Genomes Project” (J. Y. Li et al., 2014), i.e., 100 accessions of *Oryza sativa*
120 subsp. *indica* (referred to as “B”), 100 accessions of *Oryza sativa* subsp. *japonica* (referred to as
121 “C” or temperate), and 100 accessions of *Oryza sativa* with various genetic backgrounds
122 (referred to as “D” or diverse), as populations with higher diversities than the target population
123 (Table S1 in Supplemental File 1). The process used to select each of the 100 accessions is
124 described in Supplemental File 9. One accession (IRIS ID: IRIS 313-11868) was duplicated in
125 populations B and D. Among populations B, C, and D, the B population was the most

126 differentiated from A, whereas C was the most similar to A. Population D contained subsp.
127 *indica*, subsp. *japonica*, and *aus*, and *aromatic* rice accessions, which meant that the D
128 population had the highest genetic diversity. Fig. 1 is an unrooted phylogenetic tree that shows
129 the genetic relationships among accessions belonging to populations A, B, C, and D.
130



131

132 **Fig. 1. Unrooted phylogenetic tree plot for four non-mixture populations.**

133 Unrooted phylogenetic tree plot for the four non-mixture populations, which consisted of 112

134 accessions of *japonica* (A), 100 accessions of *indica* (B), 100 accessions of temperate *japonica*

135 (C), and 100 diverse accessions (D) with neighbor-joining method.

136

137 The genetic relationships among the accessions were estimated by the neighbor-joining (NJ)
138 method (Saitou and Nei, 1987) using the R package “ape” version 5.3 (Paradis et al., 2004). The
139 genetic distances were estimated according to the Jukes and Cantor (1969) model. In addition to
140 these four populations, we synthesized three populations by combining population A with
141 populations B, C, or D. The mixture populations A + B, A + C, and A + D were named “E”, “F”,
142 and “G”, respectively. We compared the QTN detection power the GWAS when the seven non-
143 mixture (A, B, C, and D) and mixture populations (E, F, and G) were used.

144

145

146

Genotype data

147 Whole genome sequencing data were available for the accessions (Jarquin et al., 2019).
148 Details about the DNA extraction and whole genome sequencing techniques are provided in a
149 previous report (Jarquin et al., 2019). The data sets deposited in the DDBJ Sequence Read
150 Archive (SRA106223, ERA358140, DRA000158, DRA000307, DRA000897, DRA000927,
151 DRA007273, DRA007256, and DRA008071) were reanalyzed. We processed the whole-genome
152 sequence data as follows so that they could be used in the GWAS. Adapters and low-quality
153 bases were removed from paired reads using the Trimmomatic v0.36 program (Bolger et al.,
154 2014). The preprocessed reads were aligned using Os-Nipponbare-Reference-IRGSP-1.0
155 (Kawahara et al., 2013) and the bwa-0.7.12 mem algorithm with the default options (H. Li, 2012).
156 The binary alignment map (BAM) files deposited in the Rice SNP-Seek database were also
157 reanalyzed. Single nucleotide polymorphism (SNP) calling was based on alignments determined
158 using the Genome Analysis Toolkit (GATK), 3.7-0-gcfe6b67 (McKenna et al., 2009; Auwera et

159 al., 2014) and Picard package V2.5.0 (<http://broadinstitute.github.io/picard>). The mapped reads
160 were realigned using RealignerTargetCreator and indelRealigner in the GATK software. The
161 SNPs and InDels were called at the population level using the UnifiedGenotyper in GATK and
162 the -glm BOTH option. We extracted bi-allelic sites in all the accessions from the variants using
163 VCFtools version 0.1.13 (Danecek et al., 2011). Then, imputations were imputed using Beagle
164 version 4.1 (Browning and Browning, 2016). Finally, we analyzed the SNPs with minor allele
165 frequencies (MAFs) ≥ 0.05 in each population. In the analysis, the genotypes were represented as
166 -1 (homozygous of the reference allele), 1 (homozygous of the alternative allele) or 0
167 (heterozygous of the reference and alternative alleles). Out of all the whole-genome sequence
168 polymorphisms, only the SNPs on chromosome 1 were analyzed. The number of SNPs on
169 chromosome 1 in each population is shown in Table 1.

170

171 **Table 1. Number of SNPs and the diversity level of non-mixture and mixture populations.**

	Population name	Number of accessions	Number of SNPs	Diversity level[†]
A.	Japonica	112	72,110	263.095
B.	Indica	100	427,943	660.416
C.	Temperate japonica	100	135,665	362.649
D.	Diverse	100	647,731	798.646
E.	A + B	212	633,507	803.064
F.	A + C	212	151,675	334.606
G.	A + D	212	684,774	859.678

172 [†] Diversity level is the index that was used to indicate the degree of genetic diversity and is
173 described in the “Degree of genetic diversity index” section below.

174

175

176

Generating phenotype data

177 Phenotypic data were simulated using the following formula:

$$178 \quad \mathbf{y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \mathbf{X}_3\beta_3 + \mathbf{u} + \mathbf{e}, \quad (\text{Eq. 1})$$

179 where \mathbf{y} is the vector that represents the simulated phenotypic values for all 411 accessions; \mathbf{X} is
180 the design matrix representing the genotypes of three quantitative trait nucleotides (QTNs) with
181 scores -1, 0, or 1; $\boldsymbol{\beta} = [\beta_1 \ \beta_2 \ \beta_3]^T$ is the vector representing the effects of the three QTNs, \mathbf{u}
182 is the vector for polygenetic effects, and \mathbf{e} is the residuals vector. Three QTN-SNPs whose MAF
183 was equal to or larger than 0.05 in all 411 accessions (672,923 SNPs in total) were randomly
184 selected from the SNPs on chromosome 1. The simulations were divided into five categories
185 (low, lower-middle, middle, higher-middle, high) based on the fixation index (F_{ST}) between
186 populations A and B for the first QTN (Fig. S1 in Supplemental File 2). We assumed that the
187 first QTN had four times greater variance than the remaining two QTNs (referred to as “QTN1”,
188 “QTN2”, and “QTN3” respectively). The remaining two QTNs were chosen randomly from
189 SNPs where the F_{ST} between A and B were low (SNPs whose F_{ST} value was in the lower 20%
190 category among the 672,923 SNPs). The F_{ST} for each marker was calculated according to Wright
191 (1965) as follows:

$$192 \quad F_{ST} = 1 - \frac{H_S}{H_T}, \quad (\text{Eq. 2})$$

193 where H_S is the average of the expected heterozygosity based on the allele frequencies of
194 populations A and B, and H_T is the expected heterozygosity based on the average allele
195 frequency of populations A and B. H_S and H_T were calculated as follows:

$$196 \quad H_S = \frac{N_A \cdot \{2p_A(1 - p_A)\} + N_B \cdot \{2p_B(1 - p_B)\}}{N_A + N_B}, \quad (\text{Eq. 3})$$

$$197 \quad H_T = 2 \left(\frac{N_A p_A + N_B p_B}{N_A + N_B} \right) \left(1 - \frac{N_A p_A + N_B p_B}{N_A + N_B} \right), \quad (\text{Eq. 4})$$

198 where p_A , p_B , N_A , and N_B are the allele frequencies and the sample sizes of populations A and B
199 respectively, and $N_A = 112$ and $N_B = 100$. The F_{ST} distribution between A and B is shown in
200 Fig. S1, which also shows the thresholds for the five F_{ST} categories.

