

1 **QTG-Finder: a machine-learning based algorithm to prioritize causal**
2 **genes of quantitative trait loci in Arabidopsis and rice**

3

4 **Fan Lin, Jue Fan, Seung Y. Rhee***

5 Department of Plant Biology, Carnegie Institution for Science, Stanford, California

6 94305, USA

7

8

9

10 **Running title:** QTL causal gene prioritization

11 **Keywords:** Arabidopsis, causal gene, machine learning, quantitative trait loci, rice

12

13

14 **Corresponding author:** Seung Y. Rhee

15 **Address:** Department of Plant Biology, Carnegie Institution for Science, Stanford,

16 California 94305, USA

17 **Phone:** (650) 739-4251

18 **E-mail:** srhee@carnegiescience.edu

19

20 **Abstract**

21 Linkage mapping is one of the most commonly used methods to identify genetic loci that
22 determine a trait. However, the loci identified by linkage mapping may contain hundreds
23 of candidate genes and require a time-consuming and labor-intensive fine mapping
24 process to find the causal gene controlling the trait. With the availability of a rich
25 assortment of genomic and functional genomic data, it is possible to develop a
26 computational method to facilitate faster identification of causal genes. We developed
27 QTG-Finder, a machine learning based algorithm to prioritize causal genes by ranking
28 genes within a quantitative trait locus (QTL). Two predictive models were trained
29 separately based on known causal genes in Arabidopsis and rice. An independent
30 validation analysis showed that the models could recall about 64% of Arabidopsis and
31 79% of rice causal genes when the top 20% ranked genes were considered. The top 20%
32 ranked genes can range from 10 to 100 genes, depending on the size of a QTL. The
33 models can prioritize different types of traits though at different efficiency. We also
34 identified several important features of causal genes including paralog copy number,
35 being a transporter, being a transcription factor, and containing SNPs that cause
36 premature stop codon. This work lays the foundation for systematically understanding
37 characteristics of causal genes and establishes a pipeline to predict causal genes based on
38 public data.

39

INTRODUCTION

40

41 As the world's population expands, food security faces a major challenge in the near
42 future. By 2050, world population is projected to grow by 34%, which will require a 70%
43 increase of global food production to meet the demand (FAO 2009). To catch up with the
44 growing global food demand, it is important to improve the efficiency of arable land
45 usage by developing better crops.

46 Many agriculturally and medically important traits are quantitative and controlled by
47 multiple genetic loci. Examples include plant height, grain yield, and flowering time in
48 plants and common disorders such as cancer, diabetes, and hypertension in humans. The
49 variation in quantitative traits allows organisms to adapt to various environments (Baxter
50 *et al.* 2010; Leinonen *et al.* 2013). Quantitative traits are determined by a combination of
51 genetic complexity and environmental factors (Mackay 2001). The genetic complexity of
52 quantitative traits comes from the involvement of multiple quantitative trait loci (QTL)
53 and the non-additive interactions among them (Carlborg and Haley 2004; Mackay 2014).
54 Causal genes of QTLs are genes whose differences in the DNA sequence or state cause a
55 phenotypic variation in parental genotypes and are supported by multiple lines of
56 experimental evidence including mutational analysis, transgenic complementation, and
57 deficiency complementation (Weigel and Nordborg 2005). To understand the
58 evolutionary forces and molecular mechanisms that shape the genetic architectures of
59 adaptive traits, we need to identify all major causal genes that contribute to phenotypic
60 variation of the traits and elucidate the molecular mechanisms of their actions. Achieving
61 this goal will facilitate rational engineering of plant traits and more accurate prediction of
62 the effects of their modifications on the engineered plant.

63 QTL linkage mapping and genome wide association study (GWAS) are two common
64 approaches used to identify QTLs, each with its own strengths and limitations. Both
65 mapping approaches are based on the co-segregation of a trait and genetic variants in a
66 population. The population for linkage mapping is usually the progeny of parental plants
67 that differ in a trait, such as an F2 population or recombinant inbred lines (Bergelson and
68 Roux 2010). GWAS mapping uses a natural population that has a heritable variation of a
69 trait. Compared to GWAS, linkage mapping does not suffer from issues like population
70 structure and suffers less from rare alleles (Bergelson and Roux 2010). For example, the
71 most significant seed dormancy QTL DOG1 identified by linkage mapping was not
72 identified by GWAS, likely due to the rarity of the strong allele in the GWAS population
73 (Bentsink *et al.* 2010; He 2014). Confounding population structure can cause a high false
74 positive rate in GWAS, though some methods have been developed to ameliorate it (Price
75 *et al.* 2010). However, efforts to correct it could result in a higher false negative rate
76 (Brachi *et al.* 2010). Linkage mapping is less prone to these issues, but it cannot identify
77 QTLs of minor effects when the sample size is small (Otto and Jones 2000; Xu 2003;
78 Martin and Orgogozo 2013; Wellenreuther and Hansson 2016).

79 For QTLs identified by linkage mapping, finding causal genes underlying them is
80 still a big bottleneck (Bergelson and Roux 2010). In a typical rice linkage mapping, the
81 size of a QTL can range from 200kb- 3Mb, which can harbor tens to hundreds of genes
82 depending on the mapping population and gene density (Bargsten *et al.* 2014; Daware *et*
83 *al.* 2017). Even in the post-genomic era where all the genes in the genome are uncovered,
84 identifying QTL causal genes is not straightforward since many QTLs either contain no
85 obvious candidate genes or too many genes potentially relevant for the trait (Nuzhdin *et*

86 *al.* 1999). Therefore, despite the many QTLs that have been reported in plants, only a few
87 have been studied at the molecular level.

88 To narrow down the range of candidate genes in a QTL region, conventional fine
89 mapping is reliable but time-consuming and labor-intensive. The basis of fine mapping is
90 to create a population that has more recombination events within a QTL in order to
91 identify a smaller genomic segment that co-segregates with the trait. However, the
92 enormous time and labor required for creating and screening a population of progenies
93 limits the usage of this method (Tuinstra *et al.* 1997). Depending on the frequency of
94 recombination, thousands of progenies may need to be genotyped to get to a gene-scale
95 resolution (Dinka *et al.* 2007). For example, 1,160 progenies were screened to identify
96 the *Pi36* gene in rice and as many as 18,994 progenies were screened to identify the
97 causal gene of Bph15 in rice (Yang *et al.* 2004; Liu *et al.* 2005). The high cost associated
98 with genotyping and phenotyping makes it challenging to apply fine mapping to all
99 QTLs.

100 Alternative approaches to refine the candidate list of causal genes include meta-
101 analysis, joint linkage-association analysis, and other computational methods including
102 machine-learning algorithms. The first two approaches require either the availability of
103 many QTL studies on similar traits or an additional association mapping experiment
104 (Buckler *et al.* 2009; Motte *et al.* 2014; Yin *et al.* 2017). Computational methods
105 including machine-learning algorithms have been developed to prioritize disease
106 associated genes and genetic variants in human (Perez-Iratxeta *et al.* 2002; Kircher *et al.*
107 2014; Ritchie *et al.* 2014; Hormozdiari *et al.* 2015). To distinguish disease-associated
108 from non-associated variants, a variety of information has been used, including the effect

109 of polymorphism (Ng and Henikoff 2003; Kircher *et al.* 2014; Gelfman *et al.* 2017),
110 sequence conservation (Pollard *et al.* 2010; Huang *et al.* 2017), regulatory information
111 (Deo *et al.* 2014), expression profile (Mordelet and Vert 2011; Deo *et al.* 2014), Gene
112 Ontology (GO) (Mordelet and Vert 2011), KEGG pathway (Mordelet and Vert 2011),
113 and publications (Perez-Iratxeta *et al.* 2002). In contrast, only two causal gene
114 prioritization approaches are available for plants. One method was developed for GWAS
115 in maize based on co-expression networks (Schaefer *et al.* 2018). Another method was
116 developed for linkage mapping based on biological process GOs (Bargsten *et al.* 2014).
117 To date, no machine-learning approaches using multiple data types have been developed
118 to address this problem.

