

1 MCRiceRepGP: a framework for identification of sexual
2 reproduction associated coding and lincRNA genes in rice.

3 Agnieszka A. Golicz¹

4 Prem L. Bhalla¹

5 Mohan B. Singh¹

6 ¹Plant Molecular Biology and Biotechnology Laboratory, Faculty of Veterinary and
7 Agricultural Sciences, University of Melbourne, Parkville, Melbourne, VIC, Australia.

8 Corresponding author: Agnieszka A. Golicz agnieszka.golicz@unimelb.edu.au

9 Corresponding author: Mohan B. Singh mohan@unimelb.edu.au

10 Number of tables: 3

11 Number of figures: 5 (all figures should be reproduced in colour)

12 Number of SI files: 2

13

14 **Significance statement**

15 Rice is a staple food crop plant for over half of the world's population and sexual reproduction
16 resulting in grain formation is a key process underpinning global food security. Despite
17 considerable research efforts, much remains to be learned about the molecular mechanisms
18 involved in rice sexual reproduction. We have developed MCRiceRepGP, a novel framework
19 which allows prediction of sexual reproduction associated genes using multi-omics data,
20 multicriteria decision analysis and machine learning. The genes identified and the methodology
21 developed will become a significant resource for the plant research community.

22 **Abstract**

23 Sexual reproduction in plants underpins global food production and evolution. It is a complex
24 process, requiring intricate signalling pathways integrating a multitude of internal and external
25 cues. However, key players and especially non-coding genes controlling plant sexual
26 reproduction remain elusive. We report the development of MCRiceRepGP a novel machine
27 learning framework, which integrates genomic, transcriptomic, homology and available
28 phenotypic evidence and employs multi-criteria decision analysis and machine learning to
29 predict coding and non-coding genes involved in rice sexual reproduction.

30 The rice genome was re-annotated using deep sequencing transcriptomic data from
31 reproduction-associated tissues/cell types identifying novel putative protein coding genes,
32 transcript isoforms and long intergenic non-coding RNAs (lincRNAs). MCRiceRepGP was
33 used for genome-wide discovery of sexual reproduction associated genes in rice; 2,275 protein-
34 coding and 748 lincRNA genes were predicted to be involved in sexual reproduction. The
35 annotation performed and the genes identified, especially the ones for which mutant lines with
36 phenotypes are available provide a valuable resource. The analysis of genes identified gives
37 insights into the genetic architecture of plant sexual reproduction. MCRiceRepGP can be used
38 in combination with other genome-wide studies, like GWAS, giving more confidence that the
39 genes identified are associated with the biological process of interest. As more data, especially
40 about mutant plant phenotypes will become available, the power of MCRiceRepGP will grow
41 providing researchers with a tool to identify candidate genes for future experiments.
42 MCRiceRepGP is available as a web application (<http://mcgplannotator.com/MCRiceRepGP/>)

43 **Key words**

44 function prediction, machine learning, *Oryza sativa*, re-annotation, rice, sexual reproduction;
45 lincRNA

46 **Introduction**

47 Sexual reproduction is a core process in the life cycle of a vast majority of eukaryotic
48 organisms. It is the main source of genetic diversity, which in turn allows for evolution and
49 adaptation. From economic perspective, sexual reproduction results in formation of edible fruit
50 and grains, underpinning crop yield and global food security. In plants, sexual reproduction is
51 initiated by the vegetative to reproductive phase transition, requiring intricate signalling
52 pathways integrating a multitude of internal and external cues. Upon commitment to flowering
53 the process involves the development of reproductive organs, successful completion of male
54 and female meiosis and fertilization, followed by embryonic development. Biological
55 processes involved in sexual reproduction consist of evolutionarily conserved core components
56 (for example, basic reproductive organ development including anthers and pistils and meiosis)
57 (Schurko and Logsdon, 2008, Wallace *et al.*, 2011, Gómez *et al.*, 2015) and a species-specific
58 regulatory level, for example the details of floral organ morphology and control of fine tuning
59 of timing of vegetative to reproductive phase transition (Jarillo and Piñeiro, 2011, Moyroud
60 and Glover, 2017). Knowledge of both, the level of conservation of core components and the
61 species-specific characteristics of reproductive processes is crucial for understanding of plant
62 fertility. Despite considerable research efforts the molecular basis of plant reproduction is not
63 yet fully understood.

64 Rice is an important cereal crop, providing staple food for over a half of the world's population.
65 It is a monocotyledonous plant species with a relatively compact genome. The rice genome
66 was one of the first plant genomes to be sequenced, providing a tremendous resource for plant
67 research community. However, despite considerable research efforts, many of the genes
68 involved in sexual reproduction remain uncharacterized (Kun *et al.*, 2013, Niu *et al.*, 2013, Fu
69 *et al.*, 2014, Rhee and Mutwil, 2014, Hu *et al.*, 2015, Yao *et al.*, 2017). Several computational
70 methods have been applied to improve understanding of gene functions. Studies of sequence
71 homology between the most extensively studied and functionally annotated proteome of
72 *Arabidopsis thaliana* and other species, including rice, allowed identification of genes with
73 conserved functions (Gómez *et al.*, 2015). Construction of co-expression networks allowed

74 identification of regulatory hubs involved in plant developmental processes, including anther
75 development (You *et al.*, 2016, de Luis Balaguer *et al.*, 2017, Lin *et al.*, 2017). Analysis of
76 expression profiles across tissues pinpointed genes with defined spatio-temporal expression
77 patterns, which could be involved in organ, tissue or cell-specific processes (Edwards and
78 Coruzzi, 1990). Studies of phenotypes of mutant lines provided annotation of genes with
79 unknown functions (Miyao *et al.*, 2003, Miyao *et al.*, 2007). Genome-wide studies of diversity
80 across hundreds of lines allowed identification of functionally important regions of increased
81 or reduced diversity, helping pinpoint genes which display high sequence conservation within
82 species (Alexandrov *et al.*, 2015, Tatarinova *et al.*, 2016).

83 Individually those approaches provide valuable insights into gene functions. The challenge is
84 to combine all the resources into a unified framework to produce a list of reliable candidate
85 genes involved in the biological process of interest (Troyanskaya *et al.*, 2003, Bradford *et al.*,
86 2010, Bargsten *et al.*, 2014). Our aim was to discover novel coding genes and lincRNAs
87 involved in rice sexual reproduction. To achieve that we have developed a set of rules to
88 prioritize the genes of interest and a novel method which combines information from gene
89 expression studies, sequence homology, known functional annotation, mutational data and
90 sequence diversity analysis. The method developed – MCRiceRepGP (Multi Criteria Rice
91 Reproductive Gene Predictor) predicts gene's potential for involvement in sexual reproduction
92 using available multi-omics data, multi-criteria decision analysis, and machine learning. We
93 applied the method to all rice genes and identified 2,275 protein coding and 748 lincRNA genes
94 involved in rice reproductive processes. The manuscript also presents the first study of
95 lincRNAs in plant gametes. A subset of the genes identified was linked to male and female-
96 specific plant fertility. Several genes linked to reproductive stage heat stress tolerance were
97 identified. For the purposes of the study, a full rice genome re-annotation using RNASeq
98 datasets from 11 tissues and cell types has been performed. MCRiceRepGP is available as a
99 web application (<http://mcgplannotator.com/MCRiceRepGP/>).

100 **Experimental procedures**

101 **Datasets used**

102 The rice genome assembly and annotation (MSU v7) and *A. thaliana* protein sequences (TAIR
103 10) were obtained from Phytozome v12.1 (Goodstein *et al.*, 2012). The RNASeq datasets were
104 downloaded from Sequence Read Archive (Table S1). To maximize mapping specificity and

105 minimize batch effects RNASeq from a minimum number of studies, covering maximum
106 number of reproductive and vegetative tissues with read length equal or longer than 100 base
107 were used. Phenotypic data for Tos17 rice mutant lines were downloaded from
108 <https://tos.nias.affrc.go.jp/> and the insertion coordinates were downloaded from
109 <http://orygenesdb.cirad.fr/>. Gene ontology (GO) annotation of *A. thaliana* genes were
110 downloaded from TAIR (ATH_GO_GOSLIM.txt, downloaded on: 20.07.2017) (Berardini *et*
111 *al.*, 2015).

112 **Parameters used**

113 Detailed commands for all the tools listed in the sections below can be found in Method S1.

114 **Genome reannotation**

115 The RNASeq reads were mapped to the reference genome using Hisat2 v2.0.5 (Kim *et al.*,
116 2015) and the parameters were adjusted for stranded libraries. Transcripts were assembled
117 separately for each library using StringTie v1.3.3b (Pertea *et al.*, 2015) and the parameters were
118 adjusted for stranded libraries. The annotations were then merged with the existing rice
119 annotation. lincRNAs were identified using procedure previously described (Golicz *et al.*,
120 2018b). In short, coding potential of genes was evaluated using Coding Potential Calculator 2
121 (Kang *et al.*, 2017) and homology to know protein coding genes. Transcripts were compared
122 using DIAMOND v0.8.24.86 (Buchfink *et al.*, 2015) blastx against NCBI RefSeq (O'Leary *et*
123 *al.*, 2016) protein database (downloaded on: 11.07.2017). A gene was considered coding if any
124 of the transcripts were classified as coding by CPC2 or had a significant match in RefSeq
125 database. The (long intergenic non-coding RNAs) lincRNAs were identified by comparing
126 positions of coding and non-coding genes using bedtools (Quinlan and Hall, 2010). All non-
127 coding genes which did not overlap any protein coding loci were classified as lincRNAs.

128 **Expression level evaluation**

129 The reads mapping to gene loci were counted using featureCounts v1.5.1 (Liao *et al.*, 2014).
130 The FPKM values were calculated as: $(10^9 * \text{fragments mapped to exons} / \text{assigned}$
131 $\text{fragments} * \text{total length of exons})$. The $\log_{1p}(\text{FPKM})$ values were adjusted for batch effects
132 using Combat v3.24.4 (Johnson *et al.*, 2007). The data used originated from three different
133 studies, which was accounted for during batch effect adjustment.

134 **Homology analysis**

135 For each coding gene representative (longest isoform) transcript was compared against the set
136 of *A. thaliana* proteins (longest isoforms) using NCBI blastx v2.6.0 (Camacho *et al.*, 2009).
137 GO annotations were transferred from *A. thaliana* genes to best matches (with lowest e-value)
138 among the rice genes.

139 **Community analysis**

140 Expression values were calculated by counting the number of reads mapping to each gene using
141 FeatureCounts v1.5.1 (Liao *et al.*, 2014). The Spearman correlations were computed using corr
142 function of psych package (Revelle, 2017). Top 5% of positive and negative correlations were
143 used to build a co-expression network using Mutual Rank method (Obayashi *et al.*, 2009) (MR
144 < 30). The Clique Percolation Method (Palla *et al.*, 2005) was used
145 (<https://sites.google.com/site/cliqueperccomp/>) to identify putative functional modules within
146 co-expression network. GO enrichment of nodes was calculated using topGO package v2.28.0
147 (Alexa *et al.*, 2006), using method 'weight' to adjust for multiple comparisons ($p < 0.01$).

148 **Diversity analysis**

149 The filtered SNP set (18 M) was downloaded from SNP-Seek database (Alexandrov *et al.*,
150 2015). The number of SNPs falling within exons of each genes was counted and divided by
151 total exon length of the gene as calculated by featureCounts. The gene was considered to be
152 low diversity if the SNP density was below half of the median SNP density calculated using
153 all genes.

154 **Process Involvement score parametrization**

155 The Process Involvement (PI) score has seven components, which are weighted differently
156 depending on their relative importance. Using knowledge of the field to supply probabilities
157 for analysis of networks has been previously successfully applied (Troyanskaya *et al.*, 2003).
158 The weights assume values between 0 and 1 and the values used were $\alpha=0.6$, $\beta=0.6$, $\gamma=0.4$,
159 $\delta=0.3$, $\epsilon=0.2$, $\zeta=0.1$. The phenotypic data (P, $\alpha=0.6$) and sequence homology with known
160 sexual reproduction regulators (H, $\beta=0.6$) were considered to be the most important pieces of
161 evidence and therefore were assigned the highest weight. Because one of the objectives of the
162 study was to uncover key regulators of sexual reproduction, participation in functional co-
163 expression modules was also considered important (CP, $\gamma=0.4$; CF, $\delta=0.3$). Sequence diversity
164 was also included, but given lower weighting. If genes had similar evidentiary support from

165 phenotypic and/or homology data and network-connectivity, genes with lower diversity are
166 hypothesized to be more likely regulators as transcription factors were shown to be the genes
167 with lowest diversity in the rice genome (Tatarinova *et al.*, 2016). Finally, the expression value
168 (EV, $\zeta=0.1$) was given lowest weighting to prevent it from over-powering the entire score.
169 Further details: Note S2.

