

1 **A modified GC-specific MAKER gene annotation method**
2 **reveals improved and novel gene predictions of high and low**
3 **GC content in *Oryza sativa***

4
5 Megan J. Bowman^{1,2}, Jane A. Pulman^{1,3,4}, Tiffany L. Liu¹, and Kevin L. Childs^{1,3}

6
7 ¹Department of Plant Biology, Michigan State University, East Lansing, MI 48824

8 ²Van Andel Institute, Grand Rapids, MI 49506

9 ³Center for Genomics Enabled Plant Science, Michigan State University, East Lansing,
10 MI 48824

11 ⁴Centre for Genomics Research, University of Liverpool, Liverpool, UK, L69 7ZB

12
13 **Abstract**

14 Accurate structural annotation depends on well-trained gene prediction programs.
15 Training data for gene prediction programs are often chosen randomly from a subset of
16 high-quality genes that ideally represent the variation found within a genome. One aspect
17 of gene variation is GC content, which differs across species and is bimodal in grass
18 genomes. We find that gene prediction programs trained on genes with random GC
19 content do not completely predict all grass genes with extreme GC content. We present a
20 new GC-specific MAKER annotation protocol to predict new and improved gene models
21 and assess the biological significance of this method in *Oryza sativa*.

22

23

24 **Background**

25 Most widely used gene prediction programs depend on Hidden Markov Models
26 (HMMs) to predict gene structure within genomic sequence [1–3]. Typically, genes are
27 modeled within HMMs using a series of hidden states that represent generic gene
28 structure. The hidden states are filled with transition probabilities based on k-mer
29 sequences taken from the genes that are used to train the HMM. It is known that gene
30 GC content can affect gene predictions. Korf found that accuracy of predicting genes in
31 one species using a SNAP HMM that was trained for a second species was more
32 correlated with the GC content of the two species' genomes than with the phylogenetic
33 distance between the two species [1]. Additionally, in mammalian genomes, which
34 contain so-called isochores, gene GC content is correlated with the GC content of the
35 surrounding genome. The AUGUSTUS gene prediction program has a feature that trains
36 multiple HMMs that are each specialized for different narrow isochore-specific GC
37 ranges in order to improve gene predictions [4–6].

38 We perceived that two factors might limit the accuracy of gene prediction in grass
39 genomes. First, in many species including most plants, the GC content of genes has a
40 relatively narrow and unimodal distribution, but in the grasses (Poaceae), the GC content
41 of genes has a broad bimodal distribution (Figure 1A; [4,7–10]). The bimodal
42 distribution of GC-content in the grasses suggests that there exist two classes of genes
43 (high GC and low GC) that the gene prediction programs are attempting to learn.
44 While gene prediction programs perform well with grasses [11], we hypothesized that the
45 accuracy of grass gene predictions could be improved by accounting for the high and low

46 GC gene classes. Furthermore, as with any supervised machine learning technique, we
47 expect that it is difficult to predict grass genes at the tails of the natural GC distribution
48 and that some grass genes may not be predicted at all using existing protocols. Second,
49 grass gene GC content is not well correlated with the surrounding genomic regions
50 (Figure 1B; [10,12,13]), and therefore, grass genomes do not contain isochores. We also
51 predict that grass genome annotation will not benefit from analysis by the isochore-
52 sensitive AUGUSTUS protocol [6]. Therefore, it is probable that gene annotation in
53 grasses can be improved further.

54 MAKER is a commonly used structural annotation engine that has been used to
55 annotate numerous plant genome assemblies [11,14–18]. The MAKER gene annotation
56 pipeline makes it very easy to train and then predict gene models from commonly used
57 *ab initio* gene prediction programs, such as SNAP and AUGUSTUS [1,6,19]. We
58 developed a new GC-specific MAKER protocol that makes use of genes with high and
59 low GC content as training data in order to derive separate versions of the SNAP and
60 AUGUSTUS HMMs that are tuned to accurately predict high and low GC genes. Using
61 this new method, we improved regular MAKER gene predictions in *Oryza sativa* (rice)
62 relative to available transcript and protein evidence. Furthermore, we identified novel
63 genes with high and low GC content that had not been predicted by the standard MAKER
64 protocol. Comparisons to the AUGUSTUS isochore-based prediction method as well as
65 to the standard MAKER protocol showed that this GC-specific MAKER protocol shifts
66 the overall GC content of predicted gene models both higher and lower than the standard
67 MAKER protocol. This new GC-specific MAKER annotation method will be of interest

68 to anyone working on structural annotation of genomes with bimodal GC content but will
69 likely improve the annotation of any genome.