119 Here, we built a supervised learning algorithm to prioritize QTL causal genes using
120 known causal genes in *Arabidopsis thaliana* (*Arabidopsis*) and *Oryza sativa* (rice) and a
121 suite of publicly available genetic and genomic data. For each species, we trained a
122 predictive model using features based on polymorphism data, function annotation, co-
123 function network, and paralog copy number. By testing the models on an independent set
124 of known causal genes, we demonstrated its efficacy in prioritizing causal genes.

125 MATERIALS AND METHODS

126 Data sources and features used in QTG-Finder

127 Twenty-eight features were extracted from published genome-scale data, which
128 included eight polymorphism features, seventeen functional annotation features, one co-
129 function network feature and two evolutionary features (Supplementary Table S3).

130 *Arabidopsis* polymorphism data of 1,135 accessions was downloaded from 1001

131 Genomes Project (<https://1001genomes.org>) (Consortium 2016) and rice polymorphism
132 data of 3,010 cultivars was downloaded from Rice SNP-Seek Database ([http://snp-
134 seek.irri.org](http://snp-
133 seek.irri.org)) (Mansueto *et al.* 2017). We used SIFT4G (v 2.4) (Ng and Henikoff 2003)
135 and SnpEff (v 4.3r) (Cingolani *et al.* 2012) to annotate the raw polymorphism data. The
136 number of non-synonymous SNP as annotated by SIFT4G was normalized to protein
137 length and used as a numeric feature (normalized_nonsyn_SNP). Non-synonymous SNPs
138 at conserved protein sequences were predicted to cause deleterious amino acid changes
139 by SIFT4G. The presence of deleterious non-synonymous SNPs in a gene was used as a
140 binary feature (is_nonsyn_deleterious). If a gene contained any deleterious non-
141 synonymous SNPs, the “is_nonsyn_deleterious” feature was set to 1, otherwise it was set
142 to 0. Other binary polymorphism features such as “is_start_lost” (start codon lost) and
143 “is_start_gained” (start codon gained) were extracted from SnpEff annotations in the
144 same way. For “is_SNP_cis”, the Position Weight Matrices of cis-elements were
145 downloaded from CIS-BP database (Build 1.02) (Weirauch *et al.* 2014) and mapped to
146 1kb upstream of all genes in the genome using FIMO (v 4.12.0) (Grant *et al.* 2011). The
147 cis-elements with a matching score above 55 were imported into SnpEff library to
148 annotate the SNPs. This matching score cutoff was determined by a cross-validation as
described later.

149 Functional annotation features were binary features based on GO (Gotz *et al.* 2008;
150 Jones *et al.* 2014) and Plant Metabolic Network (PMN) (Schlapfer *et al.* 2017).
151 Arabidopsis and rice genes were annotated by Blast2GO (BLAST+ 2.2.29) and
152 InterProScan (v 5.3-46.0). The molecular function GOs were then converted to high-level
153 functional groups such as transcription factor, receptor, kinase, transporter, and enzyme

154 to mitigate the effect of some inaccurate annotations (Jones *et al.* 2007). To assess the
155 performance of this approach, we compared aggregated high-level GO annotations from
156 Blast2GO with aggregated high-level curated GO annotations from AmiGO
157 (<http://current.geneontology.org/products/pages/downloads.html>) and un-aggregated GO
158 annotations. Genes annotated as enzymes were further classified into 13 PMN metabolic
159 domains such as carbohydrate metabolism and nucleotide metabolism (Schlapfer *et al.*
160 2017). Unclassified genes in PMN were classified as “is_other_metabolism”. Genes
161 annotated as enzymes by GO but not present in PMN databases are either enzymes
162 involved in macromolecule metabolic processes or enzymes without a specific function
163 assigned. Since the majority of them are involved in macromolecule metabolic processes,
164 we named this group as “is_macromolecule_metabolism”.

165 Co-function networks of Arabidopsis and rice were retrieved from AraNet and
166 RiceNet (Lee *et al.* 2010; Lee *et al.* 2011). The sum of all the edge weights of a gene was
167 used as the “network_weight” feature. We used the sum of edge weights because hub
168 genes have been proposed to be hotspots of phenotypic variation (Martin and Orgogozo
169 2013).

170 Paralog copy number (paralog_copy_number) and essential gene prediction
171 (is_essential_gene) were taken from a previous publication (Lloyd *et al.* 2015).

172 **Arabidopsis and rice causal genes used for training and independent validation**

173 For model training and cross-validation, curated causal genes from Martin and
174 Orgogozo were used as positives for algorithm training (Martin and Orgogozo 2013). In
175 total, 60 Arabidopsis and 45 rice causal genes were used as the initial training set
176 (Supplementary Tables S1 and S2). We curated and included gene identifiers and trait

177 categories in these tables (Supplementary Methods). For literature validation, we
178 performed a further literature curation and found eleven Arabidopsis and eighteen rice
179 causal genes, which were not included in the Martin and Orgogozo list (Supplementary
180 Methods and Supplementary Table S8).

181 The QTL regions used for independent validation were obtained from previously
182 published studies. Even though some studies fine mapped the QTLs, we still used the
183 original QTL regions instead of the fine-mapped regions since our method was developed
184 to replace fine mapping. We included all genes between the markers that were used to
185 define a QTL for prioritization. When the genome locations of the markers were not
186 provided in the publication, we searched their genome locations in Gramene marker
187 database (https://archive.gramene.org/db/markers/marker_view).

188 **Algorithm training and parameter optimization**

189 The QTG-Finder algorithm was developed in Python (v 3.6) with the ‘sklearn’
190 package (v 0.19.0) (Pedregosa *et al.* 2011). We developed an extended 5-fold cross-
191 validation framework (Figure 1A) to evaluate training performance and optimize model
192 parameters.

193 For the 5-fold cross validation, curated causal genes were used as positives and the
194 other genes from the genome were used as negatives. The positives were randomly re-
195 split into training and testing positives in a 4:1 ratio and in an iterative manner. Training
196 and testing positives were combined with different sets of negative genes that were
197 randomly selected from the rest of the genome. To increase the combination of positives
198 and negatives, we re-split the positives 50 times randomly and selected negatives 50
199 times. This number of iterations ensured greater than 99% probability that every positive

200 sample co-occurred with every negative at least once in the training or testing set during
201 the cross-validation process. The probability of co-occurrence was calculated as Equation
202 1. P_{co} is the probability of co-occurrence of a positive and a negative in a testing or
203 training set. N is the total number of negative samples. n is the number of negative
204 samples selected as testing or training samples. R is the number of iterations used to re-
205 split the positive set. C is the number of cross-validation folds that contains a positive
206 sample. C was set to 4 for the training set and 1 for the testing test. S is the number of
207 iterations to randomly select the negative set.

$$P_{co} = 1 - \left[\prod_{i=0}^n \left(1 - \frac{1}{N-i} \right) \right]^{R*C*S} \quad (1)$$

208

209 We tested different classifiers and parameters and optimized the model based on Area
210 Under the Curve of the Receiver Operating Characteristic (AUC-ROC). The average
211 AUC-ROC from all iterations was used to evaluate training performance. We tested three
212 classifiers: Random Forest, naïve Bayes, and Support Vector Machine (Cortes and
213 Vapnik 1995; Tin Kam 1998; Zhang 2004)(Supplementary Figure S1). For Random
214 Forest, we tuned the number of trees and the maximum number of features for each tree
215 based on AUC-ROC (Supplementary Figure S2). We used 100 trees and a max_feature of
216 9 for Random Forest. For Support Vector Machine, the RBF kernel was used and the C
217 parameter was tuned. Random Forest was chosen for further analysis since its
218 performance was slightly better than the other two classifiers. The ratio of positives and
219 negatives in training data was also tuned to maximize cross-validation AUC-ROC
220 (Supplementary Figure S3). The best performing positives:negatives ratio was 1:20 for

221 Arabidopsis and 1:5 for rice. For testing, a positives:negatives ratio of 1:200 was used
222 since it is close to the average ratio of causal and non-causal genes in real QTLs.

223 **Feature importance analysis**

224 We implemented a leave-one-out analysis to evaluate feature importance. This
225 method was based on the change of AUC-ROC (Δ AUC-ROC) when leaving out one
226 feature from the models. The same cross-validation framework was used for this analysis.
227 For each iteration, we calculated AUC-ROC on the original and the leave-one-out models
228 developed with the same training and testing datasets. The Δ AUC-ROC was calculated
229 by subtracting the leave-one-out AUC-ROC from the original AUC-ROC. With the
230 results from all iterations, we calculated the average Δ AUC-ROC for each feature.