201 The polygenetic effect in Eq. 5 was sampled from the multivariate normal distribution
202 whose variance-covariance matrix was proportional to the additive numerator relationship matrix
203 \mathbf{A} and was normalized so that their variance was equal to that of the three QTN effects.

$$204 \quad \mathbf{u} \sim \text{MVN}(\mathbf{0}, \mathbf{G}), \quad (\text{Eq. 5})$$

205 where $\mathbf{G} = \mathbf{A}\sigma_A^2$ is the genetic covariance matrix, and the additive genetic variance σ_A^2 was
206 automatically determined from the relationship with heritability. In this study, the additive
207 numerator relationship matrix \mathbf{A} was estimated based on the marker genotype data for 402,509
208 SNPs, which consisted of the core SNPs (defined by the Rice SNP-Seek Database as the “404k
209 CoreSNP Dataset”) in all 12 chromosomes (this marker genotype data was prepared separately
210 from the whole-genome sequence data), using the “A.mat” function in R package “rrBLUP”
211 version 4.5 (Endelman and Jannink, 2012; Endelman, 2011).

212 The residual \mathbf{e} in Eq. 6 was sampled identically and independently from the normal distribution,
213 and was then normalized so that the narrow-sense heritability was equal to 0.6. Residual \mathbf{e} was
214 calculated using the following formula:

$$215 \quad \mathbf{e} \sim \text{MVN}(0, \mathbf{I}\sigma_e^2), \quad (\text{Eq. 6})$$

216 where \mathbf{I} is an identity matrix, and the residual variance σ_e^2 was determined so that the heritability
217 was equal to 0.6.

218

219

220 **Genome-wide association study (GWAS) using simulated data**

221 We performed a GWAS on the seven non-mixture (A, B, C, D) and mixture populations
222 (E, F, and G) using the marker genotype data and the simulated phenotypic data. We fitted the
223 linear mixed model (Yu et al., 2006).

$$224 \quad \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{S}_i\alpha_i + \mathbf{Q}\mathbf{v} + \mathbf{Z}\mathbf{u} + \mathbf{e}, \quad (\text{Eq. 7})$$

225 where \mathbf{y} is the vector of phenotypic values, $\mathbf{X}\boldsymbol{\beta}$, $\mathbf{S}_i\alpha_i$, and $\mathbf{Q}\mathbf{v}$ are the fixed effects terms, $\mathbf{Z}\mathbf{u}$ is
226 the random effects term, and \mathbf{e} is the residuals vector. $\boldsymbol{\beta}$ represents all of the fixed effects other
227 than $\mathbf{S}_i\alpha_i$, and $\mathbf{Q}\mathbf{v}$, and \mathbf{X} is the incidence design matrix corresponding to $\boldsymbol{\beta}$. In this study, $\mathbf{X}\boldsymbol{\beta}$
228 was an intercept. $\mathbf{S}_i\alpha_i$ is composed of \mathbf{S}_i , which is the i_{th} marker of the genotype data, and α_i ,
229 which is the effect of that marker. $\mathbf{Q}\mathbf{v}$ is the term used to correct for the effect of population
230 structure, and in this study \mathbf{Q} was the matrix of two eigenvectors corresponding to the upper two
231 eigenvalues of the additive numerator for relationship matrix \mathbf{A} , Finally, \mathbf{u} represents the
232 polygenic effects, and \mathbf{Z} is the incidence design matrix corresponding to \mathbf{u} .

233 We used the EMMAX and P3D algorithms to reduce the computation time (Kennedy et al.,
234 1992; Kang et al., 2008; 2010; Zhang et al., 2010). The “GWAS” function in R package
235 “rrBLUP” version 4.5 (Endelman, 2011) was used to perform the GWAS described above.

236

237

238 **Evaluation of the simulation results**

239 The p -value (or $-\log_{10}(p)$) for each marker effect was estimated 100 times by the
240 GWAS in five patterns according to the size of the F_{ST} for the seven non-mixture/mixture
241 populations. In this study, the following summary statistics were mainly used to evaluate the
242 GWAS results.

243 In the 100 simulations, the QTNs were not always polymorphic in each population
244 (because the MAF of the whole population did not necessarily match the MAF of each individual
245 population). In such cases, the $-\log_{10}(p)$ value of a QTN that was not polymorphic within a
246 population could not be calculated. Therefore, when two SNPs were polymorphic within that
247 population and were adjacent to the QTN, then the statistic of the more significant SNP was used
248 as the QTN statistic. Since it was difficult to detect such QTNs using a GWAS, we calculated the
249 summary statistics by dividing two patterns depending on polymorphism patterns of QTN1, i.e.,
250 whether using all simulation results or using only results whose QTN1 was polymorphic in the
251 target population (referred to as “All” and “Polymorphic in the population”, respectively).

252

253 **Correct detection rate (CDR) and $-\log_{10}(p)$**

254 The first summary statistic was whether the $-\log_{10}(p)$ rate for each QTN exceeded the
255 threshold in each GWAS (referred to as “CDR; correct detection rate”). We assumed that QTNs
256 would be successfully detected by the GWAS when the CDR was large. The $-\log_{10}(p)$ value
257 whose false discovery rate (FDR) was 0.05 was set as the threshold using the Benjamini-
258 Hochberg method (Benjamini and Hochberg, 1995; Storey and Tibshirani, 2003). As the second
259 summary statistic, we used the $-\log_{10}(p)$ for each QTN in each GWAS, and we also assumed
260 that QTNs were successfully detected by the GWAS when this statistic was large.

261

262 **Area under the curve (AUC)**

263 We also regarded the mean of the AUC as one summary statistic. The AUC refers to the
264 area under the receiver operating characteristic (ROC) curve (Fig. S2 in Supplemental File 3),
265 which was obtained by plotting the false positive rate on the horizontal axis and the true positive
266 rate on the vertical axis when the threshold was varied (Hanley and McNeil, 1982). The AUC
267 was calculated using the following formula:

$$268 \quad \text{AUC} = \frac{1}{2}FPR_{S_1}TPR_{S_1} + \frac{1}{2} \sum_{i=2}^{m+1} (FPR_{S_i} - FPR_{S_{i-1}})(TPR_{S_i} + TPR_{S_{i-1}}), \quad (\text{Eq. 8})$$

269 where m is the number of QTNs, and $m = 3$ in this study. The $FPRs$ and $TPRs$ are the $m + 1$
270 vectors whose i_{th} elements are FPR_{S_i} and TPR_{S_i} , respectively. $FPR_{S_i} = TPR_{S_i} = 1$ when $i =$
271 $m + 1$. When $1 \leq i \leq m$, the FPR_{S_i} and TPR_{S_i} represent the false positive rate and the true
272 positive rate at the time when i QTNs exceed the threshold, respectively. They were calculated
273 using the following formula:

274
$$FPRs_i = \frac{FP_i}{FP_i + TN_i} \quad (1 \leq i \leq m), \quad (\text{Eq. 9})$$

275
$$TPRs_i = \frac{TP_i}{TP_i + FN_i} \quad (1 \leq i \leq m), \quad (\text{Eq. 10})$$

276 where TP_i , FP_i , FN_i , and TN_i are the numbers of SNPs that are the true positives (where the SNP
277 is a QTN and exceeds the threshold), the false positives (where the SNP is not a QTN but
278 exceeds the threshold), the false negatives (where the SNP is a QTN but does not exceed the
279 threshold), and the true negatives (where the SNP is not a QTN and does not exceed the
280 threshold) at the time when i QTNs exceed the threshold respectively. When we evaluated the
281 true/false positive rate, we considered the existence of linkage disequilibrium (LD) by
282 investigating SNPs with LD as one set. In this study, we defined SNPs that satisfied the
283 conditions that they were within 300 kb from the focused SNP and the condition that their
284 squares of the correlation coefficients with the focused SNP were 0.35 or more as one set when
285 considering LD. When we counted TP_i , FP_i , FN_i , and TN_i , we counted the number of the sets
286 described above instead of the number of SNPs. The value for AUC calculated in this manner
287 takes a value between 0 and 1. The GWAS is more successful when the AUC is closer to 1.
288 Using the mean of the AUC as one of the summary statistics meant that it was possible to focus
289 on each QTN and evaluate the overall results of the GWAS.