70

71 **Results**

72 **Reannotation of the *O. sativa* genome with MAKER using HMMs trained on high** 73 **and low GC content**

74 We thought that grass genes identified by gene prediction programs that are
75 trained on genes with specific GC content could both find different genes and produce
76 differing gene models at identical loci than prediction programs that are trained on genes
77 with random GC content. We tested this hypothesis by reannotating the genome of *O.*
78 *sativa*. In order to compare gene models within the *O. sativa* ssp. Nipponbare genome
79 (v7 assembly; [20]) based on the GC content of different HMM training sets, three
80 MAKER structural annotations were completed using a modified method. SNAP and
81 AUGUSTUS HMMS were trained either with training genes randomly picked without
82 regard for GC content, with training genes with low GC content or with training genes
83 with high GC content. The standard MAKER annotation using HMMs trained on
84 randomly selected training genes for SNAP and AUGUSTUS predicted 29,133 gene
85 models with transcript evidence and/or Pfam protein domains. The structural annotation
86 based on high GC HMMs produced 26,063 evidence supported gene models, and the
87 MAKER annotation based on low GC HMMs produced 26,559 evidence supported
88 models (Table 1). The average length of transcripts was very similar for the standard and
89 low GC structural annotations (Table 1). The average transcript length of the high GC

90 predictions was considerably shorter, a trend that has been previously discussed in
91 eukaryotic genomes [21].

92 The distribution of GC content of the gene predictions varied greatly (Fig. 2). The
93 standard MAKER annotation has a bimodal distribution of gene GC content with a major
94 peak at 49% and a minor peak at 68%. The GC distribution of the high GC annotation
95 has a unimodal distribution with a major peak at 68%. The low GC annotation has a
96 bimodal distribution with peaks at 47% and 67%. Notably, few low GC genes were
97 predicted by the high GC HMMs, and a lower percentage of high GC genes were
98 predicted by the low GC HMMs compared to the standard GC neutral MAKER
99 annotation.

100 The SNAP and AUGUSTUS HMMs created for the standard, high and low GC
101 MAKER structural annotations were also used together in a single MAKER run to
102 produce a six HMMs annotation (Fig. 3). For this annotation, up to six *ab initio*
103 predictions could be produced at a single locus, but when provided with multiple gene
104 predictions at a single locus, MAKER chooses the single best gene model at that locus.
105 The six HMMs annotation contained 29,942 evidence supported gene predictions (Table
106 1). The GC distribution for the six HMMs gene set was bimodal with a major peak at
107 48% and second peak at 68% (Fig. 3). In comparison to the MSU Rice Genome
108 Annotation Project (Release 7) annotation [20], 2,448 gene predictions were unique to the
109 MAKER six HMMs annotation of *O. sativa* while 7,004 gene models found in the MSU
110 annotation were missing from the SixHMMs annotation (Additional File 1).

111 To assess the impact of high and low GC specific HMMs on the structural
112 annotation of *O. sativa*, GC content and annotation edit distance (AED) scores were