231 **Independent literature validation**

232 For validation, we applied the models to an independent set of causal genes that were
233 curated from recent literature and not used for cross-validation. The models were trained
234 by all known causal genes from the initial list and negatives were randomly selected from
235 the rest of the genome. Model training was repeated 5,000 times using resampled training
236 negatives from the genome in combination with the same set of known causal genes. The
237 5,000 iterations were conducted to ensure that there was >99% probability that each gene
238 in the genome was selected at least once. We applied the models to each of the
239 independent causal gene and all other genes located within the QTL. All genes within the
240 QTL were ranked based on the frequency of being predicted as a causal gene.

241 **Model performance for multiple QTLs**

242 To understand performance of the model when it was applied to multiple QTLs of
243 the same trait, we conducted simulations. We calculated the probability of including at

244 least K causal genes within the prioritized list at a given cut-off of the rank percentile
245 when applying the models to a total of N QTLs with Equation 2. p is the probability of a
246 known causal gene to be included at a particular cutoff of the prioritized list using the
247 independent set of causal genes found in the literature. x is the number of causal genes
248 included in the cut-off.

$$249 \quad P(x \geq K) = \sum_{x=K}^N \binom{N}{x} p^x (1-p)^{N-x} \quad (2)$$

250 **Trait category analysis**

251 The trait category analysis was performed in a similar way as the independent
252 literature validation except using different training and testing sets. Each curated causal
253 gene was tested once. For each round, one curated causal gene was removed from the
254 training set. Then the model was trained and applied to rank the removed causal gene and
255 200 flanking genes.

256 **Data Availability**

257 The source code for QTG-Finder and related analyses such as cross-validation,
258 feature importance analysis, and trait category analysis are available at
259 https://github.com/carnegie/QTG_Finder. Other supplementary materials are available at
260 FigShare (<https://gsajournals.figshare.com/s/ab3d1972b290d706641e>).

261

262

RESULTS

263 **QTG-Finder: a machine-learning algorithm to prioritize causal genes**

264 We developed the QTG-Finder algorithm to accelerate finding causal genes from
265 QTL data and generated two predictive models in Arabidopsis and rice. These two

266 species were selected for model training since they have the largest number of QTL
267 causal genes (QTGs) that have been discovered by fine mapping and map-based cloning
268 in plants (Martin and Orgogozo 2013). For model training, we selected 60 Arabidopsis
269 and 45 rice causal genes as a positive set (Martin and Orgogozo, 2013, Supplementary
270 Tables S1 and S2). The negative set was a set of genes randomly selected from the rest of
271 the genome. To train the models, we used 28 Arabidopsis features and 27 rice features,
272 including polymorphism features, functional categories of genes, function interference
273 from co-function networks, gene essentiality, and paralog copy number (Supplementary
274 Tables S3, S4 and S5). These features were generally independent from each other; most
275 have a Pearson's correlation coefficient <0.2 (Supplementary Figure S4).

276 We devised an extended cross-validation framework to optimize the models (Figure
277 1A). With this framework, we evaluated the training performance with AUC-ROC and
278 optimized parameters. We used AUC-ROC for model optimization since our goal is not
279 only to identify causal genes (true positives) in the prioritized list but also reduce the
280 number of candidates by eliminating non-causal genes (true negatives) from the
281 prioritized list. To find the optimal parameters, we compared the AUC-ROC of different
282 machine-learning classifiers, modeling parameters, the ratio of positive:negative genes in
283 the training set, and different methods to generate GO features (Supplementary Figures
284 S2, S3, S4, and S5). Random Forest was selected as the classifier since it was less prone
285 to over-fitting and performed better than the other classifiers tested (Supplementary
286 Figure S1). After optimization, AUC-ROC for the Arabidopsis and rice models were 0.86
287 and 0.73, respectively (Figure 1B). The optimized models were also evaluated by
288 confusion matrix (Supplementary Table S6). The true positive and true negative rates

289 calculated from the confusion matrix indicated that the model was better at classifying
290 non-causal genes than causal genes.

291 Since the positive training set used was relatively small, we also evaluated the
292 relationship between training performance and size of the training set. The AUC-ROC
293 increased as a larger training set was used. Interestingly, maximum gain in the AUC-
294 ROC was achieved with 20 causal genes for the traits represented by the training set
295 (Supplementary Figure S6).

296 **Important features for predicting causal genes**

297 With the optimized models, we asked which features were important for causal gene
298 prediction. Since Random Forest uses features and their interactions for classification
299 (Touw *et al.* 2013), the importance of a feature cannot be measured by simple enrichment
300 or depletion of a single feature in causal genes. Therefore, we evaluated feature
301 importance based on the change of AUC-ROC (Δ AUC-ROC) when excluding a feature
302 from the model (Lloyd *et al.* 2015). When an important feature is excluded from the
303 model, the AUC-ROC should decrease.

304 For both Arabidopsis and rice models, eight features decreased AUC-ROC when
305 removed (Figure 2A and Supplementary Table S7). The six most important features for
306 Arabidopsis were paralog copy number, transporter, the number of non-synonymous
307 SNPs normalized to protein length (normalized_nonsyn_SNP), receptor, transcription
308 factor, and SNPs causing premature stop codon (is_stop_gained) (Figure 2A). The six
309 most important features for rice were paralog copy number, macromolecule metabolism,
310 network weight sum, transcription factor, transporter, and SNPs causing premature stop
311 codon (is_stop_gained). Four out of the six most important features were consistent

312 between Arabidopsis and rice models, which were paralog copy number, transporter,
313 transcription factor, and SNPs causing premature stop codon.

314 For the six most important features in Arabidopsis and rice, we examined their ratio
315 in known causal genes versus randomly selected genes in the genome (Figure 2B).
316 Compared to other genes in the genome, the causal genes in both species tended to have
317 more paralogs, higher frequency of being a transporter or a transcription factor, and
318 higher frequency of containing SNPs that cause premature stop codons. In addition,
319 Arabidopsis causal genes were more likely to be a receptor and rice causal genes were
320 more likely to be a non-hub gene.

321 The rest of the features contributed less to, but did not impair much, the model
322 performance ($\Delta\text{AUC-ROC} < 0.02$). Since there was no strong evidence that they impair
323 prediction, we did not remove them from the models for further analysis.

324 **Validating QTG-Finder by ranking an independent set of QTL genes**

325 To assess the predictability of QTG-Finder models, we searched the literature for a
326 separate set of known causal genes from the initial training set. We found eleven
327 Arabidopsis and eighteen rice genes that are likely causal genes underlying QTLs when
328 interpreting linkage mapping with additional evidence such as functional
329 complementation, fine mapping, joint linkage-association analysis or genetic analyses
330 (Supplementary Table S8). These causal genes were not used for model training or cross-
331 validation.

332 To examine model performance independently, we applied the QTG-Finder models
333 to this new set of causal genes. For each known causal gene, we ranked all the genes
334 within its QTL region, based on the frequency of being predicted as a causal gene from

335 5,000 iterations. Since the number of genes in a QTL region varies, we used a gene's
336 rank percentile for evaluation. The rank percentile of a gene indicates the percentage of
337 QTL genes that had higher ranks than the gene of interest.

338 Based on the rank percentile of these known causal genes, we evaluated model
339 performance at different cutoffs of rank percentile such as 5%, 10%, or 20%. We
340 calculated the percentage of known causal genes being recalled at different cutoffs
341 (Figure 3A). The top 20% of the ranked genes included seven Arabidopsis (~64%) and
342 fourteen rice (~79%) causal genes (Supplementary Table S8). This set included 10-100
343 non-causal genes. With a more stringent cutoff of 5%, four Arabidopsis (~27%) and three
344 rice (~26%) causal genes were prioritized. We examined the molecule types, trait
345 categories, and features of the eight known causal genes (4 Arabidopsis and 4 rice) that
346 were not prioritized within the top 20%, but did not observe any special trend.