290

291

292 **Precision, recall, and F-measure**

293 We calculated the mean of precision, the mean of recall, and the mean of the F -measure
294 as other summary statistics to evaluate the GWAS results. These summary statistics can be

295 calculated from the numbers of true positives, false positives, false negatives, and true negatives.
296 More specifically, the precision can be calculated using the following formula:

$$297 \quad \text{Precision} = \frac{TP}{TP + FP}. \quad (\text{Eq. 11})$$

298 We regarded an SNP as “positive” when the $-\log_{10}(p)$ of that SNP exceeded the
299 threshold described above. The precision represents the ratio of the detected SNPs that were
300 QTNs. The recall was defined using the following formula:

$$301 \quad \text{Recall} = \frac{TP}{TP + FN}. \quad (\text{Eq. 12})$$

302 The recall represents the proportion of QTNs detected by the GWAS. Finally, the F -
303 measure was calculated as the harmonic mean of the precision and the recall, and can be used to
304 comprehensively evaluate the GWAS results. The F -measure was calculated using the following
305 formula:

$$306 \quad F = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}. \quad (\text{Eq. 13})$$

307 The greater these summary statistics, the more accurate the GWAS results were.

308

309

310

311 **Degree of genetic diversity index**

312 In order to evaluate the relationship between genetic diversity and the CDR results, we
313 prepared an index that indicated the degree of genetic diversity in each population. The

314 Euclidean distance matrix between accessions for each population was calculated. The median
315 for the off-diagonal elements of the distance matrix was used to indicate the degree of genetic
316 diversity (referred to as the “diversity level”, Table 1). The median was chosen as the diversity
317 level because the distribution of the distances between the accessions for E and G had a double
318 peak. This was because, for mixture populations such as E and G, the distance within the
319 subpopulations was small whereas the distance between subpopulations was large. Therefore, if
320 the mean of the distances (almost the same as Nei’s gene diversity index (NEI, 1973)) is chosen
321 as the diversity level, then there is a risk of overestimating the diversity level.

322

323

324

RESULTS

325

Comparisons between the CDR and AUC for the QTN1s in each population

326

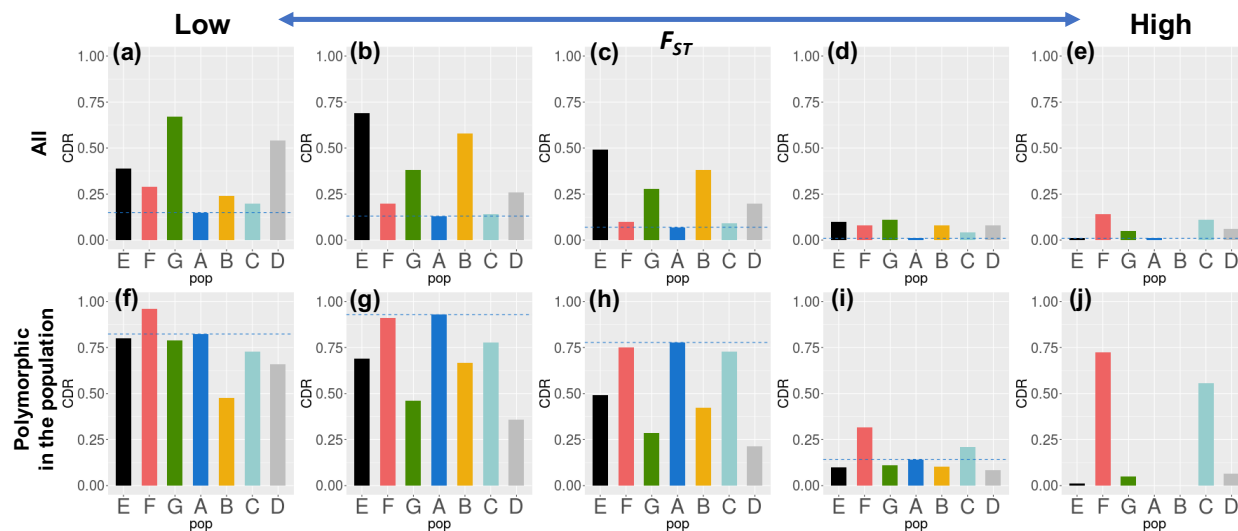
327

328

329

330

The CDRs of the QTN1s in each population were calculated under ten conditions: five levels of F_{ST} between A and B and two patterns of QTN polymorphism, i.e., whether the QTN was polymorphic or not in the target population (Fig. 2 and Table S2 in Supplemental File 5).



331

332 Fig. 2. Correct detection rate for QTN1 in each population under ten conditions.

333 The barplots of CDR of QTN1 in each population under ten conditions: five levels of F_{ST} of QTN1 and
 334 two patterns of polymorphisms of QTN. Blue horizontal dashed lines indicate the CDR in the population
 335 A for each population. A: *japonica*, B: *indica*, C: temperate *japonica*, D: diverse, E: A + B, F: A + C, G:
 336 A + D.

337

338 For almost all levels of F_{ST} , the CDRs for QTN1 in the mixture populations E, F, and G
339 were larger than in the corresponding non-mixture populations B, C, and D, regardless of the two
340 QTN polymorphism patterns (Fig. 2). The CDRs for QTN1 in the mixture populations E, F, and
341 G were always larger than in population A when all the simulation results were taken into
342 account (Fig. 2). When F_{ST} was low, and all simulation results were taken into account (Fig. 2a),
343 populations G and D, which were highly diverse populations, had a higher CDR than the other
344 populations. When F_{ST} was in the lower-middle or middle category, and all simulation results
345 were taken into account (Figs. 2b, c), population E had the highest CDR. The CDR of the highly
346 diverse populations G and D significantly decreased as F_{ST} increased. This result suggested that
347 the QTN1 effect could confound with the population structure at higher F_{ST} values, which meant
348 that it was difficult to detect QTN1 in a highly diverse population. When the F_{ST} value was in the
349 higher-middle or high level categories, and all simulation results were taken into account, (Figs.
350 2d, e), the CDR for QTN1 became quite low in all populations. In populations D, E, and G,
351 QTN1 was hardly detected because of the strong confounding effect of the population structure.
352 In the other populations, the expected heterozygosity (H_e) for QTN1 was extremely small (In A
353 and B, H_e was less than 0.1 in all 100 simulations). The small H_e may make the detection of
354 QTN1 difficult.

355 We excluded the simulations in which there were no polymorphisms in the population
356 so that the detection power of the GWAS when there were polymorphisms in an analyzed
357 population could be evaluated (Figs. 2f-j). When F_{ST} was low, population F had the highest CDR
358 and when F_{ST} was in the lower-middle or middle categories, population A had the highest CDR.
359 However, there were only 14 and 9 cases in which QTN1 was polymorphic in population A. In
360 general, the populations with low or moderate genetic diversities (A, C, and F) had higher CDRs

361 than the populations with high genetic diversities (D, E, and G). When F_{ST} was in the higher-
362 middle or high categories, the results were similar to when F_{ST} was in the lower-middle or
363 middle categories.

364 The CDRs of QTN2 and QTN3 were much lower than that of QTN1 because smaller
365 genetic variances were assigned to these QTLs than QTN1 (Table S2). As in the case of QTN1,
366 for almost all levels of F_{ST} , the CDRs of QTN2 and QTN3 were higher in the mixture
367 populations (E, F, and G) than their corresponding non-mixture populations (B, C, and D).
368 Furthermore, the CDRs for QTN2 and QTN3 in all the mixture populations were higher than for
369 population A. The CDRs for QTN2 and QTN3 were also larger when the F_{ST} for QTN1 was
370 higher.