113 plotted for each set of predicted gene models and visualized as heatmaps (Fig. 4). AED
114 scores are assigned by MAKER and can be used to assess gene prediction quality [22].
115 AED measures the concordance of a gene prediction relative to the transcript and protein
116 evidence that supports it. AED scores range between 0 and 1, where 0 indicates perfect
117 concordance between the gene prediction and evidence and 1 indicates that no transcript
118 or protein supports the prediction. Genes predicted by HMMs trained on specific GC
119 content caused a general shift in the GC distribution of predicted gene models for both
120 the high and low GC annotations, in comparison to the standard MAKER annotation
121 (Figs. 4A and 4B). In addition to this shift, standard MAKER gene predictions were
122 improved by high or low GC HMMs as determined by a decrease in AED scores between
123 overlapping gene predictions from the standard MAKER and high or low GC HMMs
124 annotations (Fig. 4E, 4F). The number of standard protocol gene models improved in the
125 six HMMs annotation was 3,740. The number or percent of genes with AED scores less
126 than 0.5 ($AED_{0.5}$) can be used for genome wide assessment of annotation quality. The
127 percentages of $AED_{0.5}$ genes were similar for all three annotations (Table 1). The high
128 percentage of well-supported gene predictions reflects the quality of transcriptome
129 evidence provided during the structural annotation process.

130

131 **Comparison of MAKER six HMMs method to alternative MAKER approaches**

132 The results of the MAKER six HMMs structural annotation were compared to
133 MAKER genome annotations where combinations of the SNAP and AUGUSTUS gene
134 prediction programs were used with alternative parameters. As AUGUSTUS can be run
135 so that it considers GC content of the genomic region (isochores) in which a gene

136 prediction is made, we also trained AUGUSTUS in its isochore-sensitive mode and used
137 it to make gene predictions within MAKER. Overall, the MAKER six HMMs annotation
138 produced more genes than any other annotation strategy tested here (Table 1, Additional
139 File 2: Table S1). MAKER run with only SNAP identified more evidence supported
140 genes than either AUGUSTUS alone trained with randomly chosen training data or the
141 isochore-specific AUGUSTUS protocol. Only a few hundred more genes were generated
142 by the isochore-specific AUGUSTUS annotation than by the randomly trained
143 AUGUSTUS HMM. Using randomly trained SNAP with either randomly trained or
144 isochore-specific AUGUSTUS produced similar numbers of gene predictions but more
145 than when MAKER is run with any of these programs alone. The number of AED_{0.5}
146 gene predictions follows a similar trend to the total number of gene predictions made by
147 any annotation protocol (Table 1; Additional File 2: Table S1; Additional File 3: Fig. S1).
148 However, as more genes are identified by a particular annotation method, the proportion
149 of AED_{0.5} genes decreases. The isochore-specific AUGUSTUS and the randomly trained
150 AUGUSTUS and SNAP gene predictions did not vary in overall GC content (Additional
151 File 3: Fig. S2).

152 For any machine learning protocol, different sets of training data can lead to
153 slightly different prediction results. To ensure that the results that we observed when we
154 trained SNAP and AUGUSTUS on high and low GC content training data sets were not
155 random, we repeated the standard MAKER annotations three times using independently
156 generated training data. The number of predicted gene models differed by less than 150
157 in the three randomly replicated standard MAKER annotations (Additional File 2: Table

158 S2), and the AED cumulative frequency plots were nearly identical (Additional File 3:
159 Fig. S3).

160

161 **Identification of novel high and low GC content genes**

162 In addition to the improved high and low GC structural annotations created with
163 the MAKER six HMMs annotation protocol, we discovered novel gene predictions
164 specific to the annotations from the high and low GC HMMs. The low GC annotation
165 contained 369 novel genes, while the high GC annotation contained 282 novel genes.
166 Interestingly, the novel genes predicted by the low GC HMMs did not always have a low
167 GC content, and some of the novel genes predicted by the high GC HMMs did not have
168 high GC content (Figs. 2, 4C, 4D). The locations of the novel high and low GC HMM
169 predictions were distributed across all twelve *O. sativa* chromosomes (Table 2;
170 Additional File 4). Of the novel high GC HMM predictions, 253 genes (90%) had some
171 level of protein or transcript evidence for the prediction, while 324 (88%) novel low GC
172 HMM predictions had protein or transcript support (Fig. 5). Overall, the AED scores
173 increased as GC content increased for the novel high GC HMM predictions and as GC
174 content decreased for the novel low GC HMM predictions (Figs. 4C and 4D). The
175 average length of the novel high GC genes was 1,439 bp, while the novel low GC genes
176 were on average 2,046 bp in length. We compared these novel gene predictions to the
177 MSU Rice Genome Annotation Project (MSU-RGAP) Release 7 gene set and found 112
178 of the low GC HMM predictions and 167 of the high GC HMM predictions were present
179 in that high-quality gene set [20].