170 **Classifiers**

171 Three classifier were tested: (1) the Naïve Bayes classifier as implemented in function
172 naiveBayes of package e1071 v1.6-8 (Meyer, 2017), (2) Classification Tree as implemented in
173 function rpart of package rpart v4.1-11 (Therneau *et al.*, 2017), (3) Logistic Regression as
174 implemented in method glm (R Core Team). Five-fold cross validation was used to measure
175 the concordance between classifier prediction and test datasets. Further details: Method S2,
176 Notes S3-S5.

177 **Test datasets**

178 Ten genes known to have confirmed crucial roles in sexual reproduction were chosen as test
179 dataset (Test Set 1) (Gómez *et al.*, 2015, Shi *et al.*, 2015a). Additionally, 781 genes implicated
180 to be involved in sexual reproduction (<https://funricegenes.github.io/>) and highly expressed in
181 reproductive tissues were used (Test Set 2).

182 **Fst score calculation**

183 The 18M SNP dataset downloaded from SNP-Seek database was used. The *japonica* sub-
184 population include temp and trop lines. The *indica* subpopulation included ind1, ind2 and ind3
185 lines. SNPs with minor allele frequency < 0.01 were remove from the dataset. Fst values were
186 calculated using vcftools, with window size of 100kb and step of 10kb. Windows which fell
187 within top 5% of highest Fst values (mean value) were retained, merged and compared with
188 positions of SexRep genes.

189 **Data availability**

190 Rice genome reannotation and files used as input for MCRiceRepGP can be found at:
191 <https://osf.io/78axs/>.

192 Source code can be obtained from: <https://github.com/agolicz/MCRiceRepGP> and
193 <https://github.com/agolicz/MCRiceRepGP-shiny>.

194 Web application can be found at: <http://mcgplannotator.com/MCRiceRepGP/>.

195 **Results**

196 **Rice genome re-annotation using RNASeq data**

197 The two available rice genome annotations (MSU-RAP and RAP-DB) were performed before
198 RNASeq data was widely available and gene evidentiary support relied mostly on ESTs, which
199 used to be derived from pools of samples, likely missing genes expressed at lower levels,
200 transiently expressed or in low abundance cell types (Note S1). This is an especially important
201 consideration while investigating sexual reproduction which depends on precise
202 spatiotemporal gene expression regulating cell fate commitment and specification involving a
203 small number of specialized cell types. Additionally, a mounting body of evidence accumulated
204 since the last rice genome annotation points to important roles of long non-coding RNAs in
205 sexual reproduction and those should also be included in the analyses (Golicz *et al.*, 2018a).
206 Long intergenic non-coding RNA (lincRNA) annotation in rice has been performed previously
207 (Zhang *et al.*, 2014, Wang *et al.*, 2015a), however the transcriptomes of egg, pollen sperm, and
208 vegetative cells were not included.

209 Accordingly, we updated the MSU-RAP annotation using RNASeq data from multiple rice
210 tissues and cell types (leaf, root, shoot, flower, seed, anther, pistils, sperm, cell, egg cell,
211 vegetative cell). The final annotation comprised 56,118 loci, including 46,149 protein-coding
212 and 9,969 lincRNA loci (Note S1, Table S2). The expression profile of newly discovered
213 putative protein-coding loci (7,218 genes, 65.9% containing open reading frame (ORF) >100
214 amino acids and 42.4% containing complete ORF > 100 amino acids) was analysed, and 80.9%
215 of genes showed highest expression levels in reproductive tissues, suggesting that a number of
216 reproduction related genes may be missing from the available MSU-RAP annotation (Fig. S1).
217 The updated annotation is well suited for the study of rice reproductive processes. It also
218 highlights the significance of including expression data from specialized organs and low
219 abundance cell types, especially those highly relevant to the study performed.

220 **The MCRiceRepGP method and its application for identification of reproduction** 221 **associated genes in rice**

222 Many publicly available rice genomic, transcriptomic and mutational datasets and databases
223 exist (Ware *et al.*, 2002, Droc *et al.*, 2006, Miyao *et al.*, 2007, Alexandrov *et al.*, 2015, Wang
224 *et al.*, 2015b). Using the updated genome annotation, these resources can be employed to help
225 identify genes associated with biological processes of interest, in this case sexual reproduction.

226 MCRiceRepGP uses information about seven features: tissue expression profile (tissue type
227 and expression levels), connectivity within co-expression network, co-expression hub
228 functional annotation, existing mutational data with phenotypic information, sequence
229 homology and single nucleotide polymorphism (SNP) diversity to calculate gene score and
230 predict whether the gene is involved in a biological process (Table 1, Fig 1).

231 **Tissue expression profile analysis and gene co-expression network construction**

232 *Tissue expression profile analysis*

233 Expression of all the rice genes across tissues was measured by quantifying number of RNASeq
234 reads mapped to each gene locus and calculating FPKM (fragments per kilo base per million)
235 value (Fig. S2). Because the dataset originated from several different studies, the expression
236 values were adjusted in order to remove batch effects (Johnson *et al.*, 2007). Genes which are
237 involved in a given process often show high or unique expression in related tissue (Wen *et al.*,
238 2016, Boyle *et al.*, 2017, Golicz *et al.*, 2018b). The samples were classified as either
239 representing vegetative or reproductive tissue (Table S4). For each gene, the tissue and tissue
240 type with highest expression levels observed was recorded. In total 72.6% (68.1% on non-
241 batch-adjusted data) genes had the highest expression level in reproductive tissue/cell type. A
242 high number of genes having peak expression in reproductive tissues is expected. Reproductive
243 processes are complex, requiring developmental transitions, cell fate decisions and formation
244 of multiple highly specialized cell types in male and female gametophytes, therefore are
245 expected to engage a multitude of genes.

246 *Co-expression network construction*

247 The FPKM expression values were used to calculate all-vs-all Spearman correlations and the
248 gene pairs within the top 5% (corresponding to minimum $\rho=0.725$ for positive correlations)
249 or bottom 5% (corresponding to maximum $\rho=-0.619$ for negative correlations) correlation
250 values were used to build a co-expression network containing 50,212 nodes and 678,548 edges.
251 Within the network, it is possible to identify sub-populations of tightly connected nodes – so
252 called communities (Acharya *et al.*, 2012). These likely correspond to functional modules
253 related to distinct biological roles. The whole network was analysed using Clique Percolation
254 Method (Palla *et al.*, 2005), detecting 5,791 communities (putative functional modules). The
255 modules were then functionally annotated using gene ontology (GO) enrichment analysis.
256 Following the procedure used in the MSU-RAP annotation, the rice genes were annotated with
257 GO terms corresponding to the most significant BLAST match in the *A. thaliana* proteome and

258 GO enrichment for each module was calculated using all genes as background. The
259 significantly enriched terms ($p < 0.01$) were assigned to modules as the functional annotation.
260 In total, 4,044 modules were annotated with at least one GO term. The assigned terms were
261 then manually inspected to identify key words/phrases associated with sexual reproduction
262 (Table S5). Nodes which were annotated with at least one GO term containing a key
263 word/phrase were annotated as associated with sexual reproduction (566 modules in total).

264 **Insertional mutant data**

265 To date, the most comprehensive rice mutant panel with a published collection of phenotypes
266 are the ~50,000 transposon Tos17 insertion lines (Miyao *et al.*, 2003, Miyao *et al.*, 2007). The
267 link between disruption of gene sequence and the observed phenotype can be indicative of gene
268 function. However, analysis of the dataset poses several challenges. Each line possesses more
269 than one transposon insertion within the genome, with up to 10 Tos17 insertions per line
270 (Miyao *et al.*, 2003). Not every insertion has a phenotypic manifestation, but in some cases, a
271 single insertion can cause multiple aberrant phenotypes. In fact, almost half of the lines showed
272 more than one phenotype (Miyao *et al.*, 2007). Because multiple Tos17 insertions within the
273 genome of one line exist establishing a correlation between insertion and phenotype is not
274 straight forward. However, if two or more lines have independent insertions in the same gene
275 and exhibit the same/similar phenotype, disruption of the gene is likely linked to the phenotype.
276 To facilitate detection of the most common phenotype associated with the insertion a more
277 fuzzy match was performed – the 49 phenotypes were split into more general categories:
278 reproductive timing, reproductive fertility, reproductive seed, reproductive organ, vegetative,
279 lethal and dwarf (Table S6).

280 The insertion sites derived from all the lines were compared with exonic positions of genes.
281 For each gene, all the lines which had an insertion within exons of the gene were extracted, and
282 the most common phenotype and phenotype category (reproductive timing, reproductive
283 fertility, reproductive seed, reproductive organ, vegetative, lethal and dwarf) were recorded. In
284 total, 3,252 genes could be assigned at least one line with phenotype, and for 1,295 the most
285 common phenotype was categorized as reproductive.

286 **Sequence homology analysis**

287 Sexual reproduction is a process conserved in eukaryotes, with a number of genes involved in
288 core processes, sharing sequence homology and conserved functions even among distantly
289 related species (Schurko and Logsdon, 2008, Wallace *et al.*, 2011, Gómez *et al.*, 2015). For

290 example, corresponding genes involved in anther and pollen development have been found
291 (Gómez *et al.*, 2015). Therefore, the functionality of *A. thaliana* homologs can help in the
292 prediction of roles of rice genes. The sequences of rice and *A. thaliana* genes were compared
293 and GO annotation was transferred from *A. thaliana* genes to best rice gene matches.
294 Additionally, the GO terms were compared with the list of key reproductive terms constructed
295 during functional annotation of the co-expression network. Genes which were annotated with
296 at least one GO term which contained a key word/term were annotated as associated with sexual
297 reproduction.

298 **Sequence diversity analysis**

299 Rice has the most extensive single nucleotide polymorphism database of any plants
300 (Alexandrov *et al.*, 2015). The database lists ~20 million SNPs discovered using genomic data
301 from ~3000 lines. Lower SNP density across genomic regions is associated with either
302 purifying selection or selective sweeps (Wollstein and Stephan, 2015). An analysis of SNP
303 diversity across the rice genome revealed that genes associated with regulation of transcription
304 have lower than average sequence diversity (Tatarinova *et al.*, 2016) and transcription factor
305 activity plays a key role in the control of biological processes. Furthermore, the known sexual
306 reproduction master regulators (Table 2) were enriched in genes with low sequence diversity
307 (Fisher exact, $p < 0.05$). The functional lncRNAs were also shown to have lower rates of
308 evolution compared to non-functional ones (Wen *et al.*, 2016). Overall, 21.23% genes were
309 identified as low diversity.

310 **Predicting gene's potential for involvement in sexual reproduction**

311 We devised a two-step approach in which we first apply Multi Criteria Decision Analysis
312 (MCDA) based Process Involvement score (PI score) and then use the top scoring genes as the
313 training dataset for Naïve Bayes classifier, which is in turn applied to the full set of genes. The
314 combination of the classification provided by Naïve Bayes and the PI score ranking allows
315 identification of most confident candidate genes involved in sexual reproduction.

316 *Process Involvement (PI) gene score*

317 The Process Involvement (PI) score is a single metric designed to measure gene's potential for
318 involvement in a biological process, in this case sexual reproduction. The score is inspired by
319 Multi Criteria Decision Analysis (MCDA), a decision-making strategy used in a variety of
320 settings from financial and urban planning to ecological risk assessment and medical
321 diagnostics (DCLG, 2009, Adunlin *et al.*, 2015, Linkov *et al.*, 2015). MCDA involves

322 combining multiple lines of evidence from different sources to aid complex problem solving.
323 A general feature of MCDA is: 1. scoring of the options 2. weighting of the scores depending
324 on their perceived importance. A similar approach can be used to evaluate the potential of
325 gene's involvement in biological process and prioritise genes with features of interest, given
326 diverse evidentiary support including expression, sequence homology, and diversity data. (Fig.
327 1, Table 1).

328 Seven features are taken into consideration and combined to provide a single score. The score
329 components were not weighted equally, ET, P and H contributing more to the score than CP,
330 CF, D, and EV (Table 1, Experimental Procedures, Note S2). Overall, the genes which scored
331 most favourably were: highly expressed in reproductive tissues, their disruption resulted in
332 reproductive phenotype, had homologues in *A. thaliana* annotated with functions in
333 reproduction, were highly connected in co-expression networks, had low sequence diversity
334 among rice lines. The score for protein coding genes and lincRNAs differed slightly. For
335 lincRNAs the homology term is ignored, as lincRNAs show little sequence conservation across
336 species and very few have functional annotation. The PI score was calculated for all rice genes,
337 resulting in a continuous distribution of scores (Fig. S3) and the genes were ordered by
338 descending PI score. The highest ranking (top scoring) genes were considered to have a high
339 potential for involvement in sexual reproduction.

340 *Using top scoring PI genes as training dataset and choosing the optimal machine learning*
341 *classifier*

342 The high and low PI scoring coding and lincRNA genes can be then used as training data for a
343 machine learning classification algorithm. The training dataset was composed of 200 coding
344 and 100 lincRNA top scoring genes (as an example of genes involved in sexual reproduction –
345 positive training dataset) and a random selection of 500 coding and 250 lincRNA genes from
346 the bottom 95% of the ranking (as an example of genes not involved in sexual reproduction –
347 negative training dataset). The GO enrichment analysis has shown the top 200 coding genes to
348 be highly enriched in functions related to sexual reproduction (Table S7), while the selection
349 of 500 genes from the bottom 95% showed no such enrichment (Table S8).