347 We also asked whether the different strengths of experimental evidence of the causal
348 genes affected the model performance. Causal genes were grouped based on the type of
349 supporting experimental evidence (Table S6). The first group included genes with weak
350 supporting evidence such as mutational analysis. The second group contained genes
351 supported by stronger evidence such as fine mapping, functional complementation, and
352 joint linkage-associate analysis. The average rank percentile of the two groups was
353 similar; 18% for the first group and 16% for the second.

354 Since most linkage mapping studies identify multiple QTLs, we asked what the
355 probability is of identifying causal genes from multiple QTLs simultaneously. We
356 calculated a theoretical model performance on multiple QTL identification as the
357 probability of identifying causal genes for at least N QTLs when applying the model to

358 all QTLs of a trait (Figure 3B and C). For example, assuming there were five QTLs
359 (N=5) of a trait identified by a linkage mapping study and each QTL contained one
360 causal gene. For the Arabidopsis model, the probability of identifying at least one causal
361 gene would be 99% when the top 20% genes of all QTLs were tested experimentally. The
362 probability of identifying all five causal genes would be 10% when the top 20% cutoff
363 was used. We further compared the performance of all three cutoffs, top 20%, top 10%,
364 and top 5%. The probability of identifying at least one out of five causal genes would be
365 no less than 80% for all three cutoffs. However, the probability to recall at least four out
366 of five causal genes at top 20% would be 40%, 14% (at top 10%), and 2% (at top 5%).
367 Therefore, a less stringent cutoff, top 20%, performs much better than a more stringent
368 cutoff if one is interested in finding most of the causal genes or causal genes of a
369 particular QTL. However, if the goal is to identify any causal gene, then screening the top
370 5% of all QTLs may be a more strategic approach.

371 To compare our results with an existing QTL prioritization method for rice (Bargsten
372 *et al.* 2014), we examined how genes in our rice validation set were prioritized in that
373 study. Only three out of eighteen genes were prioritized as candidates when the top 9%
374 genes in QTL regions were considered. For QTG-Finder, eight out eighteen genes were
375 prioritized as candidates when the top 9% genes were considered.

376 **Trait type preference of QTG-Finder models**

377 Since the training set included genes for different types of traits at an imbalanced ratio,
378 we asked how QTG-Finder models would work for each type of traits (Figure 4A). The
379 independent validation in the previous section was based on causal genes related to plant
380 development and disease resistance (Supplementary Table S8). However, this validation

381 set was not large enough for a systematic analysis and did not have any abiotic-stress-
382 related causal genes. Therefore, we performed a rank analysis for different trait categories
383 using the known causal genes from the initial training set (60 for Arabidopsis and 45 for
384 rice). For this rank analysis, each causal gene was removed from the training set once and
385 used for a rank test. The removed causal gene and its 200 neighboring genes in the
386 genome were used as a testing set. We applied the models to each testing set to obtain the
387 rank for each causal gene. Then we calculated the average rank for the causal genes in the
388 four trait categories: development, abiotic stress, biotic stress and “other”. The “other”
389 category included traits in seed hull color, oil composition, necrosis, etc. (Supplementary
390 Tables S1 and S2).

391 Model performance varied for different trait categories. Both abiotic and biotic stress
392 traits had better performance than developmental traits (Figure 4B). This could be
393 because the developmental trait category has more diversified traits and genes than the
394 abiotic and biotic stress trait category. In addition, the Arabidopsis model performed
395 slightly better than the rice model for all trait categories. This trait category analysis can
396 guide us to determine rank cutoffs when applying models to different types of traits.

397

398

DISCUSSION

399 Linkage mapping is a useful tool to identify the genomic regions responsible for many
400 agriculturally and medically important traits. However, it is not straightforward to
401 identify the genes that cause phenotypic variation of the trait from these genomic regions.
402 The discovery of causal genes still requires time-consuming and labor-intensive fine
403 mapping. In this study, we developed a machine-learning algorithm to reduce the number

404 of candidates to be tested experimentally in order to accelerate the discovery of causal
405 genes.

406 **A machine-learning algorithm to prioritize QTL causal genes**

407 Several causal variant or gene prioritization methods have been developed for human
408 data but not many in plants (Bargsten *et al.* 2014; Kircher *et al.* 2014; Jagadeesh *et al.*
409 2016; Schaefer *et al.* 2018). Most prioritization methods have been developed for GWAS
410 mapping in human, an organism where linkage mapping cannot be performed. However,
411 linkage mapping can capture rare alleles and has been broadly used to study quantitative
412 traits of livestock, crops, and model organisms. A causal gene prioritization is especially
413 helpful for large QTLs identified by linkage mapping, which can constitute tens to
414 thousands of genes. One method has been developed in rice to prioritize causal genes for
415 linkage mapping (Bargsten *et al.* 2014). This method is based on the hypothesis that
416 causal genes from multiple QTLs of the same trait are more likely to have the same
417 biological process GO terms, and therefore genes with overrepresented biological process
418 GOs were prioritized as causal genes. However, this method gives no predictions for
419 ~15% of traits and lack an unbiased performance evaluation since the same set of causal
420 genes was used to determine the cutoff and evaluate performance. We evaluated
421 performance of this GO-based method with the causal genes from the validation set used
422 in this study. The GO-based method identified less causal genes compared to QTG-
423 Finder when a similar fraction of QTL was prioritized. Another method named Camoco
424 has been developed in maize to prioritize causal genes for GWAS (Schaefer *et al.* 2018).
425 Camoco prioritizes genes based on the relative strength and degree of co-expression
426 among genes near GWAS peaks. The success of this method depends highly on the gene

427 expression dataset being appropriate for the trait of interest. For example, an expression
428 dataset of root tissues may work better for root-related traits than shoot-related traits. In
429 addition, this method may not be able to capture causal genes that are transiently
430 expressed or expressed at low levels (Moyers 2018). Since each of these approaches
431 utilizes different sets of information, they may be used in conjunction with the QTG-
432 Finder.

433 In this study, we built a supervised learning algorithm using multiple features and
434 validated its efficacy with an independent dataset from the literature. The models could
435 accelerate the discovery of causal genes by ranking all the genes in a QTL region and
436 prioritizing the top 5%, 10%, or 20% genes, which are most likely to contain the causal
437 gene, for experimental testing. Based on an assessment using independent data in the
438 literature, we calculated the performance when applying the models to all QTLs of a trait
439 and compared three cutoffs (top 5%, 10%, and 20%). The less stringent cutoff (top 20%)
440 had a higher chance to find more causal genes (Figure 3B and C) but yielded more
441 candidates that needed to be tested by experiments. The more stringent cutoff (top 5%)
442 had a lower chance to find all causal genes but yielded a smaller set of candidates to test.
443 The probability for the models to find at least one causal gene is high for all three cutoffs.
444 If the goal were to find at least one causal genes for functional studies and the particular
445 QTL regions did not matter, the 5% cutoff would be more efficient. If the goal were to
446 discover all causal genes and understand the genetic architecture of a trait, the 20% cutoff
447 would be better. Similarly, if a particular QTL were of interest for discovering the
448 underlying causal gene, the 20% cutoff would be better.

449 There are several conceptual and practical advantages of QTG-Finder algorithm. First,

450 this algorithm combines multiple types of publically available data including
451 polymorphisms, function annotations, co-function network and other genomic data,
452 which have not been applied to prioritize causal genes from linkage mapping studies.
453 Second, models were trained on causal genes from various traits and can be applied to
454 several types of traditional traits, though the prioritization efficiency was not equivalent.
455 Third, validation from the literature provides guidance on what proportion of genes to
456 prioritize in practice rather than arbitrarily selecting a threshold. Fourth, the models treat
457 each QTL independently and have the flexibility to prioritize a specific QTL of interest.

458 Two limitations of this study are the small number of known causal genes in plants
459 and the impurity of negative set used for model training. As a positive dataset, we used
460 60 Arabidopsis and 45 rice causal genes that have been verified by map-based-cloning.
461 Even though the positive dataset are of high quality, the sample size may not be large
462 enough to represent all categories or features of causal genes and therefore lead to
463 ascertainment bias. The models may perform better on over-represented gene categories
464 or features in the training set. A larger positive training set will mitigate this bias. For
465 example, the qTARO database is a useful resource to find potential new causal genes for
466 rice, though these genes would need to be curated further (Yonemaru *et al.* 2010). The
467 negative set was composed of genes randomly selected from the rest of the genome.
468 Though we excluded known causal genes, there could still be some uncharacterized
469 causal genes. As a result of these limitations, 20% cutoff will still yield ~100 candidates
470 for large QTLs, which is challenging for experimental characterization unless at least a
471 medium-throughput phenotyping method is available. Fortunately, plant science is
472 entering an era of high-throughput phenotyping with advances in automation,

473 computation and sensor technology (Fahlgren *et al.* 2015; Araus *et al.* 2018). Our study
474 establishes an extendable framework that can be easily updated with new training
475 datasets and features. As more causal genes are uncovered, the new data can be easily
476 incorporated to improve the models.