180

181 **Orthology of high and low GC novel genes to genes from other grass species**

182 Using the total predictions generated through the MAKER six HMMs annotation,
183 additional support was given to the novel predictions made by the high GC and low GC
184 HMMs by first assessing sequence homology of the novel gene predictions to the NCBI
185 non-redundant protein database [23]. Of the 651 novel predictions, 387 had a significant
186 BLASTP hit (e-values less than $1e-10$) to NCBI's non-redundant protein database.
187 Second, the homology and orthology of these genes was evaluated relative to other
188 MAKER six HMM predictions and *Brachypodium distachyon*, *Sorghum bicolor* and *Zea*
189 *mays* using OrthoMCL [24–27] (Additional File 5). Of the novel high GC predictions, 51
190 genes were placed into orthogroups, with 19 as putative homologs only with other
191 MAKER six HMMs predictions, and 32 were orthologous to genes from at least one of
192 the other grass species. Interestingly, 23 novel high GC genes represented the only rice
193 prediction in their orthogroup, and 11 novel high GC genes were single copy orthologs
194 with the other grasses. Of the novel low GC predictions, 92 genes were placed into
195 orthogroups, with 34 as putative homologs only to other MAKER six HMMs gene
196 predictions, and 58 orthologous to the other grass species. Twelve novel low GC
197 predictions were the only rice representatives in their orthogroups.

198

199 **Translating Ribosome Affinity Purification (TRAP) sequencing provides additional** 200 **evidence for novel high and low GC gene predictions**

201 In an effort to demonstrate additional support for the new GC specific gene
202 models outside of the transcript data provided during the MAKER annotation process,
203 translating ribosome affinity purification sequence (TRAP-seq) reads were pseudoaligned

204 to the MAKER six HMMs annotation [28,29], and translome enrichment indices (TEI)
205 were calculated for each of the novel genes predicted by the high and low GC HMMs.
206 The TRAP-seq samples were isolated from callus, panicle and seedling tissues of an *O.*
207 *sativa* modified RPL18 transgene [30,31]. TRAP-seq reads were aligned to 200 (71%) of
208 the 282 novel high GC HMM predictions, and 236 (64%) of the 369 novel low GC HMM
209 predictions. This indicated that in addition to the transcript data already aligned to these
210 predictions during annotation, a majority of these novel predictions are in fact being
211 actively transcribed in various tissues from *O. sativa*. The TEI is the ratio of the
212 transcripts per million (TPM) of TRAP-seq to the TPM of mRNA-seq for a specific
213 locus. High TEIs may indicate preferential translation of a transcript, while very low
214 TEIs can be indicative of limited translation [30]. The calculated TEI of each of the novel
215 genes predicted by the high and low GC HMMs that had TRAP-seq pseudoalignments
216 indicates tissue specificity (Fig. 6).

217

218 **Discussion**

219 *Ab initio* gene prediction programs employ HMMs trained on gene sets that
220 should be representative of the variation in gene nucleotide content. We hypothesized
221 that in grass genomes, where genes have a wide variation in GC content and where that
222 distribution is bimodal (Fig. 1A), gene prediction programs trained on random sets of
223 training data would be overly generalized and that this could result in poorly predicted
224 gene models with high or low GC contents. To address this, we developed a GC-specific
225 MAKER gene annotation protocol that trains gene prediction programs SNAP and
226 AUGUSTUS using training data with both high and low GC content. The resulting high-

227 GC and low-GC SNAP and AUGUSTUS HMMs were used in addition to the regularly
228 trained SNAP and AUGUSTUS HMMs to predict genes within MAKER (Fig. 3).