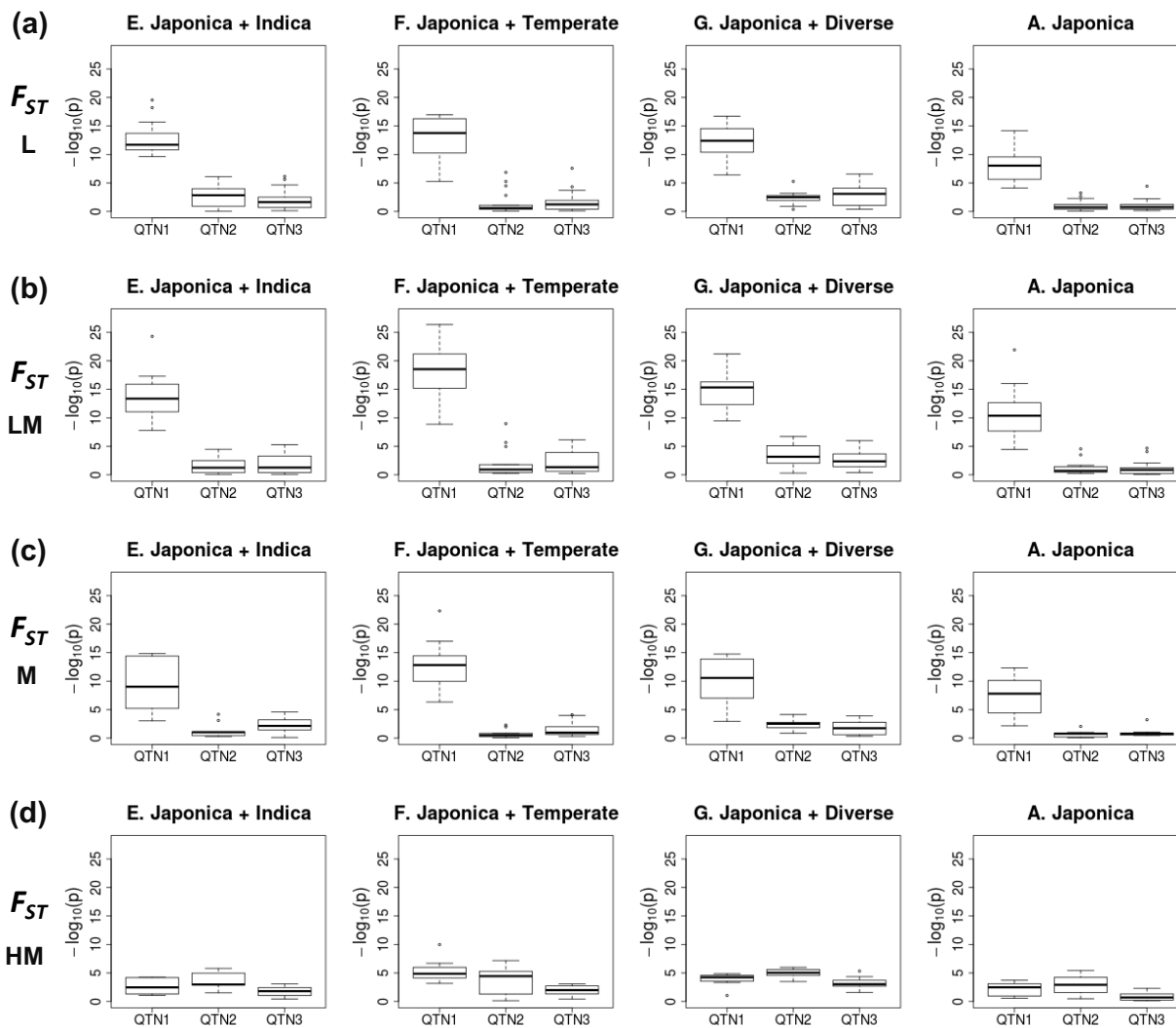
371 Populations D and G had high AUC values in all cases (Table S2). Population F had a
372 smaller AUC than populations D and G, even when the CDR was highest in population F.

373

374 **Comparisons of the $-\log_{10}(p)$ values for the GWAS on each mixture population**
375 **containing *japonica* (A)**

376 We compared the $-\log_{10}(p)$ values for each QTN between populations mixed with the
377 *japonica* population (A) to see if QTN1 was polymorphic in A (Fig. 3). Comparing these values
378 allowed us to examine whether the detection power of the GWAS improved when genetic
379 resources with higher genetic diversities were added to target population A. There is no plot for
380 the high F_{ST} values because no QTN1 was polymorphic in population A over 100 simulations
381 when the F_{ST} of QTN1 was high.

382



383

384 Fig. 3. Boxplots of $-\log_{10}(p)$ of each QTN when QTN1 was polymorphic in *japonica* (A).

385 Boxplots of $-\log_{10}(p)$ of each QTN for each mixture population and *japonica* (A) when QTN1
 386 was polymorphic in A. These plots are shown divided into four categories according to the F_{ST}
 387 value for QTN1 (a: low, b: lower-middle, c: middle, d: higher-middle).

388

389 For all of the four F_{ST} levels, the detection power improved in all mixture populations
390 compared to A (Fig. 3). Population F showed the highest detectability, and this tendency was
391 conspicuous even when F_{ST} was in the middle or higher-middle categories (Figs. 3c, d,
392 respectively). This is because the QTN1 effect is less likely to be confounded with the population
393 structure in F than in the other mixture populations (E and G). Population G had the highest
394 $-\log_{10}(p)$ values for QTN2 and QTN3, although only slightly (Fig. 3).

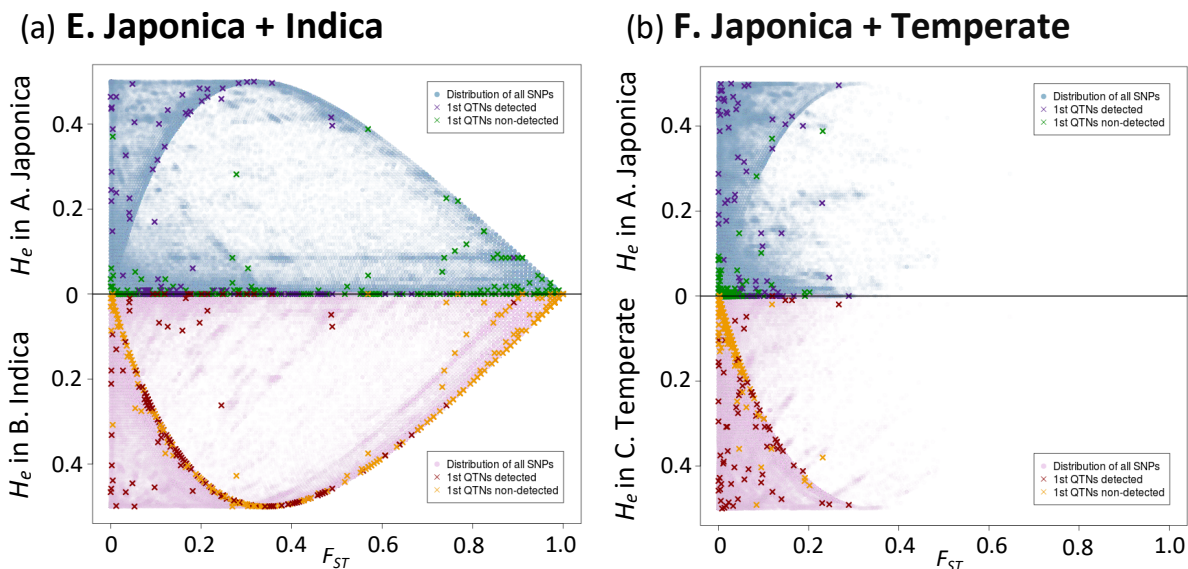
395

396

397 **Factors affecting the detection power of QTNs in the mixture populations**

398 We considered the factors related to the detection power of QTNs in the mixture populations
399 by creating a figure that represented the relationship between F_{ST} , the expected heterozygosity
400 (H_e), and the QTN1 detection power (Fig. 4 and Fig. S4 in Supplemental File 6).