477 The current models included genes in the reference genomes of rice and Arabidopsis.
478 Even though the majority of causal genes are present in the reference genomes, there are
479 exceptions. For example, SUB1A, SNORKEL1, and SNORKEL2 are causal genes absent
480 in the rice reference genome Nipponbare (Xu *et al.* 2006; Hattori *et al.* 2009). Those
481 genes cannot be predicted with the current models. In the future, this could be addressed
482 by using pan-genome gene information and presence–absence variation (Zhao *et al.*
483 2018).

484 **Important features for predicting QTL causal genes**

485 Many causal genes were repeatedly found to cause phenotypic variation of similar
486 traits, which is also known as genetic hotspots of phenotypic variation or gene reuse
487 (Martin and Orgogozo 2013). By examining 1,008 causative alleles in animals, plants,
488 and yeasts, Martin and Orgogozo found *de novo* mutations to occur repeatedly at certain
489 genes or orthologous loci and causing similar phenotypic variations either among
490 lineages or within a single lineage. The prevalence of gene reuse suggests that causal
491 genes are likely to have some genetic and genomic characteristics that allow them to be
492 repeatedly used for phenotypic variation. The mechanism for gene reuse is not clear but it
493 may be influenced by factors such as the availability of standing genetic variation,
494 mutation rate, pleiotropic constraint, and epistatic interactions of a gene (Conte *et al.*
495 2012; Conte *et al.* 2015).

496 While many QTL causal genes have been cloned, their features have not been
497 systematically examined before. Instead of evaluating each feature individually, we
498 trained Random Forest models and evaluated feature importance for all features by
499 adopting the leave-one-out strategy. Several of the most important features were
500 consistent between Arabidopsis and rice models: containing SNPs that cause a premature
501 stop codon, paralog copy number, being a transporter, and being a transcription factor.

502 We extracted polymorphism features from re-sequencing data of many accessions,
503 which provide more information about the existence of standing genetic variation in the
504 species than the polymorphism data used for linkage mapping, which typically comes
505 from two parental lines. DNA polymorphisms such as nonsense SNPs, deleterious non-
506 synonymous SNPs, SNPs at cis-regulatory elements, and SNPs at splice junctions have
507 been used as features to classify causal and non-causal variants of human diseases
508 (Kircher *et al.* 2014; Jagadeesh *et al.* 2016). These SNPs can directly affect the function
509 or expression of a gene and therefore are more likely to be causal than the rest of the
510 SNPs. Our results indicate Arabidopsis and rice causal genes were more likely to carry a
511 SNP that causes premature stop codon (nonsense SNP) than an average gene in the
512 genome. We also found Arabidopsis causal genes were more likely to have more non-
513 synonymous SNPs than an average gene in the genome. Besides the high impact SNPs in
514 coding regions, we also examined polymorphisms in non-coding regions since about 90%
515 of human trait/disease-associated SNPs are located outside of coding regions (Hindorff *et*
516 *al.* 2009). The SNPs at cis-regulatory elements did not show a high feature importance in
517 our algorithm, although this might be due to limited exploration of non-coding sequences
518 in plants. The CIS-BP database contains cis-elements of 44% of the transcription factors

519 in *Arabidopsis* (Weirauch *et al.* 2014). Developing a more accurate and complete map of
520 functional non-coding regions based on conserved noncoding sequences (Van de Velde *et*
521 *al.* 2014) will potentially make non-coding polymorphism features more useful for
522 prioritizing causal genes in the future. The SNPs linked to causal SNPs might add
523 background noise and reduce the capability to distinguish causal genes from non-causal
524 genes. This could be a reason why half of the polymorphism features were not
525 significantly enriched in the causal genes (Supplemental Table S3).

526 Paralogs contribute to the evolution of plant traits by providing functional divergence
527 that gives plants the potential to adapt to complex environments (Panchy *et al.* 2016).
528 Through evolution, genes involved in signal transduction and stress response have
529 retained more paralogs while essential genes like DNA gyrase A have retained fewer
530 paralogs (Lloyd *et al.* 2015; Panchy *et al.* 2016). By acquiring new functions or sub-
531 functions, paralogs allow plants to sense and handle different environmental conditions in
532 a more comprehensive and adjustable way. For example, the eight paralogous heavy
533 metal ATPases (HMAs) in *Arabidopsis* are all involved in heavy metal transport but have
534 different substrate preferences, tissue expression patterns, and subcellular compartment
535 locations (Takahashi *et al.* 2012). Three of them, HMA3, HMA4, HMA5, are known
536 causal genes of QTLs identified by linkage mapping. The known causal genes we
537 analyzed have more paralog copies than other genes in the genome. This suggests that
538 many plant causal genes are playing a role in providing more phenotypic tuning
539 parameters to allow plants to adapt to complex environments.

540 When the training set is small, there is a possibility of ascertainment bias. For
541 example, the GO features, `is_transporter` and `is_transcription_factor`, may be considered

542 as important features because of their enrichment in the current training set. We will have
543 more confidence of the importance of features when more known causal genes become
544 available.

545 The important features of causal genes identified by linkage mapping may be different
546 from those identified by GWAS. Given the difference of the two genetic approaches,
547 linkage mapping tends to identify large-effect alleles of protein-coding regions, while
548 GWAS tends to identify common alleles with a wider range of effect sizes at protein-
549 coding regions or non-coding regions (Singleton *et al.* 2010). Therefore, whether the
550 features used in this study can be applied to GWAS remains open. It would be interesting
551 to systematically compare causal genes identified by linkage mapping and GAWS in the
552 future.

553 Overall, QTG-Finder is a novel machine-learning pipeline to prioritize causal genes
554 for QTLs identified by linkage mapping. We trained QTG-Finder models for Arabidopsis
555 and rice based on known causal genes from each species, respectively. By utilizing
556 information like polymorphisms, function annotations, co-function networks, and paralog
557 copy numbers, the models can rank QTL genes to prioritize causal genes. Our
558 independent literature validation demonstrates that the models can recall about 64% of
559 causal genes for Arabidopsis and rice when the top 20% of ranked QTL genes were
560 considered. The algorithm is applicable to any traditional quantitative traits but the
561 performance was different for each trait type. Since QTG-Finder is a machine-learning
562 based pipeline, extending the training set and adding features can easily expand and
563 improve the models. We envision that frameworks like QTG-Finder can accelerate the

564 discovery of novel quantitative trait genes by reducing the number of candidate genes and
565 efforts of experimental testing.

566

567 **Acknowledgements**

568 We thank Dr. John Lloyd and Dr. Shin-Han Shiu for sharing the data of rice essential
569 gene prediction. We thank Kevin Radja for testing the source code and giving useful
570 comments.

571 **Funding**

572 This work was supported by the United States Department of Energy's Biological and
573 Environmental Research Program [DE-SC0008769, DE-SC0018277].