229 We tested the six HMMs protocol by reannotating the *O. sativa* genome, and we
230 identified 29,942 genes with transcript, protein or Pfam protein domain support. As
231 expected, when MAKER predicted genes in the *O. sativa* genome using either the high-
232 GC or low-GC SNAP and AUGUSTUS HMMs, the GC content of the resulting gene
233 predictions were shifted higher or lower, respectively, compared to the GC content of
234 genes predicted by the standard MAKER protocol (Fig. 2). Furthermore, the GC content
235 distribution of genes predicted by the MAKER six HMMs protocol also showed a shift of
236 the bimodal peaks to higher and lower GC values (Fig. 2). Importantly, most gene
237 predictions made by the MAKER six HMMs annotation overlapped with loci predicted
238 by the standard MAKER protocol, but in 3,740 of these cases, the predictions made by
239 the MAKER six HMMs protocol were improved over the standard MAKER predictions
240 as shown by the better evidence support (i.e. lower AED scores) (Fig. 4E, 4F). This
241 indicates that the high and low GC HMMs were often able to improve upon gene
242 predictions made by the more generally trained gene prediction programs.

243 In addition to improving the annotation of many genes, we also identified novel
244 genes using this protocol. We found 651 genes that had been identified by high-GC or
245 low-GC SNAP or AUGUSTUS HMMs but that had not been predicted using the standard
246 MAKER pipeline. Of these newly identified genes, 372 were also not found in the most
247 recent MSU-RGAP Release 7 structural annotation [20]. The 279 novel genes predicted
248 by the high-GC or low-GC HMMs that were previously found in the MSU-RGAP
249 Release 7 were likely predicted by MSU-RGAP due to the use of Fgenesh for gene

250 identification, which may have its own biases related to GC content [20,32], or due to the
251 use of different transcript and protein evidence (Additional File 1). Additionally, the
252 MSU-RGAP annotation was improved by PASA, which improves *de novo* gene
253 predictions with transcript alignment evidence, and therefore, PASA is likely not biased
254 by GC content in the same way that HMM-based gene prediction programs can be
255 affected [33]. Furthermore, 90 of the novel genes identified by the high-GC and low-GC
256 HMMs were found to be orthologous to genes from other grass species or to other
257 MAKER six HMMs gene predictions within *O. sativa* (Additional File 5). Additional
258 support for the novel gene predictions comes from examining a TRAP sequencing data
259 set that indicates that 67% of these new predictions are being actively transcribed in three
260 different tissues from *O. sativa* [30] (Figure 6). Nonetheless, as with all computational
261 gene prediction methods, the novel gene models identified by the GC-specific MAKER
262 protocol should be further vetted through additional laboratory analysis.

263 There are 7,004 genes in the MSU-RGAP Release 7 data set that were not
264 predicted by the six HMMs annotation. Of these genes, 4,327 are characterized as
265 “expressed” meaning that they have may only have transcript support. An additional
266 1,365 MSU-RGAP genes missing from the six HMMs annotation are described as
267 “hypothetical”, which indicates that they have no transcript or protein support, but they
268 may contain a conserved protein domain (Additional File 1). The expressed and
269 hypothetical MSU-RGAP genes are the genes with the weakest support from that
270 annotation project. Some of the MSU-RGAP hypothetical genes may not pass the
271 stringent evidence test that was applied to the MAKER six HMMs gene predictions
272 which all had transcript or protein support or contained a Pfam domain. Additionally, the

273 MAKER six HMMs annotation only used transcript evidence derived from StringTie
274 assemblies of a small set of RNA-seq reads, but the MSU-RGAP annotation made use of
275 EST and FL-cDNA sequences that were not used in this report. This difference in
276 evidence will have an effect on the genes predicted by MAKER [11]. Finally, the
277 MAKER six HMMs annotation was filtered to remove any predictions that had homology
278 to known transposable elements (TE) and Pfam domains. The MSU-RGAP genes were
279 also filtered to flag any genes with matches to a library of TE sequences, but these two
280 methods were necessarily different and could have resulted in the removal of different
281 subsets of TE-related gene predictions. All of these reasons can help to explain why
282 7,004 MSU-RGAP genes are not present in our six HMMs MAKER annotation.