350 Three types of classifiers were tested (1) Naïve Bayes classifier, (2) Classification Tree, (3)
351 Logistic Regression. A machine learning based classifier essentially performs the following
352 task: 'Given a set of genes A, find all the genes with similar properties in a larger set B.' The
353 classifiers were evaluated with respect to Matthews correlation coefficient (MCC), sensitivity

354 and specificity (Fig 2a). Receiver operating characteristic (ROC) curves were also generated
355 by plotting sensitivity against (1 – specificity) and the area under the curve (AUC) was
356 compared (Fig 2a). To achieve a more balanced positive to negative set ratio, the negative
357 training set was composed of randomly selected subset of a larger number of genes and the
358 effect of the repeated selection on classifier performance was also tested (Fig 2a, Notes S3-S5).
359 Overall, the Naïve Bayes classifier outperformed the other two other classifiers across all the
360 metrics for both coding and lincRNA genes and was therefore chosen to perform the analysis
361 (Fig 2b and Fig 2c). The superior performance of Naïve Bayes classifier for biological
362 classification purposes using heterogenous data has been previously observed (Troyanskaya *et*
363 *al.*, 2003, Bradford *et al.*, 2010, Sperschneider *et al.*, 2016). Additionally, Naïve Bayes
364 classifier was shown to be not sensitive to the size of negative training set (Kurczab *et al.*,
365 2014, Kurczab and Bojarski, 2017) alleviating the potential effects of introducing artificial
366 positive to negative training set ratio (Libbrecht and Noble, 2015).

367 *Applying Naïve Bayes classifier*

368 Naïve Bayes Classifier identified, 2,275 coding genes and 748 lincRNAs as involved in sexual
369 reproduction (the genes identified by Naïve Bayes Classifier as involved in sexual reproduction
370 were termed SexRep genes, Table S9). Again, GO analysis of SexRep genes showed strong
371 enrichment of genes associated with sexual reproduction (Table S10). The number of genes
372 involved in different reproduction related processes improved markedly when comparing the
373 top 200 genes identified by PI score and the genes identified by Naïve Bayes classifier (for
374 example, 54 and 347 genes respectively annotated as possibly involved in flower development;
375 addition of 293 genes, addition of ~51 genes would be expected at random). The SexRep genes
376 include 198 genes for which Tos17 mutant phenotypes were available (162 coding genes and
377 36 lincRNAs) (Fig. 2d). The four most common phenotypes were low fertile, sterile,
378 germination rate and dwarf. This is consistent with observations that fertility and dwarf
379 phenotypes are highly correlated (Miyao *et al.*, 2007).

380 *Testing MCRiceRepGP predictions*

381 The classifier has been trained to prioritize certain features including: high expression in
382 reproductive tissues, homology to know *A. thaliana* proteins involved in reproduction and high
383 connectivity in co-expression network. We have compared the results with a set of genes
384 known to be crucial in rice sexual reproduction (Gómez *et al.*, 2015, Shi *et al.*, 2015a), which
385 broadly fit into the criteria set while training the classifier (Test Set 1, Table 2). The genes

386 represent a number of functional classes, including transcription factors, protein kinase, DNA
387 de-methylase, Polycomb group protein and an lincRNA mi-RNA sponge (Nonomura *et al.*,
388 2003, Ono *et al.*, 2012, Yun *et al.*, 2013, Pan *et al.*, 2014, Wang *et al.*, 2017) and are involved
389 in diverse processes including floral organ identity specification, floral patterning,
390 sporogenesis, gamete fusion, endosperm and embryonic development. The method has
391 classified all of those genes, including the lincRNA, as involved in reproduction. Additionally,
392 we have tested the results against a database of 781 genes implicated to be involved in sexual
393 reproduction (Test Set 2). Twenty eight percent of the Test Set 2 genes overlapped with SexRep
394 genes and such an overlap is unlikely to occur by chance alone (permutation test, $p < 0.01$),
395 confirming the suitability of the method for discovery of genes associated with sexual
396 reproduction. Disregarding genes found in Test Sets 1 and 2 the method identified 2,060 coding
397 and 747 lincRNA novel genes potentially involved in sexual reproduction.

398 **Characterization of genes predicted to be involved in sexual reproduction**

399 **Overall properties of genes predicted to be involved in sexual reproduction**

400 The 3,023 SexRep genes (2,275 protein-coding genes and 748 lincRNAs) were analysed in
401 more detail. Both coding and lincRNA SexRep genes showed an even distribution across
402 chromosomes (Fig. 3a). The protein coding genes had higher overall expression levels when
403 compared to lincRNAs, which is consistent with observations in rice and other plant species
404 (Fig. 3b and Fig. 3c) (Zhang *et al.*, 2014, Wang *et al.*, 2015a). Analysis of tissue expression
405 patterns of coding and lincRNA SexRep genes revealed that the highest proportion of genes
406 had peak expression in egg and sperm cells respectively (Fig. 3b and Fig. 3c). Molecular
407 function enrichment (Table S11) of the protein coding-genes showed them to be involved in
408 protein binding, transcription factor activity, kinase activity and chromatin binding. Overall,
409 54.4% of the protein SexRep genes had no detectable similarity to *A. thaliana* genes involved
410 in sexual reproduction, but 59.2% were found in communities annotated with reproductive
411 functions. Similarly, 61% of lincRNAs were found in communities annotated with
412 reproductive functions (Fig 2c).

413 **Top candidate SexRep genes have diverse functional annotation**

414 The SexRep genes can be ranked by PI score to identify most confident candidates. Top 10
415 SexRep genes (as ranked by PI score) were investigated in more detail (Table 3). Analysis of
416 *A. thaliana* homologs suggests a diversity of molecular functions including protein kinases,
417 transcription factor, UDP-glucose phosphorylase, histidinol dehydrogenase and ferritin. The

418 genes appear to be involved in a range of processes from floral organ specification, cell cycle
419 regulation, pollen maturation to pollen tube guidance. The most common phenotype found
420 among the top ten genes was low fertility. To our knowledge four of the genes
421 (LOC_Os01g68870, LOC_Os02g02560, LOC_Os06g08380 and LOC_Os12g10540) have
422 already been characterized, confirming their involvement in sexual reproduction and influence
423 on fertility (Yao *et al.*, 2017).

424 **SexRep genes have distinct tissue expression profiles**

425 Genes which show unique or high activity in a given tissue are considered to be likely to
426 contribute to the relevant biological processes (Wen *et al.*, 2016, Boyle *et al.*, 2017, Golicz *et*
427 *al.*, 2018b). We investigated overall expression profiles of SexRep genes which show peak
428 expression in a given tissue/cell type (Fig. 4a). Principal components analysis (PCA) shows
429 clear clustering of both coding and lincRNA genes with peak expression in flower bud/flower,
430 egg cells, pollen sperm cells and vegetative cells (Fig. 4a) suggesting that the genes may be
431 involved in common biological processes. Protein coding SexRep genes have overall lower
432 expression specificities (show broad expression across tissues/cell types), when compared to
433 lincRNAs (lincRNAs are expressed in a limited number of tissues/cell types, Fig. 4b), which
434 again is consistent with observations in other species (Golicz *et al.*, 2018a). For example, sperm
435 cell SexRep protein-coding genes have one of the lowest median values of expression
436 specificity index, while the lincRNA genes have the highest.

437 **Expression profile of SexRep genes suggests genes involved in male and female fertility**

438 Sexual reproduction requires formation of reproductive structures including flower, anthers
439 and pistils as well as successful male and female gametophyte development and fertilization.
440 Defects which are sex specific will result in aberrant male or female fertility. We have
441 investigated expression patterns of SexRep genes associated with fertility phenotype (Fig. 4c).
442 Majority of the genes show sex-specific preferential expression. The genes associated with
443 fertility phenotype show a clear split into three groups (1) genes with preferential expression
444 in anthers and vegetative cells (2) genes with preferential expression in sperm cells and (3)
445 genes with preferential expression in pistils and egg cells. Genes with preferential expression
446 in male or female organs are potential contributors to sex-specific fertility.

447 **A subset of SexRep genes shows population differentiation between *japonica* and *indica***
448 **genotypes**

449 In rice, there is an ancient and well-established divergence between two subspecies *japonica*
450 and *indica* and the subpopulations are easily distinguishable based on their DNA sequence
451 (Garris *et al.*, 2005). The subspecies also display phenotypic differences. For example, the
452 *indica* lines being overall more heat tolerant than the *japonica* lines (Jagadish *et al.*, 2007, Zhao
453 *et al.*, 2016), although heat tolerant lines exist in both sub-populations. Heat stress is known to
454 reduce rice fertility with flowering (anthesis and fertilization) being the most susceptible stages
455 of development (Jagadish *et al.*, 2007). The large polymorphism database available for rice
456 (Alexandrov *et al.*, 2015) allows detailed genome-wide studies of differences between
457 subspecies. The pairwise differentiation index (Fst) can be calculated between subpopulations,
458 used to pinpoint regions of highest sequence diversity and find loci contributing to differences
459 in phenotypes (Zhou *et al.*, 2015). In total, 288 SexRep protein coding genes fell within
460 genomic regions corresponding to the top 5% of Fst values calculated between *japonica* and
461 *indica* genotypes (Fig. 3a). GO enrichment analysis of those genes points to significant
462 enrichment of genes associated with anther dehiscence ($p = 0.0048$, Table S12 and Table S13).
463 Poor anther dehiscence is in turn known to be the leading cause of spikelet sterility induced by
464 high temperatures due to poor efficiency of pollen delivery to stigma (Jagadish *et al.*, 2010,
465 Zhao *et al.*, 2016). High differentiation of anther dehiscence related genes is consistent with
466 observations of differential heat tolerance of *indica* and *japonica* sub-species.

467 **Several SexRep genes overlap loci associated with sterility in rice**

468 The method used for detection of SexRep genes can also be used to enhance findings of genome
469 wide association studies (GWAS). GWAS have been successfully used to uncover genomic
470 regions containing loci associated with agronomic traits (Huang *et al.*, 2010, Yano *et al.*, 2016).
471 Although high density SNP maps give good resolution to GWAS studies, usually several
472 candidate genes within the region of interest are identified (Dingkuhn *et al.*, 2017). Usage of
473 additional lines of evidence such as the ones used for identification of SexRep genes can help
474 point to more confident candidates within the sections of the genome identified by GWAS. We
475 have compared the genomic locations of recently identified SNPs linked to heat stress
476 associated sterility in rice (Dingkuhn *et al.*, 2017) with coordinates of SexRep genes and
477 identified six genes potentially related to sterility (Table S14). The number of SexRep genes
478 found in vicinity of sterility associated SNPs (closer to the SNP than any other gene) was higher
479 than it would be expected by chance (Chi Square, $p < 0.01$).

480 **Discussion**

481 Despite considerable research efforts genes controlling sexual reproduction in plants remain
482 enigmatic. Computational biology approaches can provide new insights by combining and
483 analysing large-scale data from a number of sources, including genomic, transcriptomic and
484 mutational datasets. The main challenge is the effective integration of all the information
485 available. In this study, Process Involvement (PI) score and Naïve Bayes Classifier were
486 applied to identify genes involved in sexual reproduction. MCRiceRepGP depends on seven
487 features which describe the gene in terms of expression profile, biological network
488 connectivity, homology with known sexual reproduction regulators and overall sequence
489 diversity. MCRiceRepGP was applied to protein coding genes as well as non-coding RNA loci
490 and identified three thousand protein coding genes and lincRNA loci involved in sexual
491 reproduction. Analysis of all protein coding genes predicted to be involved in sexual
492 reproduction (SexRep genes) highlighted genes involved in protein binding, transcription
493 factor and kinase activity. The most common mutant phenotype associated with both coding
494 and lincRNA SexRep genes was low fertility. The top SexRep protein coding genes had diverse
495 functional annotations and are implicated in processes from floral organ specification, pollen
496 development to pollen tube guidance. The genes identified are valuable resource providing
497 potential targets for further experiments, including many long non-coding RNAs. Previous
498 studies have shown long non-coding RNAs to play active roles in reproductive processes and
499 the candidates identified in this study can open new avenues for rice research.

500 In this analysis MCRiceRepGP was parametrized to favour genes highly or specifically
501 expressed in reproductive organs and with sequence homology to *A. thaliana* genes. However,
502 alternative parameters can be chosen depending on the experimental goals. Other mutant lines
503 can also be utilized. Recently a comprehensive library of neutron mutants became available,
504 although no phenotypes have yet been recorded (Li *et al.*, 2017). Additionally, looking at
505 individual components of the PI score can also point to genes of interest. For instance, looking
506 only at genes which do not have homologs in *A. thaliana*, could help uncover rice specific
507 regulators.