574 **Conflict of interest**

575 The authors declare no conflict of interest.

576

577 **References**

- 578 Araus, J.L., S.C. Kefauver, M. Zaman-Allah, M.S. Olsen, and J.E. Cairns, 2018 Translating High-
579 Throughput Phenotyping into Genetic Gain. *Trends Plant Sci.* 23 (5):451-466.
- 580 Bargsten, J.W., J.P. Nap, G.F. Sanchez-Perez, and A.D. van Dijk, 2014 Prioritization of candidate genes in
581 QTL regions based on associations between traits and biological processes. *BMC Plant Biol.*
582 14:330-311.
- 583 Baxter, I., J.N. Brazelton, D. Yu, Y.S. Huang, B. Lahner *et al.*, 2010 A coastal cline in sodium
584 accumulation in *Arabidopsis thaliana* is driven by natural variation of the sodium transporter
585 *AtHKT1;1*. *PLoS Genet.* 6 (11):e1001193.
- 586 Bentsink, L., J. Hanson, C.J. Hanhart, H. Blankestijn-de Vries, C. Coltrane *et al.*, 2010 Natural variation for
587 seed dormancy in *Arabidopsis* is regulated by additive genetic and molecular pathways. *Proc.*
588 *Natl. Acad. Sci. U. S. A.* 107 (9):4264-4269.
- 589 Bergelson, J., and F. Roux, 2010 Towards identifying genes underlying ecologically relevant traits in
590 *Arabidopsis thaliana*. *Nat. Rev. Genet.* 11 (12):867-879.
- 591 Brachi, B., N. Faure, M. Horton, E. Flahauw, A. Vazquez *et al.*, 2010 Linkage and association mapping of
592 *Arabidopsis thaliana* flowering time in nature. *PLoS Genet.* 6 (5):e1000940.
- 593 Buckler, E.S., J.B. Holland, P.J. Bradbury, C.B. Acharya, P.J. Brown *et al.*, 2009 The Genetic Architecture
594 of Maize Flowering Time. *Science* 325 (5941):714-718.
- 595 Carlborg, O., and C.S. Haley, 2004 Epistasis: too often neglected in complex trait studies? *Nat. Rev. Genet.*
596 5 (8):614-618.
- 597 Cingolani, P., A. Platts, L.L. Wang, M. Coon, T. Nguyen *et al.*, 2012 A program for annotating and
598 predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of
599 *Drosophila melanogaster* strain w(1118); iso-2; iso-3. *Fly* 6 (2):80-92.

600 Consortium, T.G., 2016 1,135 Genomes Reveal the Global Pattern of Polymorphism in *Arabidopsis*
601 *thaliana*. *Cell* 166 (2):481-491.

602 Conte, G.L., M.E. Arnegard, J. Best, Y.F. Chan, F.C. Jones *et al.*, 2015 Extent of QTL Reuse During
603 Repeated Phenotypic Divergence of Sympatric Threespine Stickleback. *Genetics* 201 (3):1189-
604 1200.

605 Conte, G.L., M.E. Arnegard, C.L. Peichel, and D. Schluter, 2012 The probability of genetic parallelism and
606 convergence in natural populations. *Proc Biol Sci* 279 (1749):5039-5047.

607 Cortes, C., and V. Vapnik, 1995 Support-vector networks. *Machine Learning* 20 (3):273-297.

608 Daware, A.V., R. Srivastava, A.K. Singh, S.K. Parida, and A.K. Tyagi, 2017 Regional Association
609 Analysis of MetaQTLs Delineates Candidate Grain Size Genes in Rice. *Front Plant Sci* 8:807.

610 Deo, R.C., G. Musso, M. Tasan, P. Tang, A. Poon *et al.*, 2014 Prioritizing causal disease genes using
611 unbiased genomic features. *Genome Biol.* 15 (12):534.

612 Dinka, S.J., M.A. Campbell, T. Demers, and M.N. Raizada, 2007 Predicting the size of the progeny
613 mapping population required to positionally clone a gene. *Genetics* 176 (4):2035-2054.

614 Fahlgren, N., M.A. Gehan, and I. Baxter, 2015 Lights, camera, action: high-throughput plant phenotyping
615 is ready for a close-up. *Curr. Opin. Plant Biol.* 24:93-99.

616 FAO, 2009 How to feed the world in 2050.

617 Gelfman, S., Q. Wang, K.M. McSweeney, Z. Ren, F. La Carpia *et al.*, 2017 Annotating pathogenic non-
618 coding variants in genic regions. *Nat Commun* 8 (1):236.

619 Gotz, S., J.M. Garcia-Gomez, J. Terol, T.D. Williams, S.H. Nagaraj *et al.*, 2008 High-throughput functional
620 annotation and data mining with the Blast2GO suite. *Nucleic Acids Res.* 36 (10):3420-3435.

- 621 Grant, C.E., T.L. Bailey, and W.S. Noble, 2011 FIMO: scanning for occurrences of a given motif.
622 *Bioinformatics* 27 (7):1017-1018.
- 623 Hattori, Y., K. Nagai, S. Furukawa, X.-J. Song, R. Kawano *et al.*, 2009 The ethylene response factors
624 *SNORKEL1* and *SNORKEL2* allow rice to adapt to deep water. *Nature* 460:1026.
- 625 He, H., 2014 Environmental Regulation of Seed Performance. Dissertation. Wageningen University,
626 Wageningen University.
- 627 Hindorf, L.A., P. Sethupathy, H.A. Junkins, E.M. Ramos, J.P. Mehta *et al.*, 2009 Potential etiologic and
628 functional implications of genome-wide association loci for human diseases and traits. *Proc. Natl.*
629 *Acad. Sci. U. S. A.* 106 (23):9362-9367.
- 630 Hormozdiari, F., G. Kichaev, W.Y. Yang, B. Pasaniuc, and E. Eskin, 2015 Identification of causal genes
631 for complex traits. *Bioinformatics* 31 (12):206-213.
- 632 Huang, Y.F., B. Gulko, and A. Siepel, 2017 Fast, scalable prediction of deleterious noncoding variants
633 from functional and population genomic data. *Nat Genet* 49 (4):618-624.
- 634 Jagadeesh, K.A., A.M. Wenger, M.J. Berger, H. Guturu, P.D. Stenson *et al.*, 2016 M-CAP eliminates a
635 majority of variants of uncertain significance in clinical exomes at high sensitivity. *Nat Genet* 48
636 (12):1581-1586.
- 637 Jones, C.E., A.L. Brown, and U. Baumann, 2007 Estimating the annotation error rate of curated GO
638 database sequence annotations. *BMC Bioinform.* 8:170.
- 639 Jones, P., D. Binns, H.Y. Chang, M. Fraser, W. Li *et al.*, 2014 InterProScan 5: genome-scale protein
640 function classification. *Bioinformatics* 30 (9):1236-1240.
- 641 Kircher, M., D.M. Witten, P. Jain, B.J. O'Roak, G.M. Cooper *et al.*, 2014 A general framework for
642 estimating the relative pathogenicity of human genetic variants. *Nat Genet* 46 (3):310-315.

643 Lee, I., B. Ambaru, P. Thakkar, E.M. Marcotte, and S.Y. Rhee, 2010 Rational association of genes with
644 traits using a genome-scale gene network for *Arabidopsis thaliana*. *Nat. Biotechnol.* 28 (2):114-
645 149.

646 Lee, I., Y.S. Seo, D. Coltrane, S. Hwang, T. Oh *et al.*, 2011 Genetic dissection of the biotic stress response
647 using a genome-scale gene network for rice. *Proc. Natl. Acad. Sci. U. S. A.* 108 (45):18548-18553.

648 Leinonen, P.H., D.L. Remington, J. Leppala, and O. Savolainen, 2013 Genetic basis of local adaptation and
649 flowering time variation in *Arabidopsis lyrata*. *Mol. Ecol.* 22 (3):709-723.

650 Liu, X.Q., L. Wang, S. Chen, F. Lin, and Q.H. Pan, 2005 Genetic and physical mapping of Pi36(t), a novel
651 rice blast resistance gene located on rice chromosome 8. *Mol. Genet. Genomics* 274 (4):394-401.

652 Lloyd, J.P., A.E. Seddon, G.D. Moghe, M.C. Simenc, and S.-H. Shiu, 2015 Characteristics of Plant
653 Essential Genes Allow for within- and between-Species Prediction of Lethal Mutant Phenotypes.
654 *Plant Cell* 27 (8):2133-2147.

655 Mackay, T.F.C., 2001 The genetic architecture of quantitative traits. *Annu. Rev. Genet.* 35:303-339.

656 Mackay, T.F.C., 2014 Epistasis and Quantitative Traits: Using Model Organisms to Study Gene-Gene
657 Interactions. *Nat. Rev. Genet.* 15 (1):22-33.

658 Mansueto, L., R.R. Fuentes, F.N. Borja, J. Detras, J.M. Abriol-Santos *et al.*, 2017 Rice SNP-seek database
659 update: new SNPs, indels, and queries. *Nucleic Acids Res.* 45 (D1):D1075-D1081.