401



402

403 Fig. 4. Relationship between F_{ST} , H_e , and the detection power of QTN1.

404 The distribution of each marker is plotted thinly with between subpopulation F_{ST} on the

405 horizontal axis and H_e of each subpopulation on the vertical axis. The dark X marks on the plot

406 show the SNPs selected as QTN1s in this study. Red and purple marks were detected by GWAS,

407 and green and yellow ones were not detected by GWAS.

408

409 Detection of the QTNs by the GWAS was generally difficult when the between-
410 subpopulation F_{ST} value was high, or H_e was low (Fig. 4). There seemed to be a significant
411 difference between plots F and E or G (Fig. 4a, Fig. S4, and Fig. 4b). However, in population F,
412 because A and C are genetically close, the F_{ST} between the subpopulations was not high.
413 Therefore, the relationships between F_{ST} , H_e , and the GWAS detection power applied to all
414 mixture populations.

415 Some of the QTNs were detected by the GWAS when F_{ST} was in the medium category, and
416 H_e in one of the subpopulations was close to 0 (Fig. 4a and Fig. S4). This suggested that even if
417 the QTN was fixed in one subpopulation, the QTN may still be detected by the GWAS if another
418 subpopulation was added to the analysis.

419

420

421

422

423 **Comparisons among the precision, recall, and F -measure values for each population**

424 The three summary statistics (the mean of precision, the mean of recall, and the mean of the
425 F -measure) were also calculated under ten conditions (Fig. S5 in Supplemental File 7). The
426 precision of the mixture populations was better than the precision value for population A for
427 almost all F_{ST} categories when all simulation results were taken into account (Fig. S5a). However,
428 it is not necessarily true that the precision of the mixture populations outperformed that of their
429 original genetic resources (compare E with B, F with C, and G with D). The recall values of the
430 mixture populations were larger than for their original genetic resources under all conditions.

431 Finally, a comparison of the F -measure for each population showed that there seemed to be no
432 tendency associated with F_{ST} . Therefore, it was difficult to conclude which population was
433 suitable for a GWAS when the F -measure is used. These results indicated that using mixture
434 populations for a GWAS led to the detection of more SNPs, including QTNs.

435

436 **Relationship between the CDR results and genetic diversity**

437 The relationship between the CDR results for QTN1 and the degree of genetic diversity was
438 evaluated under the two QTN polymorphism patterns, i.e., whether or not QTN was polymorphic
439 in the population (Fig. S6 in Supplemental File 8). The CDRs for the mixture populations were
440 usually larger than for the non-mixture populations if their diversity levels were close (Fig. S6a,
441 b). A comparison of the results for the different F_{ST} categories showed that when F_{ST} was low,
442 the populations with the highest diversities, such as D or G, had the highest CDRs, and when F_{ST}
443 was in the lower-middle or middle categories, the populations with the second-highest diversities,
444 such as B or E, had the highest CDRs. Finally, when F_{ST} was in the higher-middle or high
445 categories, the populations with relatively low diversities, such as C or F, had the highest CDRs
446 (Fig. S6a). However, when the simulations in which there were no polymorphisms in the
447 population were excluded, the populations with relatively low diversities, such as A, C, or F, had
448 the highest CDRs in almost all the F_{ST} categories (Fig. S6b).

449

450

451

452

DISCUSSION

453

Relationship between F_{ST} and QTN detection

454 One of the main results of this study was that the detection of QTNs was difficult in populations
455 with high genetic diversities, such as D, E, and G, when the F_{ST} for QTN1 between *japonica* (A)
456 and *indica* (B) was high. This was because the QTN effect confounds with the effect of
457 population structure in these populations. We also examined the reasons why the CDRs for
458 QTN2 and QTN3 were high when the QTN1 F_{ST} value was high.

459 In this study, phenotypic values were simulated using the following expression:

$$460 \quad \mathbf{y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \mathbf{X}_3\beta_3 + \mathbf{u} + \mathbf{e}, \quad (\text{Eq. 3})$$

461 where \mathbf{u} is the polygenetic effect, and is the term that reflects differences between accessions and
462 thus differences between subpopulations. Therefore, if the degree of QTN1 genetic
463 differentiation between *japonica* (A) and *indica* (B) is high, it can be assumed that there is a high
464 correlation between $\mathbf{X}_1\beta_1$ and \mathbf{u} . In this study, we generated phenotypic values using a certain
465 variance ratio under the assumption that each term is independent. Therefore, if there is a
466 correlation between $\mathbf{X}_1\beta_1$ and \mathbf{u} , and the variance between these two terms is considered as one
467 unit, it can be assumed that the variance is smaller than the total value of the two variances under
468 the assumption of independence. Therefore, the variance of these two terms ($\mathbf{X}_1\beta_1 + \mathbf{u}$) in the
469 total phenotypic variance becomes smaller, whereas the variances caused by the terms $\mathbf{X}_2\beta_2$ and
470 $\mathbf{X}_3\beta_3$ become greater than those when it is assumed that each term is independent.

471 The GWAS model used in this study was

$$472 \quad \mathbf{y} = \mathbf{X}\beta + \mathbf{S}_i\alpha_i + \mathbf{Q}\mathbf{v} + \mathbf{Z}\mathbf{u} + \mathbf{e}, \quad (\text{Eq. 7})$$

473 where \mathbf{Qv} is the term used to correct the effect of population structure, and \mathbf{Zu} shows the
474 polygenetic effect. In this GWAS model, $\mathbf{S}_i\alpha_i$ and \mathbf{Qv} or \mathbf{Zu} have some correlation when $\mathbf{S}_i =$
475 \mathbf{X}_1 . This correlation results in the underestimation of α_i by the terms originally used to correct
476 the effects of population structure or family relatedness, such as \mathbf{Qv} and \mathbf{Zu} . Therefore, QTN
477 detection is quite difficult when a GWAS is performed on mixture populations. For QTN2 and
478 QTN3, where $\mathbf{S}_i = \mathbf{X}_2$ or $\mathbf{S}_i = \mathbf{X}_3$, there is generally no correlation between $\mathbf{S}_i\alpha_i$ and \mathbf{Qv} or \mathbf{Zu} .
479 Therefore, the detection of these QTNs is not related to these terms. Furthermore, the variances
480 represented by the terms $\mathbf{X}_2\beta_2$ and $\mathbf{X}_3\beta_3$ are considered to be higher when the QTN1 genetic
481 differentiation is not high. Therefore, the CDRs of QTN2 and QTN3 were high when the F_{ST} for
482 QTN1 was high (Fig. 2 and Table S1). It has been suggested by Atwell et al. (2010) that a bias
483 may occur in the GWAS results when the QTN correlates with population structure or family
484 relatedness.

485

486

487 **Relationship between H_e and the QTN detection**

488 The detection of QTNs by a GWAS was difficult when the expected heterozygosity (H_e) in the
489 population was low. When H_e in the population was low, the MAF was low, and alleles and
490 mutations with low allele frequencies are known as "rare alleles" or "rare variants". In such cases,
491 the QTN effect when the H_e values are low may be concealed by the QTN effect when H_e is not
492 low or by the environmental effect because there are few accessions with one allele. For this
493 reason, it is generally challenging to detect QTNs in such cases, but a method to deal with this
494 problem has been developed (Wu et al., 2011).

495 One example of a “rare variant” that is common in plants is the haplotype condition. It has
496 been reported that haplotypes are difficult to detect using a GWAS (Stram, 2014). This is
497 because haplotypes are often “rare variants” and their H_e values in the population are low.
498 Another problem caused by "rare variants" is that the non-causal SNP whose LD is strong with
499 the “rare variant” may have a higher $-\log_{10}(p)$ value than the “rare variant”. This occurrence,
500 known as “synthetic association”, often happens when the minor allele frequency of the SNP is
501 higher than that of the “rare variant” (Dickson et al., 2010). These "synthetic associations" were
502 often detected in this simulation study.