508 The method can be used in conjunction with other genome wide analyses. A number of
509 genome-wide screens which help identify genomic regions associated with traits exist. These
510 include genome-wide association studies (GWAS) to identify loci linked to traits of interest or
511 calculation of fixation index (Fst) between sub-populations and identification of genomic

512 regions with high and low differentiation. However, regions identified usually contain multiple
513 genes, and it is not clear which one affects the trait. For example, in a recent GWAS study
514 genes within $\pm 100\text{kb}$ of associated polymorphism were considered (Dingkuhn *et al.*, 2017)
515 and the F_{st} values are also calculated for $\sim 100\text{kb}$ windows (Zhou *et al.*, 2015). Often sequence
516 homology only is used, but combining multiple lines of evidence can give more confident
517 candidate gene predictions. Comparison of genome coordinates of SexRep genes against
518 regions of high differentiation between *japonica* and *indica* genotypes revealed
519 overrepresentation of genes associated with pollen release from anther, while comparison with
520 GWAS data identified six genes potentially related to sterility.

521 A web application which implements MCRiceRepGP has been made available (Fig. 5). The
522 application allows building of new classifiers by varying of PI score parameters, key words
523 and classifier features. The results can be browsed online and are available for download.

524 **Conclusion**

525 We have developed MCRiceRepGP – a method which combines evidence from heterogenous
526 data sources for identification of novel genes involved in rice sexual reproduction. An easy to
527 use web application has been made available and allows building of different classifiers.
528 Additionally, for this study an updated rice genome annotation has been generated using deep
529 sequencing data from reproductive tissues and cell types. The methodology developed, the
530 putative reproduction associated genes and especially lincRNAs identified using
531 MCRiceRepGP as well as the new rice genome annotation provide a valuable resource for
532 further studies of rice sexual reproduction. Identification of previously unannotated genes from
533 sexual reproduction specific tissues highlights the importance of including expression data
534 from specialized organs and low abundance cell types in the genome annotation efforts. The
535 novel sexual reproduction associated genes and lincRNAs described in the study provide
536 targets for future research efforts. The method described may become an inspiration and an
537 example of how different types of data can be integrated to predict most confident candidate
538 genes and future research targets.

539 **Acknowledgements**

540 This research was supported by Melbourne Bioinformatics at the University of Melbourne,
541 project UOM0033. The research was supported by ARC Discovery grant DP0988972 and the
542 University of Melbourne McKenzie Postdoctoral Fellowship.

543 **Authors' contributions**

544 AAG designed and performed the experiments, wrote the manuscript. MBS conceived
545 research, designed the experiments, wrote the manuscript. PLB conceived research.

546 Table 1 Features taken into account when evaluating the PI score.

Feature	General feature description	As applied for identification of sexual reproduction genes	Possible values	Parameter (feature weight)	Parameter values
Expression type (ET)	Is the highest expression recorded in a relevant tissue type?	Is the highest expression recorded in reproductive tissue?	0 – no 1 – yes	α	0.6
Phenotype category (P)	Is the most common mutant phenotype consistent with the highest expression tissue type?	Is the most common phenotype associated with the transposable element insertion in the gene reproductive only?	0 – no 1 – yes		
Sequence homology (H)	Is the homologous <i>A. thaliana</i> gene annotated with relevant functions?	Is the homologous <i>A. thaliana</i> gene annotated with reproductive function?	0 – no 1 – yes	β	0.6
Community participation (CP)	Is the gene found within a community in co-expression network?	Is the gene found within a community in co-expression network?	0 – no 1 – yes	γ	0.4
Community function (CF)	Does the gene belong to a community annotated with relevant functions?	Does the gene belong to a community annotated with reproductive function?	0 – no 1 – yes	δ	0.3
Sequence diversity (D)	Does the gene display low sequence diversity within the species?	Does the gene display low sequence diversity within the species?	0 – no 1 – yes	ϵ	0.2
Expression value (EV)	The FPKM value for the gene in the tissue with highest expression	The FPKM value for the gene in the tissue with highest expression	EV	ζ	0.1

547

548

549 Table 2 PI scores for known genes involved in sexual reproduction. Rep – predicted to be involved in sexual reproduction by MCRiceRepGP.

550 MCRiceRepGP was not tested on LDMAR, a lincRNA known to be involved in rice sexual reproduction as it was not found in the annotation.

Gene name	Gene ID	MSU-RAP Locus ID	Function	ET	P	H	CP	CF	D	Log1p(EV)	PI	MCRiceRepGP
MADS3	OSATST00001046	LOC_Os01g10504	Transcription factor	1	0	1	0	0	1	5.16	1.316	Rep
MADS58	OSATST00036993	LOC_Os05g11414	Transcription factor	1	0	1	1	1	0	4.85	1.785	Rep
MADS15	OSATST00045376	LOC_Os07g01820	Transcription factor	1	0	1	1	0	0	3.68	1.368	Rep
MADS1	OSATST00025799	LOC_Os03g11614	Transcription factor	1	0	1	0	0	1	5.21	1.321	Rep
DL	OSATST00025795	LOC_Os03g11600	Transcription factor	1	0	1	0	0	1	4.85	1.285	Rep
MSP1	OSATST00006904	LOC_Os01g68870	Kinase/Signalling	1	1	1	1	1	1	2.41	2.341	Rep
OsRos1a	OSATST00001213	LOC_Os01g11900	DNA demethylation	1	0	1	1	0	1	3.71	1.571	Rep
OsFIE2	OSATST00049934	LOC_Os08g04270	Polycomb silencing	1	0	1	1	1	0	3.54	1.654	Rep
HAP2	OSATST00037441	LOC_Os05g18730	Gamete fusion	1	0	1	1	1	0	4.55	1.755	Rep
Osa-eTM160	NC_OSATST00025950	N/A	miRNA sponge	1	0	N/A	1	0	1	2.02	0.802	Rep

551

552

553 Table 3 Top ten SexRep genes with highest PI scores.

	Gene ID	MSU-RAP Locus ID	PI score	Arabidopsis homolog	Arabidopsis protein name	Arabidopsis protein function/mutant phenotype	Confirmed function in rice	Rice mutant phenotype
Coding	OSATST00018594	LOC_Os02g02560	2.702	AT5G17310	UGP2	UDP-glucose phosphorylase [TAIR website (Berardini <i>et al.</i> , 2015)].	Preferentially expressed in pollen and plays key role during pollen maturation (Mu <i>et al.</i> , 2009). Shows differential expression in varieties differing in male fertility (Pan <i>et al.</i> , 2014).	Sterile
	OSATST00027308	LOC_Os03g24170	2.52	AT3G56960	PIP5K4	Phosphatidylinositol-4-phosphate 5-kinase activity. Key role in pollen tip growth. Overexpression of this gene leads to altered pollen tube morphology [TAIR website].		Low fertile
	OSATST00041382	LOC_Os06g08380	2.513	AT2G13680	CALS5	Responsible for the synthesis of callose deposited at the primary cell wall of meiocytes, tetrads and microspores. Required for exine layer formation during microgametogenesis and for pollen viability. Highest expression in meiocytes, tetrads, microspores and mature pollen [TAIR website].	OsGSL5 plays essential role in rice male fertility (Shi <i>et al.</i> , 2015b).	Low fertile
	OSATST00001380	LOC_Os01g13190	2.47	AT5G63890	HISN8	Histidinol dehydrogenase. Identified in screen of male gametophytic mutants [TAIR website].		Low fertile

	OSATST00023893	LOC_Os02g53720	2.447	AT2G26330	ER	Homologous to receptor protein kinases. Involved in specification of organs originating from the shoot apical meristem [TAIR website].		Low fertile
	OSATST00015635	LOC_Os12g10540	2.407	AT4G18960	AG	Floral homeotic gene encoding a MADS domain transcription factor [TAIR website]. Specifies floral meristem and carpel and stamen identity [TAIR website].	Controls ovule identity in rice (Dreni <i>et al.</i> , 2007). Involved in meristem determinacy (Dreni <i>et al.</i> , 2011).	Low fertile
	OSATST00018657	LOC_Os02g03060	2.384	AT3G48750	CDC2	A-type cyclin-dependent kinase. Loss of function phenotype has reduced fertility [TAIR website].		Low fertile
	OSATST00010933	LOC_Os11g01530	2.372	AT2G40300	FER4	Ferritins are essential to protect cells against oxidative damage, but they do not constitute the major iron pool [TAIR website].		
	OSATST00006904	LOC_Os01g68870	2.341	AT5G07280	EMS1	A putative leucine-rich repeat receptor protein kinase. Controls somatic and reproductive cell fates in anther [TAIR website].	Necessary to restrict the number of cells entering into male and female sporogenesis and to initiate anther wall formation in rice (Nonomura <i>et al.</i> , 2003).	Sterile
	OSATST00035108	LOC_Os04g52450	2.334	AT3G22200	HER1	Mediates pollen tube guidance [TAIR website].		Low fertile
Non-coding	NC_OSATST00049879		1.494	N/A	N/A	N/A	N/A	Low fertile
	NC_OSATST00032850	LOC_Os04g31740	1.466	N/A	N/A	N/A	N/A	Low fertile
	NC_OSATST00036949	LOC_Os05g10910	1.416	N/A	N/A	N/A	N/A	Germination rate
	NC_OSATST00000613	LOC_Os01g06620	1.326	N/A	N/A	N/A	N/A	
	NC_OSATST00032898	LOC_Os04g32180	1.293	N/A	N/A	N/A	N/A	Low fertile
	NC_OSATST00056573		1.293	N/A	N/A	N/A	N/A	Low fertile

	NC_OSATST00041036		1.29	N/A	N/A	N/A	N/A	Low fertile
	NC_OSATST00042178		1.284	N/A	N/A	N/A	N/A	Low fertile
	NC_OSATST00036564		1.283	N/A	N/A	N/A	N/A	Vivipary
	NC_OSATST00038962	LOC_Os05g36994	1.261	N/A	N/A	N/A	N/A	

554

555 Fig. 1. **MCRiceRepGP method overview.** Seven features (ET – expression type, P –
556 phenotype category, H – sequence homology, CP – community participation, CF – community
557 function, D – sequence diversity, EV – expression value) are used when evaluating a gene’s
558 potential for involvement in a biological process. The features are scored and weighted and the
559 Process Involvement (PI) score is calculated. The top and bottom scoring genes are used as
560 positive and negative training set to build Naïve Bayes classifiers for coding and lincRNA genes.
561 The classifiers are then used to identify a final set of genes involved in a given process. The
562 values of parameters used to identify genes involved in sexual reproduction are presented in
563 square brackets.

564 Fig. 2. **Comparison of the tested classifiers and the characteristics of the final Naïve Bayes**
565 **classifier used for the analysis.** (a) Three popular classifiers were tested: Naïve Bayes
566 classifier, Classification Tree, and Logistic Regression. The performance measures used to
567 assess the classifiers were: area under the receiver operating characteristic (ROC) curve (AUC)
568 – interpreted as the ability of the classifier to distinguish between the two cases, MCC –
569 Matthews correlation coefficient, sensitivity and specificity. The Naïve Bayes classifier was
570 the top performing algorithm. The positive training sets were the top 200 and 100 coding and
571 lincRNA genes as ranked by PI score. The negative training sets were the 500 and 250 genes
572 randomly drawn from the bottom 95% of PI score gene ranking. In total, 200 negative training
573 sets for coding and lincRNA genes were drawn and 5-fold cross validation for each negative
574 set was performed (3×5×200 classifiers built for coding and lincRNAs genes). (b) The ROC
575 curves along with other performance measures for the final Naïve Bayes classifiers for coding
576 genes (200 top PI scoring coding genes as positive training set, random selection of 500 coding
577 genes from the bottom 95% of PI score gene ranking as negative training set) and lincRNA
578 genes (100 top PI scoring lincRNA genes as positive training set, random selection of 250
579 lincRNA genes from the bottom 95% of PI score gene ranking as negative training set). The
580 performance of the classifiers was tested using 5-fold cross validation, the values provided are
581 means. (c) Proportion of coding and lincRNA genes in positive training set, negative training
582 set and final predicted SexRep genes which had value ‘1’ for the six binary features listed in
583 Table 1 (ET – expression type, P – phenotype category, H – sequence homology, CP –
584 community participation, CF – community function, D – sequence diversity). Vast majority of
585 coding and lincRNA genes had peak expression in a reproductive tissue/cell type and belonged
586 co-expression module(s). Many coding genes showed homology to known sexual reproduction
587 regulators and had low sequence diversity. (d) Ten most common insertional mutant

588 phenotypes for coding and lincRNA SexRep genes. The most common phenotype was low
589 fertility and sterility.

590 **Fig. 3. The landscape of SexRep genes.** (a) Circular plot presenting SexRep gene distribution
591 along the rice genome. From the outside ring: (1) coding SexRep genes, (2) lincRNA SexRep
592 genes, (3) Fst index between *japonica* and *indica* sub-populations, calculated for 100 kb
593 overlapping windows with a step of 10kb, (4) SexRep genes falling within regions of 5%
594 highest Fst values (5) SexRep genes overlapping sterility associated loci identified in GWAS.
595 (b,c) Heatmaps presenting expression of SexRep genes across tissues/cell types. Coding genes
596 have higher overall expression values. Many of the lincRNA genes are expressed in sperm
597 cells. Pie charts on top of heatmaps summarize the number of genes with peak expression in a
598 given tissue/cell type.