660 Martin, A., and V. Orgogozo, 2013 The Loci of repeated evolution: a catalog of genetic hotspots of
661 phenotypic variation. *Evolution* 67 (5):1235-1250.

662 Mordelet, F., and J.P. Vert, 2011 ProDiGe: Prioritization Of Disease Genes with multitask machine
663 learning from positive and unlabeled examples. *BMC Bioinform.* 12:389.

664 Motte, H., A. Vercauteren, S. Depuydt, S. Landschoot, D. Geelen *et al.*, 2014 Combining linkage and
665 association mapping identifies *RECEPTOR-LIKE PROTEIN KINASE1* as an essential Arabidopsis
666 shoot regeneration gene. *Proc. Natl. Acad. Sci. U. S. A.* 111 (22):8305-8310.

667 Moyers, B.T., 2018 Camoco: A Net for the Sea of Candidate Genes. *Plant Cell* 30 (12):2889.

668 Ng, P.C., and S. Henikoff, 2003 SIFT: predicting amino acid changes that affect protein function. *Nucleic*
669 *Acids Res.* 31 (13):3812-3814.

670 Nuzhdin, S.V., C.L. Dilda, and T.F. Mackay, 1999 The genetic architecture of selection response.
671 Inferences from fine-scale mapping of bristle number quantitative trait loci in *Drosophila*
672 *melanogaster*. *Genetics* 153 (3):1317-1331.

673 Otto, S.P., and C.D. Jones, 2000 Detecting the undetected: Estimating the total number of loci underlying a
674 quantitative trait. *Genetics* 156 (4):2093-2107.

675 Panchy, N., M. Lehti-Shiu, and S.H. Shiu, 2016 Evolution of Gene Duplication in Plants. *Plant Physiol* 171
676 (4):2294-2316.

677 Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion *et al.*, 2011 Scikit-learn: Machine
678 Learning in Python. *J Mach Learn Res* 12:2825-2830.

679 Perez-Iratxeta, C., P. Bork, and M.A. Andrade, 2002 Association of genes to genetically inherited diseases
680 using data mining. *Nat. Genet.* 31:316.

681 Pollard, K.S., M.J. Hubisz, K.R. Rosenbloom, and A. Siepel, 2010 Detection of nonneutral substitution
682 rates on mammalian phylogenies. *Genome Res.* 20 (1):110-121.

683 Price, A.L., N.A. Zaitlen, D. Reich, and N. Patterson, 2010 New approaches to population stratification in
684 genome-wide association studies. *Nat. Rev. Genet.* 11 (7):459-463.

685 Ritchie, G.R., I. Dunham, E. Zeggini, and P. Flicek, 2014 Functional annotation of noncoding sequence
686 variants. *Nat Methods* 11 (3):294-296.

687 Schaefer, R., J.-M. Michno, J. Jeffers, O.A. Hoekenga, B.P. Dilkes *et al.*, 2018 Integrating co-expression
688 networks with GWAS to prioritize causal genes in maize. *Plant Cell*.

689 Schlapfer, P., P. Zhang, C. Wang, T. Kim, M. Banf *et al.*, 2017 Genome-Wide Prediction of Metabolic
690 Enzymes, Pathways, and Gene Clusters in Plants. *Plant Physiol* 173 (4):2041-2059.

691 Singleton, A.B., J. Hardy, B.J. Traynor, and H. Houlden, 2010 Towards a complete resolution of the
692 genetic architecture of disease. *Trends Genet* 26 (10):438-442.

693 Takahashi, R., K. Bashir, Y. Ishimaru, N.K. Nishizawa, and H. Nakanishi, 2012 The role of heavy-metal
694 ATPases, HMAs, in zinc and cadmium transport in rice. *Plant Signal Behav* 7 (12):1605-1607.

695 Tin Kam, H., 1998 The random subspace method for constructing decision forests. *IEEE Transactions on*
696 *Pattern Analysis and Machine Intelligence* 20 (8):832-844.

697 Touw, W.G., J.R. Bayjanov, L. Overmars, L. Backus, J. Boekhorst *et al.*, 2013 Data mining in the Life
698 Sciences with Random Forest: a walk in the park or lost in the jungle? *Brief Bioinform* 14 (3):315-
699 326.

700 Tuinstra, M.R., G. Ejeta, and P.B. Goldsbrough, 1997 Heterogeneous inbred family (HIF) analysis: a
701 method for developing near-isogenic lines that differ at quantitative trait loci. *Theor. Appl. Genet.*
702 95 (5):1005-1011.

703 Van de Velde, J., K.S. Heyndrickx, and K. Vandepoele, 2014 Inference of Transcriptional Networks in
704 Arabidopsis through Conserved Noncoding Sequence Analysis. *Plant Cell* 26 (7):2729-2745.

705 Weigel, D., and M. Nordborg, 2005 Natural variation in Arabidopsis. How do we find the causal genes?
706 *Plant Physiol.* 138 (2):567-568.

707 Weirauch, M.T., A. Yang, M. Albu, A.G. Cote, A. Montenegro-Montero *et al.*, 2014 Determination and
708 Inference of Eukaryotic Transcription Factor Sequence Specificity. *Cell* 158 (6):1431-1443.

709 Wellenreuther, M., and B. Hansson, 2016 Detecting Polygenic Evolution: Problems, Pitfalls, and Promises.
710 *Trends Genet* 32 (3):155-164.

711 Xu, K., X. Xu, T. Fukao, P. Canlas, R. Maghirang-Rodriguez *et al.*, 2006 Sub1A is an ethylene-response-
712 factor-like gene that confers submergence tolerance to rice. *Nature* 442 (7103):705-708.

713 Xu, S., 2003 Theoretical basis of the Beavis effect. *Genetics* 165 (4):2259-2268.

714 Yang, H.Y., A.Q. You, Z.F. Yang, F. Zhang, R.F. He *et al.*, 2004 High-resolution genetic mapping at the
715 Bph15 locus for brown planthopper resistance in rice (*Oryza sativa* L.). *Theor. Appl. Genet.* 110
716 (1):182-191.

717 Yin, Z.G., H.D. Qi, Q.S. Chen, Z.G. Zhang, H.W. Jiang *et al.*, 2017 Soybean plant height QTL mapping
718 and meta-analysis for mining candidate genes. *Plant Breed.* 136 (5):688-698.

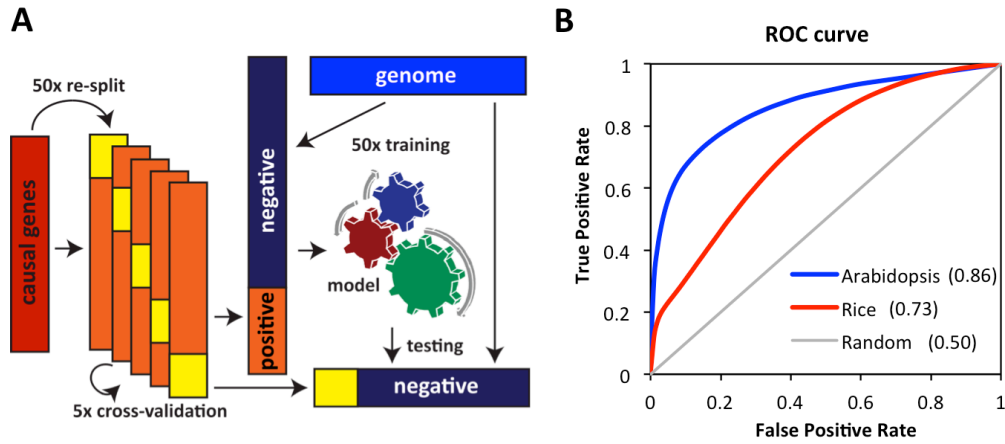
719 Yonemaru, J.-i., T. Yamamoto, S. Fukuoka, Y. Uga, K. Hori *et al.*, 2010 Q-TARO: QTL Annotation Rice
720 Online Database. *Rice* 3 (2):194-203.