503

504

505 **Summary and further discussion on each result**

506 Generally, the CDRs of the QTNs showed that the populations suitable for a GWAS were
507 different depending on whether all the QTNs were to be detected or only the polymorphic QTNs
508 in the target population. Specifically, if all QTNs are to be detected when the degree of genetic
509 differentiation between QTNs is low, then it is optimal to use a population with high genetic
510 diversity that has as many polymorphisms as possible. However, as the degree of genetic
511 differentiation becomes more extensive, a population with high genetic diversity is not suitable
512 for a GWAS because the QTN effect is more likely to confound with the population structure. In
513 contrast, a population with moderate genetic diversity, such as population F, was suitable for a
514 GWAS, regardless of the degree of genetic differentiation. This was partly because the QTN1
515 effect was less likely to confound with the population structure in F than in E or G, even when
516 F_{ST} was high. However, in either case, when the degree of genetic differentiation is extensive, it
517 is difficult to detect the QTNs in any population. Therefore a GWAS analysis is not suitable,

518 which means that another approach, such as biparental QTL mapping, must be used to identify
519 genes (Lander and Botstein, 1989).

520 Population F had a smaller AUC than populations D and G, even when the CDR for
521 population F was the highest. From its definition, AUC is more dependent on how low
522 $-\log_{10}(p)$ of the QTN with the lowest $-\log_{10}(p)$ value is than on how high the $-\log_{10}(p)$ of
523 the QTN with the highest $-\log_{10}(p)$ value is. Furthermore, in this study, the number of markers
524 for the GWAS differed (Table 1). When $-\log_{10}(p)$ values for the QTNs were similar among the
525 different populations, the larger number of markers meant that the true negative rate increased,
526 and the false positive rate decreased in a population, which resulted in an increase in the AUC of
527 a population with a larger number of markers, e.g. D and G.

528 A comparison of the mixture populations and *japonica* (A) using $-\log_{10}(p)$ showed
529 that when the QTNs are polymorphic in a target population with low genetic diversity, genetic
530 resources with higher genetic diversities should be added to the target population. However, in
531 order to avoid cases where the degree of genetic differentiation among the QTNs is extensive
532 between the target population and genetic resources, it is desirable to use populations that are
533 genetically close to the target population.

534 Finally, the results suggested that the F_{ST} differences between the subpopulations and
535 the expected heterozygosity (H_e) of each subpopulation greatly influenced QTN detection by the
536 GWAS in the mixture populations (Fig. 4 and Fig. S4). This result was in agreement with the
537 above finding that QTN detection using a GWAS was generally difficult when F_{ST} was high, or
538 H_e were low. However, these situations frequently happened when the F_{ST} between the
539 subpopulations was moderate. Therefore, even if a QTN is fixed in one subpopulation, it may be

540 possible to detect the QTN by adding another population to the analysis because when the H_e of
541 the QTN is low in one population and F_{ST} is moderate, it can be assumed that H_e is relatively
542 high in the other population. Therefore, the H_e of the mixture population as a whole becomes
543 larger and the detection of a QTN is possible unless the confounding of the effect of that QTN
544 with the population structure is extensive. Although this situation is not difficult to interpret, it is
545 extremely important that SNPs with high F_{ST} and low H_e values must exist in large numbers
546 among populations. After taking this fact into account, a GWAS with a mixture population can
547 be useful. Therefore, creating the proposed diagram shown in Fig. 4 and Fig. S4, will lead to a
548 quantitative understanding of what kind of SNPs can be detected by a GWAS in mixture
549 populations of interest.

550

551

552 **Relationships with using whole-genome sequences**

553 One of the major factors related to the QTN detection power was the fixation index F_{ST}
554 differences among subpopulations. When the F_{ST} difference between the *japonica* (A) and *indica*
555 (B) subpopulations was low, the CDR of the mixture populations was high. One example of such
556 markers is that mutations may have occurred at the same position in both populations after they
557 differentiated. Since such variants are relatively new variants, the LD relationship between these
558 variants and surrounding markers will be weak. Therefore, these variants cannot be detected
559 using marker genotype data with a small number of markers, such as an SNP array. However, the
560 use of whole-genome sequences will increase the marker density, which improves the possibility
561 of detecting such variants with a GWAS. In summary, using whole-genome sequences improves

562 the possibility of detecting QTNs with low F_{ST} values and the use of mixture populations should
563 further improve the QTN detection power. In this study, there were cases where SNPs in a low
564 LD region were selected as QTNs when F_{ST} was low.

565

566

567

568

CONCLUSION

569 In this study, we examined a way of selecting a population that was suitable for a GWAS by
570 conducting simulations using populations with various genetic backgrounds. We evaluated the
571 results of the simulations by dividing them into ten patterns according to two criteria: the degree
572 of genetic differentiation (F_{ST}) between two main subpopulations and QTN polymorphism in a
573 target population. When the QTNs are polymorphic in a target population, increasing the
574 population size by adding available genotypes to the target population improves the detection
575 power. We suggest that a population genetically similar to a target population is desirable. After
576 investigating F_{ST} and expected heterozygosity H_e as factors that may substantially influence the
577 detection power of a GWAS, the results showed that SNPs with high F_{ST} and low H_e values were
578 less likely to be detected by a GWAS that used mixture populations. These results indicated that
579 the detection power of a GWAS was improved by using mixture populations with different
580 genetic backgrounds. Furthermore, the use of publicly available whole-genome sequences meant
581 it was possible to increase the population size and to use polymorphic markers that were present
582 in high numbers, which should also improve the detection power of the GWAS.

583

584

585

586

587

588

ACKNOWLEDGMENTS

589 This study was supported by Grant-in-Aid for Scientific Research(A) (25252002), Grant-
590 in-Aid for Scientific Research(B) (Grant number 15H04436), JST, PRESTO (Grant number
591 JPMJPR15O6), the Cross-ministerial Strategic Innovation Promotion Program (SIP), the
592 “Technologies for creating next-generation agriculture, forestry and fisheries” (funding agency:
593 Bio-oriented Technology Research Advancement Institution, NARO), and JST CREST (Grant
594 Number JPMJCR16O).

595 We are also grateful for the advice given by Dr. Ryohei Tanaka.

596

SUPPLEMENTAL MATERIAL

597 Supplemental File 1: **Table S1.** Information about the 299 rice accessions used in this study.

598 Supplemental File 2: **Fig. S1.** Histogram showing the F_{ST} differences between *japonica* (A) and *indica*
599 (B).

600 Supplemental File 3: **Fig. S2.** Example of a ROC curve and the AUC.

601 Supplemental File 4: **Fig. S3.** Principal components analysis results for chromosome 1 and all the
602 chromosomes.

603 Supplemental File 5: **Table S2.** Correct detection rate rates for all QTNs and the AUC in each population.

604 Supplemental File 6: **Fig. S4.** Relationship between F_{ST} , H_e , and the QTN1 detection power for the
605 population G.

606 Supplemental File 7: **Fig. S5.** Bar plots of the precision, the recall and the F -measure results.

607 Supplemental File 8: **Fig. S6.** Relationship between the diversity level and the CDR of QTN1.

608 Supplemental File 9: **Supplementary Note.** Additional information about the materials used in this study.

609 **OPTIONAL SECTIONS**

610 **Availability of data and material**

611 Whole genome sequencing data are available of 112 accessions of *Oryza sativa* subsp.
612 *japonica* in the DDBJ Sequence Read Archive (SRA106223, ERA358140, DRA000158,
613 DRA000307, DRA000897, DRA000927, DRA007273, DRA007256, and DRA008071). Whole
614 genome sequencing data for all the other accessions are available in the "Rice SNP-Seek
615 Database".