599 **Fig. 4. Overall expression patterns of SexRep genes with peak expression in a given**
600 **tissue/cell type.** (a) PCA analysis of coding and lincRNA gene expression values across
601 tissues, each point corresponds to a gene and is coloured according to tissue/cell type in which
602 the gene had peak expression value. Genes with common peak expression tissue/cell type
603 cluster together – show similar overall expression patterns, which suggests involvement in
604 common pathways/biological processes. (b) The box plots present tissue specificity index
605 (*Tau*) of genes having highest expression point in a given tissue/cell type. Difference between
606 coding and lincRNA genes can be observed. For example, the protein coding genes with peak
607 expression in sperm cells have the lowest tissue expression specificities, while the lincRNA
608 genes have the highest. The nested nature of sampling (for example, sperm cells are found
609 within anthers, which in turn are found within flowers) could affect specificity calculations.
610 Therefore, specificity indices were calculated twice, first using all samples (classic) and then
611 adjusting for sample structure (adjusted). However, in both cases similar patterns were
612 observed. (c) Heatmap presenting expression patterns of SexRep genes associated with fertility
613 phenotype. The genes show sex-specific expression.

614 **Fig. 5. Screen shot of MCRiceRepGP web app results.** The panel on the left side allows the
615 user to control gene type, key words, PI score parameters and features to be included in the
616 classifier. Results are displayed on the right-side panel. Results include classifier statistics,
617 classifier ROC curve, control classifier ROC curve and the table with the final results.

618 References

- 619 **Acharya, L., Judeh, T. and Zhu, D.** (2012) A survey of computational approaches to
620 reconstruct and partition biological networks. In *Statistical and Machine Learning*
621 *Approaches for Network Analysis*: John Wiley & Sons, Inc., pp. 1-43.
- 622 **Adunlin, G., Diaby, V., Montero, A.J. and Xiao, H.** (2015) Multicriteria decision analysis in
623 oncology. *Health expectations : an international journal of public participation in*
624 *health care and health policy*, **18**, 1812-1826.
- 625 **Alexa, A., Rahnenführer, J. and Lengauer, T.** (2006) Improved scoring of functional groups
626 from gene expression data by decorrelating GO graph structure. *Bioinformatics*, **22**,
627 1600-1607.
- 628 **Alexandrov, N., Tai, S., Wang, W., Mansueto, L., Palis, K., Fuentes, R.R., Ulat, Victor J.,**
629 **Chebotarov, D., Zhang, G., Li, Z., Mauleon, R., Hamilton, Ruaraidh S. and**
630 **McNally, K.L.** (2015) SNP-Seek database of SNPs derived from 3000 rice genomes.
631 *Nucleic Acids Research*, **43**, D1023-D1027.
- 632 **Bargsten, J.W., Severing, E.I., Nap, J.-P., Sanchez-Perez, G.F. and van Dijk, A.D.J.**
633 (2014) Biological process annotation of proteins across the plant kingdom. *Current*
634 *Plant Biology*, **1**, 73-82.
- 635 **Berardini, T.Z., Reiser, L., Li, D., Mezheritsky, Y., Muller, R., Strait, E. and Huala, E.**
636 (2015) The Arabidopsis information resource: making and mining the 'Gold Standard'
637 annotated reference plant genome. *Genesis (New York, N.Y. : 2000)*, **53**, 474-485.
- 638 **Boyle, E.A., Li, Y.I. and Pritchard, J.K.** (2017) An expanded view of complex traits: from
639 polygenic to omnigenic. *Cell*, **169**, 1177-1186.
- 640 **Bradford, J.R., Needham, C.J., Tedder, P., Care, M.A., Bulpitt, A.J. and Westhead, D.R.**
641 (2010) GO-At: in silico prediction of gene function in *Arabidopsis thaliana* by
642 combining heterogeneous data. *The Plant Journal*, **61**, 713-721.
- 643 **Buchfink, B., Xie, C. and Huson, D.H.** (2015) Fast and sensitive protein alignment using
644 DIAMOND. *Nat Meth*, **12**, 59-60.
- 645 **Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K. and**
646 **Madden, T.L.** (2009) BLAST+: architecture and applications. *BMC Bioinformatics*,
647 **10**, 421-421.
- 648 **DCLG** (2009) Multi-criteria analysis: a manual. London: Department for Communities and
649 Local Government.
- 650 **de Luis Balaguer, M.A., Fisher, A.P., Clark, N.M., Fernandez-Espinosa, M.G., Möller,**
651 **B.K., Weijers, D., Lohmann, J.U., Williams, C., Lorenzo, O. and Sozzani, R.** (2017)
652 Predicting gene regulatory networks by combining spatial and temporal gene
653 expression data in Arabidopsis root stem cells. *Proceedings of the National Academy*
654 *of Sciences*.
- 655 **Dingkuhn, M., Pasco, R., Pasuquin, J.M., Damo, J., Soulié, J.-C., Raboin, L.-M.,**
656 **Dusserre, J., Sow, A., Manneh, B., Shrestha, S. and Kretschmar, T.** (2017) Crop-
657 model assisted phenomics and genome-wide association study for climate adaptation
658 of indica rice. 2. Thermal stress and spikelet sterility. *Journal of Experimental Botany*,
659 **68**, 4389-4406.
- 660 **Dreni, L., Jacchia, S., Fornara, F., Fornari, M., Ouwerkerk, P.B.F., An, G., Colombo, L.**
661 **and Kater, M.M.** (2007) The D-lineage MADS-box gene OsMADS13 controls ovule
662 identity in rice. *The Plant Journal*, **52**, 690-699.
- 663 **Dreni, L., Pilatone, A., Yun, D., Erreni, S., Pajoro, A., Caporali, E., Zhang, D. and Kater,**
664 **M.M.** (2011) Functional Analysis of All AGAMOUS Subfamily Members in Rice
665 Reveals Their Roles in Reproductive Organ Identity Determination and Meristem
666 Determinacy. *The Plant Cell*, **23**, 2850-2863.

- 667 **Droc, G., Ruiz, M., Larmande, P., Pereira, A., Piffanelli, P., Morel, J.B., Dievart, A.,**
668 **Courtois, B., Guiderdoni, E. and Périn, C.** (2006) OryGenesDB: a database for rice
669 reverse genetics. *Nucleic Acids Research*, **34**, D736-D740.
- 670 **Edwards, J.W. and Coruzzi, G.M.** (1990) Cell-specific gene expression in plants. *Annual*
671 *Review of Genetics*, **24**, 275-303.
- 672 **Fu, Z., Yu, J., Cheng, X., Zong, X., Xu, J., Chen, M., Li, Z., Zhang, D. and Liang, W.**
673 (2014) The rice basic helix-loop-helix transcription factor TDR INTERACTING
674 PROTEIN2 is a central switch in early anther development. *The Plant Cell*, **26**, 1512-
675 1524.
- 676 **Garris, A.J., Tai, T.H., Coburn, J., Kresovich, S. and McCouch, S.** (2005) Genetic
677 Structure and Diversity in *Oryza sativa* L. *Genetics*, **169**, 1631-1638.
- 678 **Golicz, A.A., Singh, M.B. and Bhalla, P.L.** (2018a) LncRNAs in plant and animal sexual
679 reproduction. *Trends in Plant Science*, **23**, 195-205.
- 680 **Golicz, A.A., Singh, M.B. and Bhalla, P.L.** (2018b) The long intergenic non-coding
681 (lincRNA) landscape of the soybean genome. *Plant Physiology*.
- 682 **Gómez, J.F., Talle, B. and Wilson, Z.A.** (2015) Anther and pollen development: a conserved
683 developmental pathway. *Journal of Integrative Plant Biology*, **57**, 876-891.
- 684 **Goodstein, D.M., Shu, S., Howson, R., Neupane, R., Hayes, R.D., Fazo, J., Mitros, T.,**
685 **Dirks, W., Hellsten, U., Putnam, N. and Rokhsar, D.S.** (2012) Phytozome: a
686 comparative platform for green plant genomics. *Nucleic Acids Research*, **40**, D1178-
687 D1186.
- 688 **Hu, Y., Liang, W., Yin, C., Yang, X., Ping, B., Li, A., Jia, R., Chen, M., Luo, Z., Cai, Q.,**
689 **Zhao, X., Zhang, D. and Yuan, Z.** (2015) Interactions of OsMADS1 with floral
690 homeotic genes in rice flower development. *Molecular Plant*, **8**, 1366-1384.
- 691 **Huang, X., Wei, X., Sang, T., Zhao, Q., Feng, Q., Zhao, Y., Li, C., Zhu, C., Lu, T., Zhang,**
692 **Z., Li, M., Fan, D., Guo, Y., Wang, A., Wang, L., Deng, L., Li, W., Lu, Y., Weng,**
693 **Q., Liu, K., Huang, T., Zhou, T., Jing, Y., Li, W., Lin, Z., Buckler, E.S., Qian, Q.,**
694 **Zhang, Q.-F., Li, J. and Han, B.** (2010) Genome-wide association studies of 14
695 agronomic traits in rice landraces. *Nat Genet*, **42**, 961-967.
- 696 **Jagadish, S.V.K., Craufurd, P.Q. and Wheeler, T.R.** (2007) High temperature stress and
697 spikelet fertility in rice (*Oryza sativa* L.). *Journal of Experimental Botany*, **58**, 1627-
698 1635.
- 699 **Jagadish, S.V.K., Muthurajan, R., Oane, R., Wheeler, T.R., Heuer, S., Bennett, J. and**
700 **Craufurd, P.Q.** (2010) Physiological and proteomic approaches to address heat
701 tolerance during anthesis in rice (*Oryza sativa* L.). *Journal of Experimental Botany*, **61**,
702 143-156.
- 703 **Jarillo, J.A. and Piñeiro, M.** (2011) Timing is everything in plant development. The central
704 role of floral repressors. *Plant Science*, **181**, 364-378.
- 705 **Johnson, W.E., Li, C. and Rabinovic, A.** (2007) Adjusting batch effects in microarray
706 expression data using empirical Bayes methods. *Biostatistics*, **8**, 118-127.
- 707 **Kang, Y.-J., Yang, D.-C., Kong, L., Hou, M., Meng, Y.-Q., Wei, L. and Gao, G.** (2017)
708 CPC2: a fast and accurate coding potential calculator based on sequence intrinsic
709 features. *Nucleic Acids Research*, **45**, W12-W16.
- 710 **Kim, D., Langmead, B. and Salzberg, S.L.** (2015) HISAT: a fast spliced aligner with low
711 memory requirements. *Nat Meth*, **12**, 357-360.
- 712 **Kun, W., Xiaoju, P., Yanxiao, J., Yang, P., Yingguo, Z. and Li, S.** (2013) Gene, protein,
713 and network of male sterility in rice. *Frontiers in Plant Science*, **4**, 92.
- 714 **Kurczab, R. and Bojarski, A.J.** (2017) The influence of the negative-positive ratio and
715 screening database size on the performance of machine learning-based virtual
716 screening. *PLoS ONE*, **12**, e0175410.