283 Interestingly, there may be additional unrecognized parameters that could be used
284 to improve gene prediction besides our strategy of training gene prediction HMMs in a
285 GC-specific fashion. In the six HMMs annotation, some low GC predictions were
286 generated by the high GC HMMs, and some high GC predictions came from the low GC
287 HMMs (Fig 4A, 4B). While these could be cases of identical gene models being created
288 by two or more HMMs at a particular locus with MAKER randomly retaining only one
289 prediction as the final model for the locus, we also observed novel low GC predictions
290 created by high GC HMMs as well as novel high GC predictions arising from low GC
291 HMMs (Figs. 3, 4C, 4D). This suggests that some unrecognized gene features besides
292 simple GC content were present in the high and low GC HMMs that allowed the
293 prediction of novel low and high GC genes, respectively.

294 It has been known that the GC content of genes used to train gene prediction
295 HMMs can affect the accuracy of gene predictions [1,6]. The AUGUSTUS gene finder

296 has an isochore-sensitive protocol that was developed in order to more accurately predict
297 mammalian genes. Despite the fact that isochores do not exist in plants (Fig. 1; [12,13]),
298 we used the isochore-sensitive AUGUSTUS protocol to predict genes in *O. sativa*, but
299 we did not see a substantial difference in the number or quality of predicted gene models
300 or a change in overall GC content distribution of those gene predictions (Additional File
301 3: Figs. S1, S2). This result was expected as gene GC content is not well correlated with
302 the GC content of the surrounding genomic region, and therefore, partitioning the training
303 data before training the gene prediction programs was found to be more effective at
304 improving gene annotations in *O. sativa*.

305 Given the importance of accurate gene prediction to downstream genomics
306 applications, the GC-specific MAKER protocol described here will be of use to those
307 working on the structural annotation of any species with a bimodal distribution of GC
308 content. MAKER is a powerful tool that enables research groups of any size to pursue
309 structural annotation of sequenced genomes and, with the addition of this protocol, will
310 aid in more accurate gene prediction.

311

312 **Conclusions**

313 In this paper we presented a new GC-specific MAKER annotation protocol that
314 was used to successfully identify new evidence supported gene models in *Oryza sativa*
315 with high and low GC content. This new method also improved 13% of gene models
316 produced by the standard MAKER protocol. Comparisons of this method to the standard
317 training protocols for the SNAP and AUGUSTUS *ab initio* gene prediction programs as
318 well as the isochore-sensitive AUGUSTUS gene prediction method showed that by

319 training gene prediction HMMs with data representing multiple ranges of GC content and
320 allowing MAKER to pick the best *ab initio* gene prediction generated by multiple gene
321 prediction HMMs, it is possible to create a final gene annotation set that includes large
322 numbers of both improved and novel gene predictions. The novel gene predictions are
323 supported by various forms of evidence including transcript and protein alignments and
324 membership in ortholog groups with genes from other grass species. Additionally,
325 TRAP-sequencing has shown that a majority of these new predictions are being actively
326 transcribed in *O. sativa*. MAKER is a widely used structural annotation program that
327 allows researchers to produce quality genome annotations. This new method will be an
328 important addition to those interested in the prediction of genes in regions of extreme GC
329 content in Poaceae genomes but will probably be generally applicable for species with
330 narrow, unimodal gene GC distributions as well.

331

332 **Methods**

333 **Processing, quality assessment and assembly of evidence**

334 Thirty-one paired end RNA-seq datasets for *O. sativa* grown from different stress
335 environments and tissues were downloaded from the National Center for Biotechnology
336 Information Sequence Read Archive (NCBI-SRA) (Additional File 2: Table S3) using
337 SRAToolkit v. 2.3.4-2 [34]. Raw read quality was assessed with FastQC v. 0.10.1 and
338 Illumina adapters were trimmed using Trimmomatic v. 0.32. Transcripts were assembled
339 using StringTie v. 1.3.0, and these transcript assemblies were subsequently used as EST
340 evidence for all MAKER runs. The SwissProt plant protein dataset was downloaded
341 (ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/knowledgebase/taxonomic_di

342 visions/uniprot_sprot_plants.dat.gz), and all *O. sativa* protein sequences were removed.
343 The remaining protein sequences not from *O. sativa* were used as protein evidence during
344 the MAKER annotation.