616

617 **Competing interests**

618 The authors declare that they have no competing interests.

619

620 **Author's contributions**

621 KH, HKK, and HI conceived and designed the study. KH and HI performed the
622 mathematical and statistical analysis. KH, HKK, MY, EK, SY and HN contributed to marker

623 genotyping. KH, HKK, and HI wrote the manuscript in consultation with MY, EK, SY, and HN.
624 All authors read and approved the final manuscript.

625

626

REFERENCES

- 627 Alexandrov, N., S. Tai, W. Wang, L. Mansueto, K. Palis, R.R. Fuentes, V.J. Ulat, et al. 2015.
628 SNP-Seek Database of SNPs Derived from 3000 Rice Genomes. *Nucleic Acids Res.*
629 43(D1):D1023–27. doi: 10.1093/nar/gku1039.
- 630 Atwell, S., Y.S. Huang, B.J. Vilhjálmsson, G. Willems, M. Horton, Y. Li, D. Meng, et al. 2010.
631 Genome-Wide Association Study of 107 Phenotypes in Arabidopsis Thaliana Inbred Lines.
632 *Nature* 465(7298):627–31. doi: 10.1038/nature08800.
- 633 Auwera, G.A. Van Der, M.O. Carneiro, C. Hartl, R. Poplin, A. Levy-moonshine, T. Jordan, K.
634 Shakir, et al. 2014. From FastQ Data to High Confidence Variant Calls: The Genome
635 Analysis Toolkit Best Practices Pipeline. *Curr Protoc Bioinformatics* 11(10):1–33. doi:
636 10.1002/0471250953.bi1110s43.From.
- 637 Begum, H., J.E. Spindel, A. Lalusin, T. Borromeo, G. Gregorio, J. Hernandez, P. Virk, B.
638 Collard, and S.R. McCouch. 2015. Genome-Wide Association Mapping for Yield and Other
639 Agronomic Traits in an Elite Breeding Population of Tropical Rice (*Oryza Sativa*). *PLoS*
640 *One* 10(3):1–19. doi: 10.1371/journal.pone.0119873.
- 641 Benjamini, Y., and Y. Hochberg. 1995. Controlling the False Discovery Rate: A Practical and
642 Powerful Approach to Multiple Testing. *J. R. Stat. Soc.* 57(1):289–300. doi:
643 10.2307/2346101.
- 644 Bolger, A.M., M. Lohse, and B. Usadel. 2014. Trimmomatic: A Flexible Trimmer for Illumina
645 Sequence Data 30(15):2114–20. doi: 10.1093/bioinformatics/btu170.
- 646 Browning, B.L., and S.R. Browning. 2016. Genotype Imputation with Millions of Reference
647 Samples. *Am. J. Hum. Genet.* 98(1):116–26. doi: 10.1016/j.ajhg.2015.11.020.
- 648 Danecek, P., A. Auton, G. Abecasis, C.A. Albers, E. Banks, M.A. DePristo, R.E. Handsaker, et
649 al. 2011. The Variant Call Format and VCFtools 27(15):2156–58. doi:
650 10.1093/bioinformatics/btr330.
- 651 Dickson, S.P., K. Wang, I. Krantz, H. Hakonarson, and D.B. Goldstein. 2010. Rare Variants
652 Create Synthetic Genome-Wide Associations. *PLoS Biol.* 8(1):e1000294. doi:
653 10.1371/journal.pbio.1000294.
- 654 Endelman, J.B. 2011. Ridge Regression and Other Kernels for Genomic Selection with R
655 Package RrBLUP. *Plant Genome J.* 4(3):250. doi: 10.3835/plantgenome2011.08.0024.
- 656 Endelman, J.B., and J.L. Jannink. 2012. Shrinkage Estimation of the Realized Relationship
657 Matrix. *G3 (Bethesda)* 2(11):1405–13. doi: 10.1534/g3.112.004259.

- 658 Hanley, J.A., and B.J. McNeil. 1982. The Meaning and Use of the Area under a Receiver
659 Operating Characteristic (ROC) Curve. *Radiology* 142(1):29–36. doi:
660 10.1080/02634938208400381.
- 661 Huang, X., Y. Zhao, X. Wei, C. Li, A. Wang, Q. Zhao, W. Li, et al. 2012. Genome-Wide
662 Association Study of Flowering Time and Grain Yield Traits in a Worldwide Collection of
663 Rice Germplasm. *Nat. Genet.* 44(1):32–39. doi: 10.1038/ng.1018.
- 664 Jarquin, D., H. Kajiya-Kanegae, C. Taishen, S. Yabe, R. Persa, J. Yu, H. Nakagawa, M.
665 Yamazaki, and H. Iwata. 2019. Coupling Day Length Data and Genomic Prediction Tools
666 for Predicting Time-Related Traits under Complex Scenarios, 703488.
- 667 Kang, H.M., J.H. Sul, S.K. Service, N.A. Zaitlen, S.Y. Kong, N.B. Freimer, C. Sabatti, and E.
668 Eskin. 2010. Variance Component Model to Account for Sample Structure in Genome-
669 Wide Association Studies. *Nat. Genet.* 42(4):348–54. doi: 10.1038/ng.548.
- 670 Kang, H.M., N.A. Zaitlen, C.M. Wade, A. Kirby, D. Heckerman, M.J. Daly, and E. Eskin. 2008.
671 Efficient Control of Population Structure in Model Organism Association Mapping.
672 *Genetics* 178(3):1709–23. doi: 10.1534/genetics.107.080101.
- 673 Kawahara, Y., M. de la Bastide, J.P. Hamilton, H. Kanamori, W.R. McCombie, S. Ouyang, D.C.
674 Schwartz, et al. 2013. Improvement of the *Oryza Sativa* Nipponbare Reference Genome
675 Using next Generation Sequence and Optical Map Data 6(1):1–10. doi: 10.1186/1939-8433-
676 6-1.
- 677 Kennedy, B.W., M. Quinton, and J.A. van Arendonk. 1992. Estimation of Effects of Single
678 Genes on Quantitative Traits. *J. Anim. Sci.* 70(7):2000–2012. doi:
679 10.1016/j.alcohol.2011.07.002.
- 680 Koboldt, D.C., K.M. Steinberg, D.E. Larson, R.K. Wilson, and E.R. Mardis. 2013. The Next-
681 Generation Sequencing Revolution and Its Impact on Genomics. *Cell* 155(1):27–38. doi:
682 10.1016/j.cell.2013.09.006.
- 683 Korte, A., and A. Farlow. 2013. The Advantages and Limitations of Trait Analysis with GWAS:
684 A Review. *Plant Methods* 9(1):29. doi: 10.1186/1746-4811-9-29.
- 685 Lander, E.S., and D. Botstein. 1989. Mapping Mendelian Factors Underlying Quantitative Traits
686 Using RFLP Linkage Maps. *Genetics* 121:185–99.
687 <http://www.ncbi.nlm.nih.gov/pubmed/2563713>.
- 688 Li, H. 2012. Exploring Single-Sample Snp and Indel Calling with Whole-Genome de Novo
689 Assembly 28(14):1838–44. doi: 10.1093/bioinformatics/bts280.
- 690 Li, J.Y., J. Wang, and R.S. Zeigler. 2014. The 3,000 Rice Genomes Project: New Opportunities
691 and Challenges for Future Rice Research. *Gigascience* 3(1):1–3. doi: 10.1186/2047-217X-
692 3-8.
- 693 Mansueto, L., R.R. Fuentes, F.N. Borja, J. Detras, J.M. Abrio-Santos, D. Chebotarov, M.
694 Sanciango, et al. 2017. Rice SNP-Seek Database Update: New SNPs, Indels, and Queries.
695 *Nucleic Acids Res.* 45(D1):D1075–81. doi: 10.1093/nar/gkw1135.
- 696 Mansueto, L., R.R. Fuentes, D. Chebotarov, F.N. Borja, J. Detras, J.M. Abriol-Santos, K. Palis,
697 et al. 2016. SNP-Seek II: A Resource for Allele Mining and Analysis of Big Genomic Data
698 in *Oryza Sativa*. *Curr. Plant Biol.* 7–8:16–25. doi: 10.1016/j.cpb.2016.12.003.