- 717 **Kurczab, R., Smusz, S. and Bojarski, A.J.** (2014) The influence of negative training set size
718 on machine learning-based virtual screening. *Journal of Cheminformatics*, **6**, 32-32.
- 719 **Li, G., Jain, R., Chern, M., Pham, N.T., Martin, J.A., Wei, T., Schackwitz, W.S., Lipzen,**
720 **A.M., Duong, P.Q., Jones, K.C., Jiang, L., Ruan, D., Bauer, D., Peng, Y., Barry,**
721 **K.W., Schmutz, J. and Ronald, P.C.** (2017) The sequences of 1,504 mutants in the
722 model rice variety Kitaake facilitate rapid functional genomic studies. *The Plant Cell*.
- 723 **Liao, Y., Smyth, G.K. and Shi, W.** (2014) featureCounts: an efficient general purpose
724 program for assigning sequence reads to genomic features. *Bioinformatics*, **30**, 923-
725 930.
- 726 **Libbrecht, M.W. and Noble, W.S.** (2015) Machine learning applications in genetics and
727 genomics. *Nature Reviews Genetics*, **16**, 321.
- 728 **Lin, H., Yu, J., Pearce, S.P., Zhang, D. and Wilson, Z.A.** (2017) RiceAntherNet: a gene co-
729 expression network for identifying anther and pollen development genes. *The Plant*
730 *Journal*, **92**, 1076-1091.
- 731 **Linkov, I., Massey, O., Keisler, J., Rusyn, I. and Hartung, T.** (2015) From "Weight of
732 Evidence" to quantitative data integration using Multicriteria Decision Analysis and
733 Bayesian Methods. *ALTEX*, **32**, 3-8.
- 734 **Meyer, D.** (2017) Misc Functions of the Department of Statistics (e1071), TU Wien.
- 735 **Miyao, A., Iwasaki, Y., Kitano, H., Itoh, J.-I., Maekawa, M., Murata, K., Yatou, O.,**
736 **Nagato, Y. and Hirochika, H.** (2007) A large-scale collection of phenotypic data
737 describing an insertional mutant population to facilitate functional analysis of rice
738 genes. *Plant Molecular Biology*, **63**, 625-635.
- 739 **Miyao, A., Tanaka, K., Murata, K., Sawaki, H., Takeda, S., Abe, K., Shinozuka, Y.,**
740 **Onosato, K. and Hirochika, H.** (2003) Target site specificity of the Tos17
741 retrotransposon shows a preference for insertion within genes and against insertion in
742 retrotransposon-rich regions of the genome. *The Plant Cell*, **15**, 1771-1780.
- 743 **Moyroud, E. and Glover, B.J.** (2017) The Evolution of Diverse Floral Morphologies. *Current*
744 *Biology*, **27**, R941-R951.
- 745 **Mu, H., Ke, J., Liu, W., Zhuang, C. and Yip, W.** (2009) UDP-glucose pyrophosphorylase2
746 (OsUgp2), a pollen-preferential gene in rice, plays a critical role in starch accumulation
747 during pollen maturation. *Chinese Science Bulletin*, **54**, 234.
- 748 **Niu, N., Liang, W., Yang, X., Jin, W., Wilson, Z.A., Hu, J. and Zhang, D.** (2013) EAT1
749 promotes tapetal cell death by regulating aspartic proteases during male reproductive
750 development in rice. *Nat Commun*, **4**, 1445.
- 751 **Nonomura, K.-I., Miyoshi, K., Eiguchi, M., Suzuki, T., Miyao, A., Hirochika, H. and**
752 **Kurata, N.** (2003) The MSP1 gene is necessary to restrict the number of cells entering
753 into male and female sporogenesis and to initiate anther wall formation in rice. *The*
754 *Plant Cell*, **15**, 1728-1740.
- 755 **O'Leary, N.A., Wright, M.W., Brister, J.R., Ciufu, S., Haddad, D., McVeigh, R., Rajput,**
756 **B., Robbertse, B., Smith-White, B., Ako-Adjei, D., Astashyn, A., Badretdin, A.,**
757 **Bao, Y., Blinkova, O., Brover, V., Chetvernin, V., Choi, J., Cox, E., Ermolaeva,**
758 **O., Farrell, C.M., Goldfarb, T., Gupta, T., Haft, D., Hatcher, E., Hlavina, W.,**
759 **Joardar, V.S., Kodali, V.K., Li, W., Maglott, D., Masterson, P., McGarvey, K.M.,**
760 **Murphy, M.R., O'Neill, K., Pujar, S., Rangwala, S.H., Rausch, D., Riddick, L.D.,**
761 **Schoch, C., Shkeda, A., Storz, S.S., Sun, H., Thibaud-Nissen, F., Tolstoy, I., Tully,**
762 **R.E., Vatsan, A.R., Wallin, C., Webb, D., Wu, W., Landrum, M.J., Kimchi, A.,**
763 **Tatusova, T., DiCuccio, M., Kitts, P., Murphy, T.D. and Pruitt, K.D.** (2016)
764 Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion,
765 and functional annotation. *Nucleic Acids Research*, **44**, D733-D745.

- 766 **Obayashi, T., Hayashi, S., Saeki, M., Ohta, H. and Kinoshita, K.** (2009) ATTED-II
767 provides coexpressed gene networks for Arabidopsis. *Nucleic Acids Research*, **37**,
768 D987-D991.
- 769 **Ono, A., Yamaguchi, K., Fukada-Tanaka, S., Terada, R., Mitsui, T. and Iida, S.** (2012) A
770 null mutation of ROS1a for DNA demethylation in rice is not transmittable to progeny.
771 *The Plant Journal*, **71**, 564-574.
- 772 **Palla, G., Derenyi, I., Farkas, I. and Vicsek, T.** (2005) Uncovering the overlapping
773 community structure of complex networks in nature and society. *Nature*, **435**, 814-818.
- 774 **Pan, Y., Li, Q., Wang, Z., Wang, Y., Ma, R., Zhu, L., He, G. and Chen, R.** (2014) Genes
775 associated with thermosensitive genic male sterility in rice identified by comparative
776 expression profiling. *BMC Genomics*, **15**, 1114.
- 777 **Pertea, M., Pertea, G.M., Antonescu, C.M., Chang, T.-C., Mendell, J.T. and Salzberg,
778 S.L.** (2015) StringTie enables improved reconstruction of a transcriptome from RNA-
779 seq reads. *Nat Biotech*, **33**, 290-295.
- 780 **Quinlan, A.R. and Hall, I.M.** (2010) BEDTools: a flexible suite of utilities for comparing
781 genomic features. *Bioinformatics*, **26**, 841-842.
- 782 **Revelle, W.** (2017) psych: procedures for personality and psychological research. Evanston,
783 Illinois, USA: Northwestern University.
- 784 **Rhee, S.Y. and Mutwil, M.** (2014) Towards revealing the functions of all genes in plants.
785 *Trends in Plant Science*, **19**, 212-221.
- 786 **Schurko, A.M. and Logsdon, J.M.** (2008) Using a meiosis detection toolkit to investigate
787 ancient asexual “scandals” and the evolution of sex. *BioEssays*, **30**, 579-589.
- 788 **Shi, J., Dong, A. and Shen, W.-H.** (2015a) Epigenetic regulation of rice flowering and
789 reproduction. *Frontiers in Plant Science*, **5**, 803.
- 790 **Shi, X., Sun, X., Zhang, Z., Feng, D., Zhang, Q., Han, L., Wu, J. and Lu, T.** (2015b)
791 GLUCAN SYNTHASE-LIKE 5 (GSL5) plays an essential role in male fertility by
792 regulating callose metabolism during microsporogenesis in rice. *Plant and Cell
793 Physiology*, **56**, 497-509.
- 794 **Sperschneider, J., Gardiner, D.M., Dodds, P.N., Tini, F., Covarelli, L., Singh, K.B.,
795 Manners, J.M. and Taylor, J.M.** (2016) EffectorP: predicting fungal effector proteins
796 from secretomes using machine learning. *New Phytologist*, **210**, 743-761.
- 797 **Tatarinova, T.V., Chekalin, E., Nikolsky, Y., Bruskin, S., Chebotarov, D., McNally, K.L.
798 and Alexandrov, N.** (2016) Nucleotide diversity analysis highlights functionally
799 important genomic regions. *Proceedings of the National Academy of Sciences*, **6**,
800 35730.
- 801 **Therneau, T., Atkinson, B. and Ripley, B.** (2017) Recursive Partitioning and Regression
802 Trees.
- 803 **Troyanskaya, O.G., Dolinski, K., Owen, A.B., Altman, R.B. and Botstein, D.** (2003) A
804 Bayesian framework for combining heterogeneous data sources for gene function
805 prediction (in *Saccharomyces cerevisiae*). *Proceedings of the National Academy of
806 Sciences*, **100**, 8348-8353.
- 807 **Wallace, S., Fleming, A., Wellman, C.H. and Beerling, D.J.** (2011) Evolutionary
808 development of the plant and spore wall. *AoB Plants*, **2011**, plr027.
- 809 **Wang, H., Niu, Q.-W., Wu, H.-W., Liu, J., Ye, J., Yu, N. and Chua, N.-H.** (2015a) Analysis
810 of non-coding transcriptome in rice and maize uncovers roles of conserved lncRNAs
811 associated with agriculture traits. *The Plant Journal*, **84**, 404-416.
- 812 **Wang, J., Qi, M., Liu, J. and Zhang, Y.** (2015b) CARMO: a comprehensive annotation
813 platform for functional exploration of rice multi-omics data. *The Plant Journal*, **83**,
814 359-374.

- 815 **Wang, M., Wu, H.-J., Fang, J., Chu, C. and Wang, X.-J.** (2017) A long noncoding RNA
816 involved in rice reproductive development by negatively regulating osa-miR160.
817 *Science Bulletin*, **62**, 470-475.
- 818 **Ware, D., Jaiswal, P., Ni, J., Pan, X., Chang, K., Clark, K., Teytelman, L., Schmidt, S.,**
819 **Zhao, W., Cartinhour, S., McCouch, S. and Stein, L.** (2002) Gramene: a resource
820 for comparative grass genomics. *Nucleic Acids Research*, **30**, 103-105.
- 821 **Wen, K., Yang, L., Xiong, T., Di, C., Ma, D., Wu, M., Xue, Z., Zhang, X., Long, L., Zhang,**
822 **W., Zhang, J., Bi, X., Dai, J., Zhang, Q., Lu, Z.J. and Gao, G.** (2016) Critical roles
823 of long noncoding RNAs in Drosophila spermatogenesis. *Genome Research*, **26**, 1233-
824 1244.
- 825 **Wollstein, A. and Stephan, W.** (2015) Inferring positive selection in humans from genomic
826 data. *Investigative Genetics*, **6**, 5.
- 827 **Yano, K., Yamamoto, E., Aya, K., Takeuchi, H., Lo, P.-c., Hu, L., Yamasaki, M., Yoshida,**
828 **S., Kitano, H., Hirano, K. and Matsuoka, M.** (2016) Genome-wide association study
829 using whole-genome sequencing rapidly identifies new genes influencing agronomic
830 traits in rice. *Nat Genet*, **48**, 927-934.
- 831 **Yao, W., Li, G., Yu, Y. and Ouyang, Y.** (2017) funRiceGenes dataset for comprehensive
832 understanding and application of rice functional genes. *GigaScience*, gix119-gix119.
- 833 **You, Q., Zhang, L., Yi, X., Zhang, K., Yao, D., Zhang, X., Wang, Q., Zhao, X., Ling, Y.,**
834 **Xu, W., Li, F. and Su, Z.** (2016) Co-expression network analyses identify functional
835 modules associated with development and stress response in *Gossypium arboreum*. *Nat*
836 *Reports*, **6**, 38436.
- 837 **Yun, D., Liang, W., Dreni, L., Yin, C., Zhou, Z., Kater, M.M. and Zhang, D.** (2013)
838 OsMADS16 genetically interacts with OsMADS3 and OsMADS58 in specifying floral
839 patterning in rice. *Molecular Plant*, **6**, 743-756.
- 840 **Zhang, Y.-C., Liao, J.-Y., Li, Z.-Y., Yu, Y., Zhang, J.-P., Li, Q.-F., Qu, L.-H., Shu, W.-S.**
841 **and Chen, Y.-Q.** (2014) Genome-wide screening and functional analysis identify a
842 large number of long noncoding RNAs involved in the sexual reproduction of rice.
843 *Genome Biology*, **15**, 512.
- 844 **Zhao, L., Lei, J., Huang, Y., Zhu, S., Chen, H., Huang, R., Peng, Z., Tu, Q., Shen, X. and**
845 **Yan, S.** (2016) Mapping quantitative trait loci for heat tolerance at anthesis in rice using
846 chromosomal segment substitution lines. *Breeding Science*, **66**, 358-366.
- 847 **Zhou, Z., Jiang, Y., Wang, Z., Gou, Z., Lyu, J., Li, W., Yu, Y., Shu, L., Zhao, Y., Ma, Y.,**
848 **Fang, C., Shen, Y., Liu, T., Li, C., Li, Q., Wu, M., Wang, M., Wu, Y., Dong, Y.,**
849 **Wan, W., Wang, X., Ding, Z., Gao, Y., Xiang, H., Zhu, B., Lee, S.-H., Wang, W.**
850 **and Tian, Z.** (2015) Resequencing 302 wild and cultivated accessions identifies genes
851 related to domestication and improvement in soybean. *Nat Biotech*, **33**, 408-414.