721 Zhang, H., 2004 The Optimality of Naive Bayes. *Proc. FLAIRS*.

722 Zhao, Q., Q. Feng, H. Lu, Y. Li, A. Wang *et al.*, 2018 Pan-genome analysis highlights the extent of
723 genomic variation in cultivated and wild rice. *Nat Genet* 50 (2):278-284.
724
725

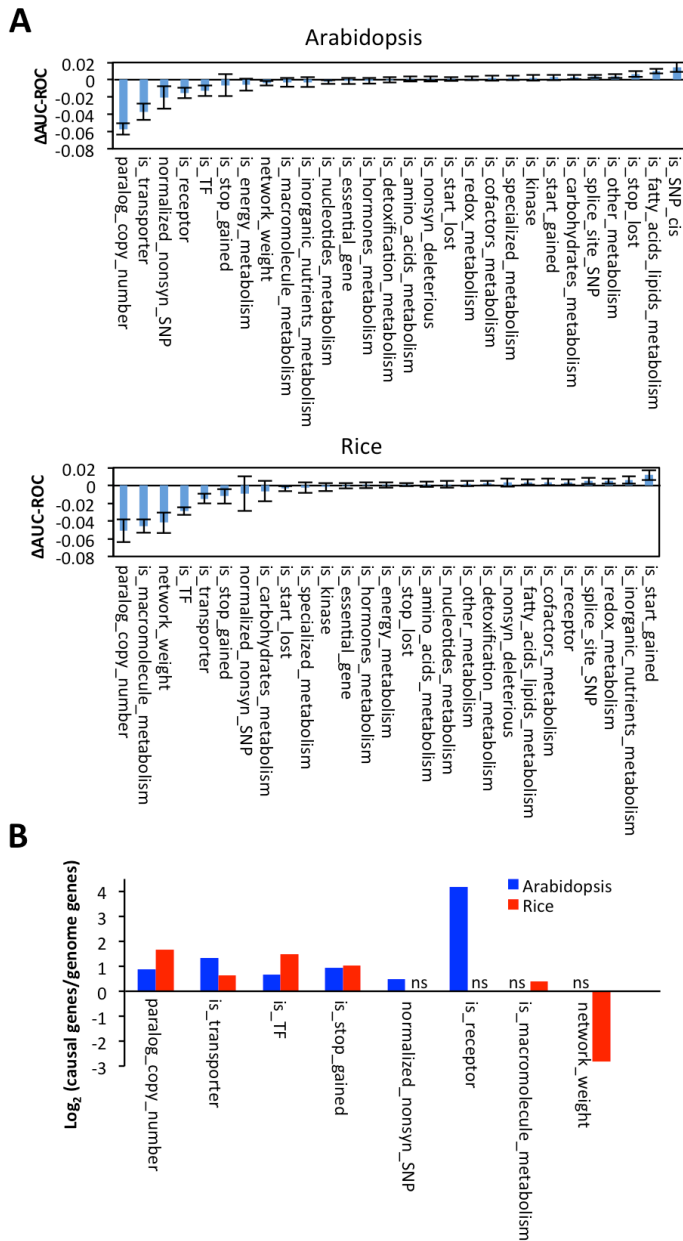
726

727 **Figures**



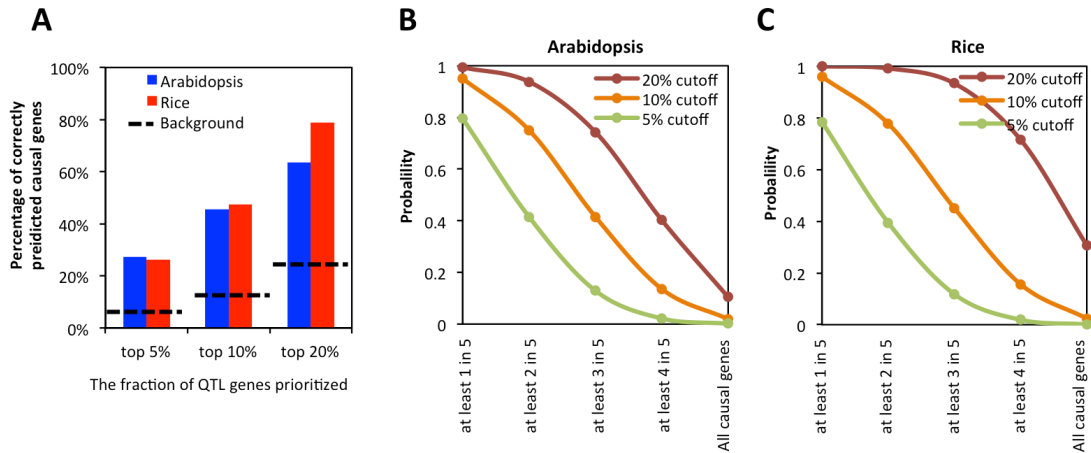
728

729 **Figure 1** Model training and optimization based on cross-validation. (A) model training
730 and cross-validation framework. We randomly selected negatives from the genome and
731 iterated to maximize the combinations of training and testing data. (B) The ROC curve of
732 Arabidopsis and rice models after parameter optimization. True and false positive rates
733 were based on the average of all iterations. The grey diagonal line indicates the expected
734 performance based on random guessing. The number in parentheses indicates Area Under
735 the ROC Curve (AUC-ROC)

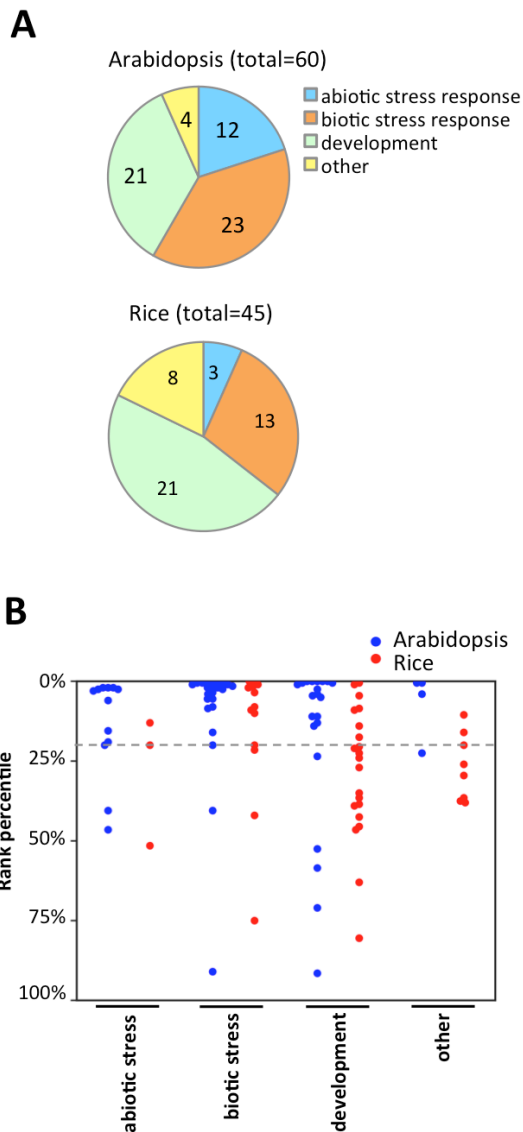


736
 737 **Figure 2** Important features of causal genes and their enrichment or depletion relative to
 738 the genome background. (A) Feature importance as indicated by the change of AUC-
 739 ROC (Δ AUC-ROC) when excluding each feature. The Δ AUC-ROC indicates the average
 740 value of all iterations. Error bars indicate standard deviation. The features with a name
 741 that starts with “is_” are binary variables. (B) The enrichment or depletion of the top 6
 742 features in Arabidopsis and rice models. The enrichment/depletion were indicated by the

743 ratio of causal genes to genome background. ns, not shown because the feature is not one
744 of the top 6 features in that species



745 **Figure 3** Model performance at different thresholds. (A) Percentage of recalled causal
746 genes of a single QTL at different rank thresholds. Dashed lines indicate the background
747 of random selections. (B-C) The probability of causal gene recall when analyzing
748 multiple QTLs simultaneously.



750 **Figure 4** Performance comparison across trait categories. (A) Trait categories of known
 751 causal genes from the training set. (B) The rank percentile of causal genes of different
 752 trait categories. Each causal gene and 200 neighboring genes were used as testing set
 753 only once. All other known causal genes were used for training. Each dot indicates a
 754 known causal gene. The grey dashed line indicates 20% rank percentile. The trait
 755 categories of causal genes are defined in Tables S1 and S2.

757