- 699 McKenna, A., M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis, A. Kernytsky, K. Garimella, et
700 al. 2009. The Genome Analysis Toolkit: A MapReduce Framework for Analyzing next-
701 Generation DNA Sequencing Data. *Proc. Int. Conf. Intellect. Capital, Knowl. Manag.*
702 *Organ. Learn.* 20:254–60. doi: 10.1101/gr.107524.110.20.
- 703 Metzker, M.L. 2010. Sequencing Technologies the next Generation. *Nat. Rev. Genet.* 11(1):31–
704 46. doi: 10.1038/nrg2626.
- 705 Misra, G., S. Badoni, R. Anacleto, A. Graner, N. Alexandrov, and N. Sreenivasulu. 2017. Whole
706 Genome Sequencing-Based Association Study to Unravel Genetic Architecture of Cooked
707 Grain Width and Length Traits in Rice. *Nat. Sci. Reports* 7(1):12478. doi: 10.1038/s41598-
708 017-12778-6.
- 709 NEI, M. 1973. Analysis of Gene Diversity in Subdivided Populations. *Proc. Nat. Acad. Sci. USA*
710 70(12):3321–23. doi: 10.1016/j.jasrep.2018.01.028.
- 711 Ott, J., J. Wang, and S.M. Leal. 2015. Genetic Linkage Analysis in the Age of Whole-Genome
712 Sequencing. *Nat. Rev. Genet.* 16(5):275–84. doi: 10.1038/nrg3908.
- 713 Paradis, E., J. Claude, and K. Strimmer. 2004. APE: Analyses of Phylogenetics and Evolution in
714 R Language *20(2):289–90.* doi: 10.1093/bioinformatics/btg412.
- 715 Saitou, N., and M. Nei. 1987. The Neighbor-Joining Method - a New Method for Reconstructing
716 Phylogenetic Trees. *Mol. Biol. Evol.* 4(4):406–25.
- 717 Storey, J.D., and R. Tibshirani. 2003. Statistical Significance for Genomewide Studies. *Proc.*
718 *Natl. Acad. Sci.* 100(16):9440–45. doi: 10.1073/pnas.1530509100.
- 719 Stram, D. 2014. *Design, Analysis, and Interpretation of Genome-Wide Association Scans.*
720 Heidelberg, New York: Springer Science+Business Media.
- 721 Wang, W., R. Mauleon, Z. Hu, D. Chebotarov, S. Tai, Z. Wu, M. Li, et al. 2018. Genomic
722 Variation in 3,010 Diverse Accessions of Asian Cultivated Rice. *Nature* 557(7703):43–49.
723 doi: 10.1038/s41586-018-0063-9.
- 724 Wright, S. 1965. THE INTERPRETATION OF POPULATION STRUCTURE BY F-
725 STATISTICS WITH SPECIAL REGARD TO SYSTEMS OF MATING. *Evolution (N. Y.)*.
726 19:395–420.
- 727 Wu, M.C., S. Lee, T. Cai, Y. Li, M. Boehnke, and X. Lin. 2011. Rare-Variant Association
728 Testing for Sequencing Data with the Sequence Kernel Association Test. *Am. J. Hum.*
729 *Genet.* 89(1):82–93. doi: 10.1016/j.ajhg.2011.05.029.
- 730 Yabe, S., M. Yamasaki, K. Ebana, T. Hayashi, and H. Iwata. 2016. Island-Model Genomic
731 Selection for Long-Term Genetic Improvement of Autogamous Crops. *PLoS One* 11(4):1–
732 21. doi: 10.1371/journal.pone.0153945.
- 733 Yano, K., E. Yamamoto, K. Aya, H. Takeuchi, P.C. Lo, L. Hu, M. Yamasaki, et al. 2016.
734 Genome-Wide Association Study Using Whole-Genome Sequencing Rapidly Identifies
735 New Genes Influencing Agronomic Traits in Rice. *Nat. Genet.* 48(8):927–34. doi:
736 10.1038/ng.3596.

- 737 Yu, J., G. Pressoir, W.H. Briggs, I. Vroh Bi, M. Yamasaki, J.F. Doebley, M.D. McMullen, et al.
738 2006. A Unified Mixed-Model Method for Association Mapping That Accounts for
739 Multiple Levels of Relatedness. *Nat. Genet.* 38(2):203–8. doi: 10.1038/ng1702.
- 740 Zhang, Z., E. Ersoz, C.Q. Lai, R.J. Todhunter, H.K. Tiwari, M.A. Gore, P.J. Bradbury, et al.
741 2010. Mixed Linear Model Approach Adapted for Genome-Wide Association Studies. *Nat.*
742 *Genet.* 42(4):355–60. doi: 10.1038/ng.546.
- 743 Zhao, K., C.-W. Tung, G.C. Eizenga, M.H. Wright, M.L. Ali, A.H. Price, G.J. Norton, et al.
744 2011. Genome-Wide Association Mapping Reveals a Rich Genetic Architecture of
745 Complex Traits in *Oryza Sativa*. *Nat. Commun.* 2:467. doi: 10.1038/ncomms1467.
- 746

747 FIGURES AND TABLES

748 Fig. 1. **Unrooted phylogenetic tree plot for four non-mixture populations.**

749 Unrooted phylogenetic tree plot for the four non-mixture populations, which consisted of 112
750 accessions of *japonica* (A), 100 accessions of *indica* (B), 100 accessions of temperate *japonica*
751 (C), and 100 diverse accessions (D) with neighbor-joining method.

752

753 Fig. 2. **Correct detection rate for QTN1 in each population under ten conditions.**

754 The barplots of CDR of QTN1 in each population under ten conditions: five levels of F_{ST} of QTN1 and
755 two patterns of polymorphisms of QTN. Blue horizontal dashed lines indicate the CDR in the population
756 A for each population. A: *japonica*, B: *indica*, C: temperate *japonica*, D: diverse, E: A + B, F: A + C, G:
757 A + D.

758

759 Fig. 3. **Boxplots of $-\log_{10}(p)$ of each QTN when QTN1 was polymorphic in *japonica* (A).**

760 Boxplots of $-\log_{10}(p)$ of each QTN for each mixture population and *japonica* (A) when QTN1
761 was polymorphic in A. These plots are shown divided into four categories according to the F_{ST}
762 value for QTN1 (a: low, b: lower-middle, c: middle, d: higher-middle).

763

764 **Fig. 4. Relationship between F_{ST} , H_e , and the detection power of QTN1.**

765 The distribution of each marker is plotted thinly with between subpopulation F_{ST} on the
766 horizontal axis and H_e of each subpopulation on the vertical axis. The dark X marks on the plot
767 show the SNPs selected as QTN1s in this study. Red and purple marks were detected by GWAS,
768 and green and yellow ones were not detected by GWAS.

769

770 **Table 1. Number of SNPs and the diversity level of non-mixture and mixture populations.**

	Population name	Number of accessions	Number of SNPs	Diversity level[†]
A.	Japonica	112	72,110	263.095
B.	Indica	100	427,943	660.416
C.	Temperate japonica	100	135,665	362.649
D.	Diverse	100	647,731	798.646
E.	A + B	212	633,507	803.064
F.	A + C	212	151,675	334.606
G.	A + D	212	684,774	859.678

771 [†] Diversity level is the index to indicate the degree of genetic diversity, which is described in the

772 Materials and Method section.