852

853

854 **Supporting Information**

855 **Fig. S1** Distribution of tissues/cell types with peak expression levels of putative protein coding
856 genes not found in MSU-RAP annotation

857 **Fig. S2** Heatmap representing expression of all coding and non-coding genes

858 **Fig. S3** Distribution of PI scores for coding and non-coding genes

859 **Fig. S4** Comparison of classifier performance for three different biological processes

860 **Fig. S5** Comparison of classifier performance for three different biological processes with
861 scrambled labels

862 **Fig. S6** Number of shared SexRep genes identified by Naïve Bayes Classifier while varying
863 alpha and zeta parameters of PI score

864 **Fig. S7** Overlap between results of five MCRiceRepGP runs using different negative training
865 sets for coding and lincRNA genes

866 **Table S1** Datasets used in the analysis

867 **Table S2** Summary of annotation statistics

868 **Table S3** Comparison between current annotation and existing lincRNA annotations

869 **Table S4** Classification of samples as reproductive or vegetative

870 **Table S5** Dictionary of GO phrases associated with sexual reproduction

871 **Table S6** Classification of phenotypes as vegetative or reproductive

872 **Table S7** Biological process GO enrichment of 200 top genes as ranked by PI score

873 **Table S8** Biological process GO enrichment of 500 randomly chosen genes from bottom 95%
874 as ranked by PI score

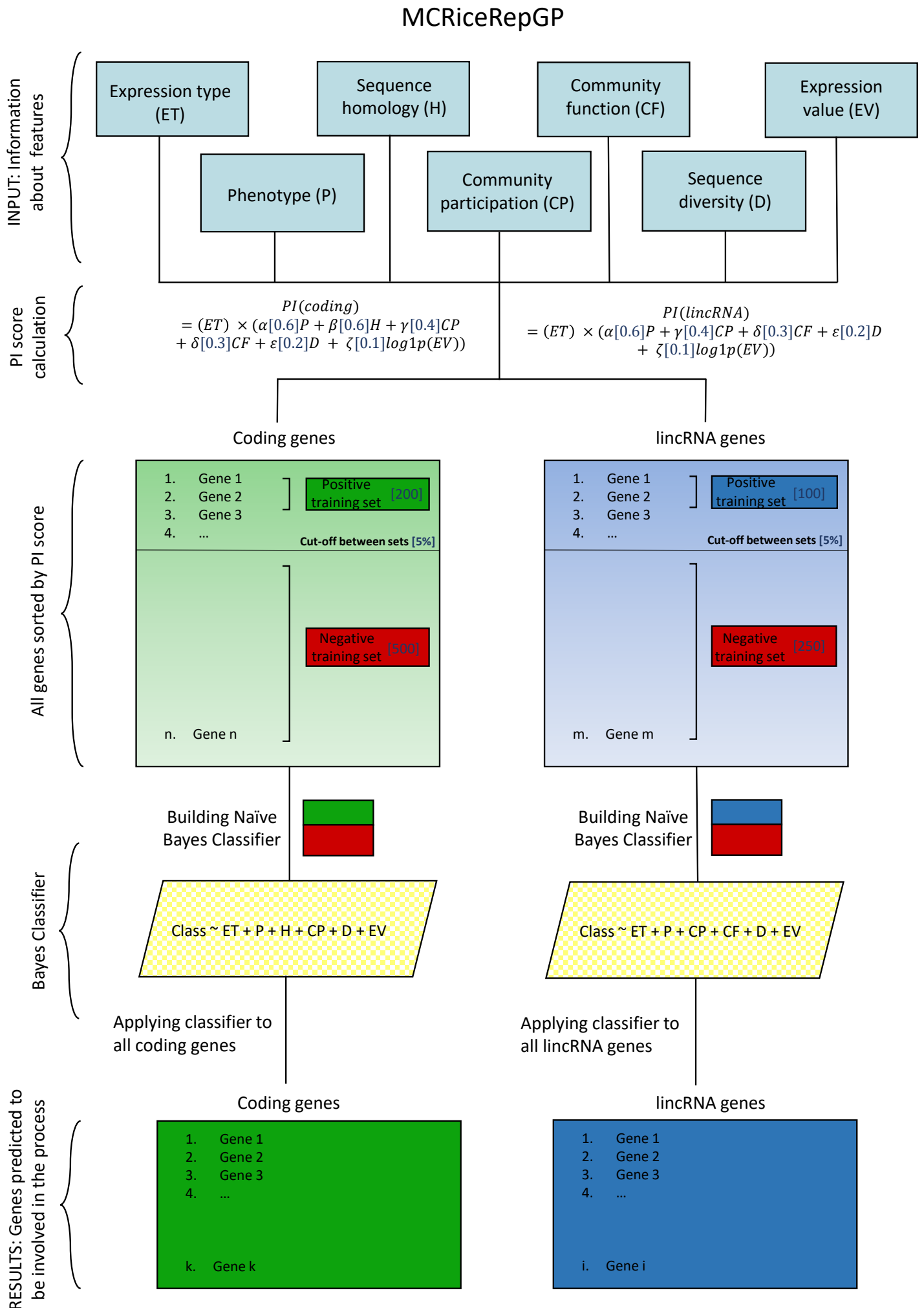
875 **Table S9** All SexRep genes identified

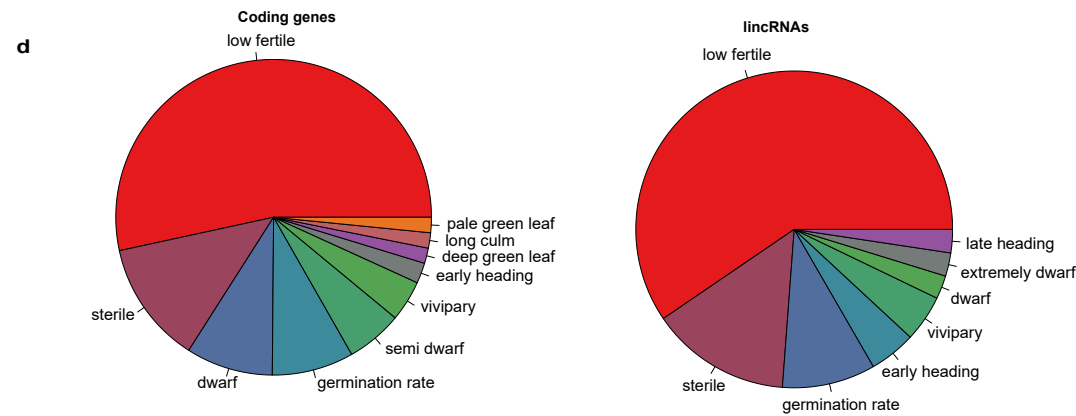
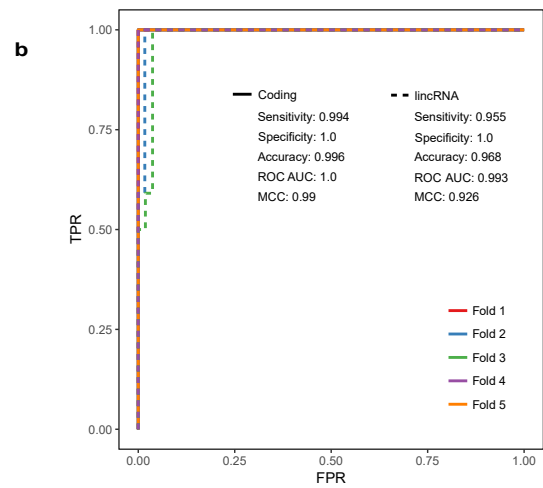
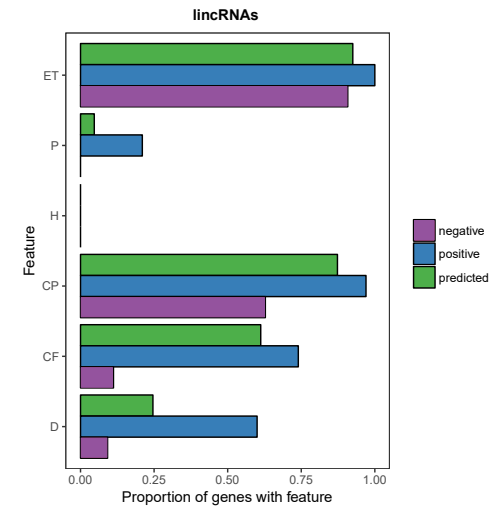
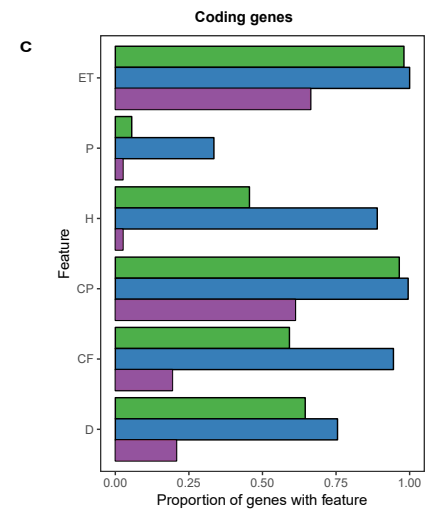
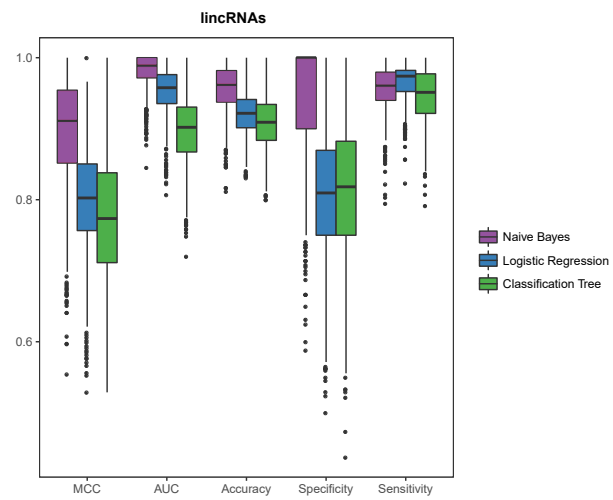
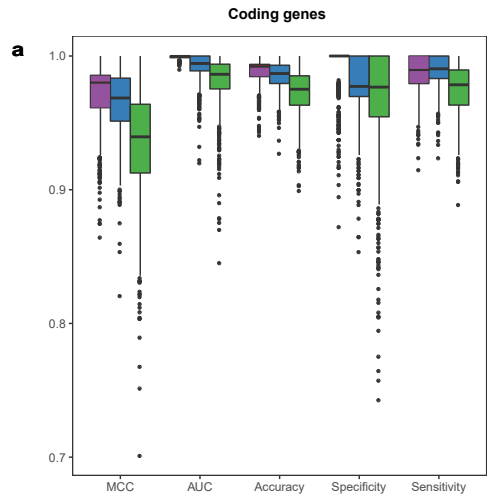
876 **Table S10** Biological processes GO enrichment of SexRep genes

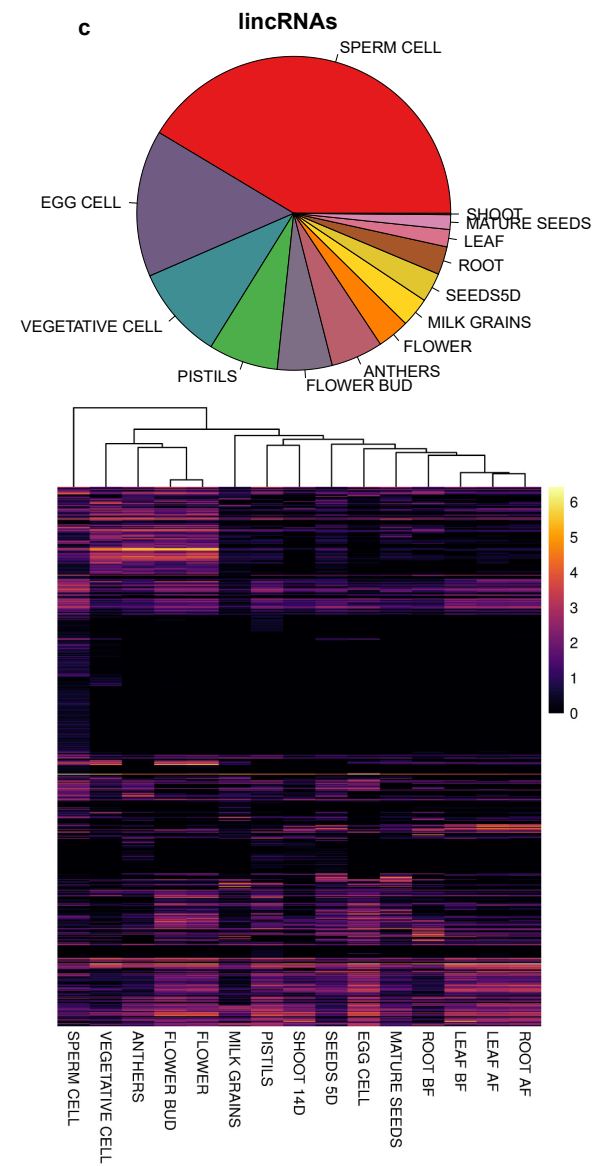
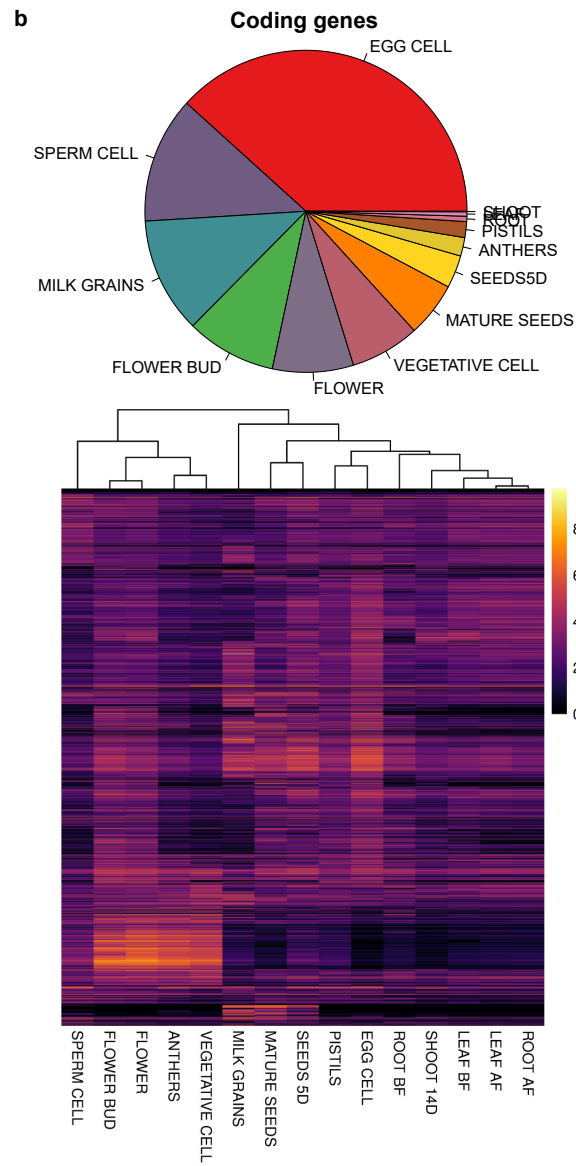
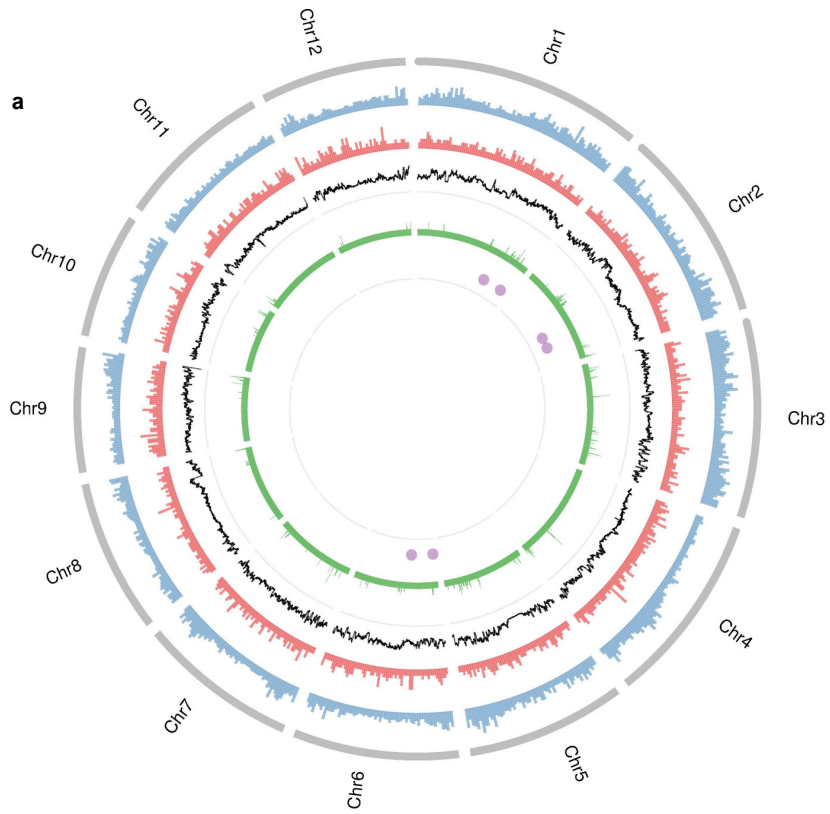
877 **Table S11** Molecular function GO enrichment of SexRep genes

878 **Table S12** Biological processes GO enrichment of SexRep genes falling within top 5% of most
879 differentiated genomic regions between indica and japonica lines

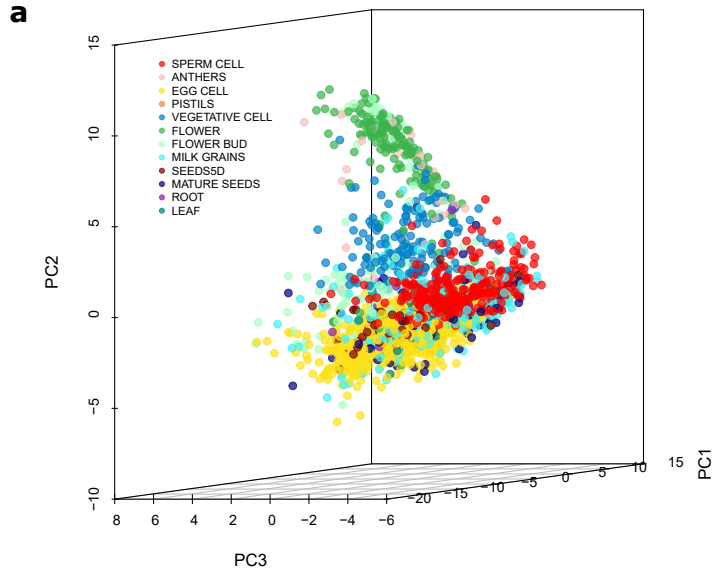
- 880 **Table S13** Anther dehiscence associated genes found in Supplementary table 12
- 881 **Table S14** SexRep gene overlapping sterility associated SNPs
- 882 **Method S1** Commands used for external software packages
- 883 **Method S2** Details of classifier implementation
- 884 **Note S1** Rice genome re-annotation
- 885 **Note S2** PI score parametrization for sexual reproduction
- 886 **Note S3** Comparison of classifiers used for identification of genes involved in three distinct
887 biological processes
- 888 **Note S4** Naïve Bayes Classifier Performance with varying Process Involvement score
889 parameters
- 890 **Note S5** Concordance between genes identified by classifiers built using different negative
891 training sets



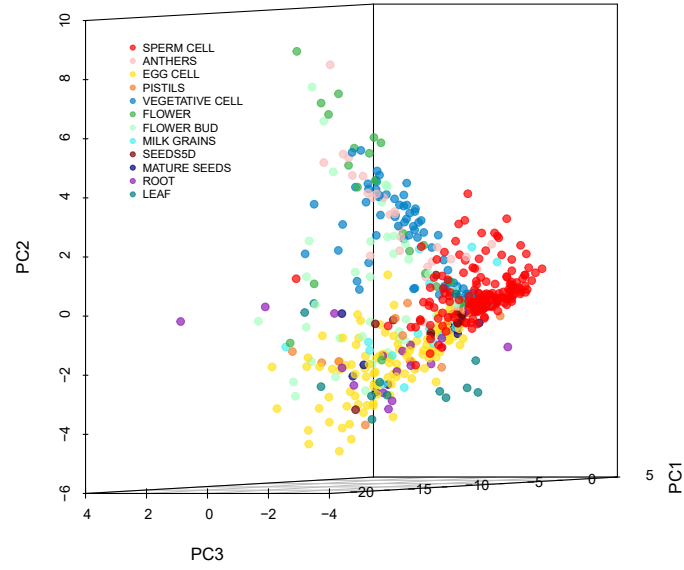




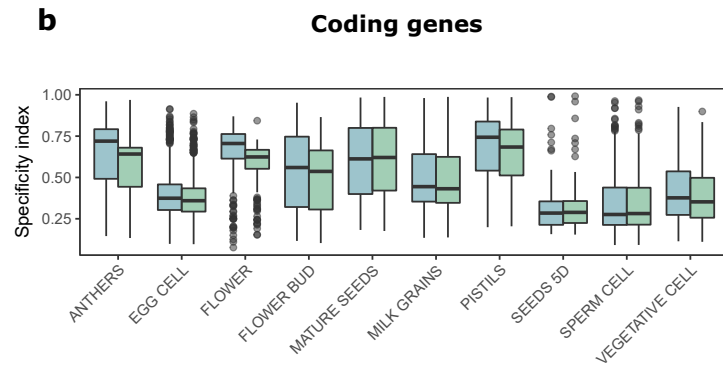
Coding genes



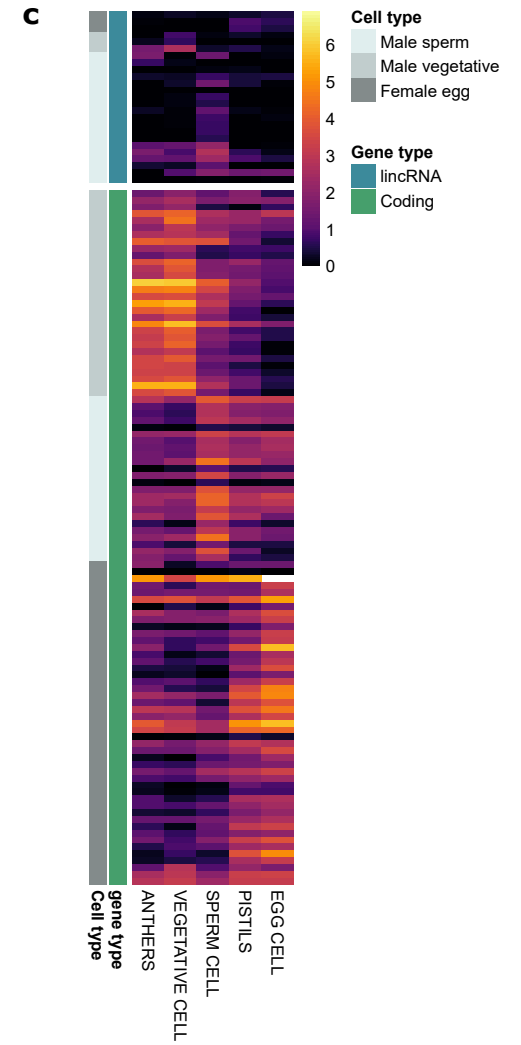
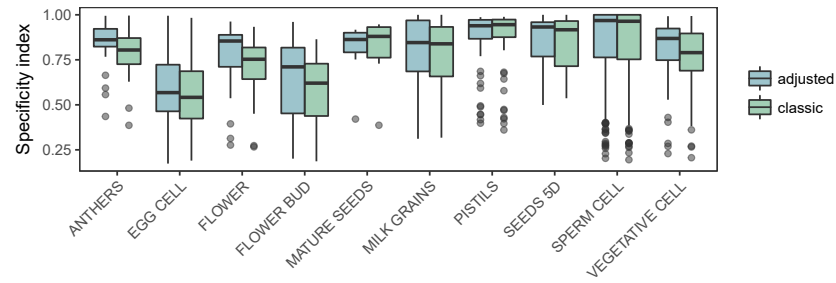
lincRNAs



Coding genes



lincRNAs



MCRiceRepGP

Expression type

Key words

PI score parameters

Gene type

Classifier features to include

Tissue type (ET):

Key words or phrases to include

Key words or phrases to exclude

Phenotype (P):

Homology (H):

Community participation (CP):

Community function (CF):

Sequence diversity (D):

Expression value (EV):

Size of positive training set:

Size of negative training set:

Percentage cut-off between positive and negative set

Coding or non-coding

Features to be used in the classifier

- ET
- P
- H
- CP
- CF
- D
- EV

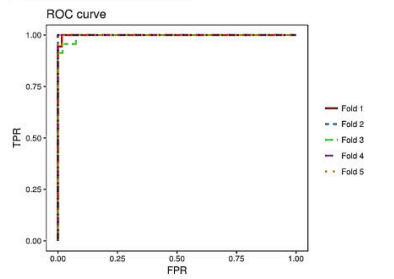
[Run MCRiceRepGP](#) [Download results](#)

Analysis Help

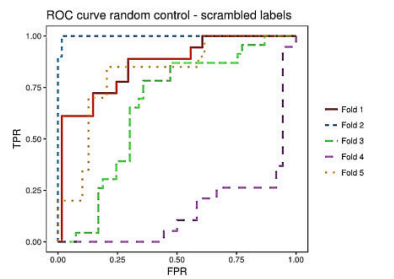
Session id: MCRiceRepGP24Feb2018031544

Sensitivity: 0.9967213(0.00733137)
 Specificity: 0.9614976(0.03786712)
 Accuracy: 0.9872044(0.01291534)
 AUC: 0.9989975(0.001776838)
 MCC: 0.9665674(0.03367194)

Classifier statistics



ROC curve for classifier



ROC curve for control classifier

	GeneID	ET	EV	P	H	CP	CF	D	PI	MSU-RAP
1	OSATST00018594	1	6.02	1	1	1	1	1	2.70	LOC_Os02g02560
2	OSATST00027308	1	4.20	1	1	1	1	1	2.52	LOC_Os03g24170
3	OSATST00041382	1	4.13	1	1	1	1	1	2.51	LOC_Os06g08380
4	OSATST00035108	1	4.34	1	1	1	1	0	2.33	LOC_Os04g52450
5	OSATST00003715	1	3.98	1	1	1	1	0	2.30	LOC_Os01g42060
6	OSATST00011253	1	3.59	1	1	1	1	0	2.26	LOC_Os11g04840
7	OSATST00027307	1	3.92	1	1	1	0	1	2.19	LOC_Os03g24160
8	OSATST00033548	1	6.92	1	0	1	1	1	2.19	LOC_Os04g38560

Table with results