345 **MAKER standard *de novo* structural annotation of *O. sativa***

346 The MAKER-P (r1128) genome annotation pipeline was used to annotate the Os-
347 Nipponbare-Reference-IGRSP-1.0 v7 genome assembly. A custom repeat library was
348 created for *O. sativa* using a method described previously [18];
349 (http://weatherby.genetics.utah.edu/MAKER/wiki/index.php/Repeat_Library_Construction-Advanced), and the custom repeat library was used by RepeatMasker within the
350 MAKER pipeline to mask repetitive elements. Transcript assemblies and protein
351 sequences described above were used as evidence to aid gene predictions.
352

353 A complete description for running MAKER has been provided previously
354 [19,35] and that protocol provides details about ancillary scripts and example command
355 calls. An abbreviated description of the standard MAKER pipeline is given here, and
356 details about the extended GC-specific MAKER pipeline are given below. As MAKER
357 is run iteratively, repeat masking and evidence alignment was performed during an initial
358 MAKER run, and the resulting GFF3 file containing masked regions and protein and
359 transcript alignments was used during all subsequent MAKER runs. The initial MAKER
360 run generates data that aids in the training of the gene predictions programs SNAP
361 (version 2013-11-29) [1] and AUGUSTUS (version 2.6.1) [6] (Figure 3). During the
362 initial MAKER run, the parameter `est2genome` was used to cause MAKER to promote
363 transcript alignments to gene models. High-quality transcript-derived gene models (AED
364 ≤ 0.2) were used to train SNAP and AUGUSTUS. Instructions for training SNAP can

365 be found elsewhere [1,19,36]. We use a custom shell script, `train_augustus.sh`, which
366 trains the AUGUSTUS HMM in only a few hours for most species.

367

```
368 train_augustus.sh <path to working directory for training>  
369 <path to MAKER gff3 output from initial MAKER run> <species  
370 name for AUGUSTUS HMM directory> <path to single fasta file  
371 with all transcript assemblies>
```

372 The `train_augustus.sh` shell script prepares training and testing data sets and makes use of
373 the `autoAug.pl` training script from AUGUSTUS to create the appropriate HMM files.

374 This training script is relatively fast, as it only requires the transcript evidence to be
375 aligned to the genomic regions that contain training and testing gene models instead of
376 aligning those sequences to the entire genome. The working directory is used for writing
377 a number of intermediate files and directories during the AUGUSTUS training process.

378 All transcript sequences that were used as evidence during the initial MAKER run must
379 be placed into a single transcript fasta file and provided here as those sequences will be
380 used during the AUGUSTUS HMM training. The species name provided for the HMM

381 training will be used to name the directory that holds all of the files for the new HMM

382 and is also used to specify the AUGUSTUS HMM in the `maker_opts.ctl` file. It is

383 necessary to have write permissions in the `/config/species` directory within AUGUSTUS
384 installation directory in order for this script to work as that is where the AUGUSTUS

385 writes the species-specific HMM directory. On a shared compute system, it may be

386 necessary to make a local installation of AUGUSTUS and to then point MAKER to that

387 installation by updating the path in the `maker_exe.ctl` file. After training SNAP and

388 AUGUSTUS HMMs, MAKER was then run one last time using only the SNAP and

389 AUGUSTUS HMMs to predict genes. During the final MAKER run, the parameters
390 keep_preds was set to 1.

391 To identify the high-quality gene set, the MAKER accessory scripts gff_merge
392 and fasta_merge, which are included in the MAKER installation, were used to generate a
393 GFF3 file with all gene predictions and evidence data and the transcript and protein fasta
394 files for those predictions. Pfam domains were identified within the predicted proteins
395 using hmmscan[36].

396

```
397 hmmscan --domE 1e-5 -E 1e-5 --tblout <MAKER max predictions  
398 hmmscan output file> <path to Pfam-A.hmm> <path to predicted  
399 protein fasta file>
```

400

401 The annotation GFF3 file, the transcript and protein fasta files and the hmmscan results
402 file were used to generate a quality MAKER standard gene set.

403

```
404 generate_maker_standard_gene_list.pl --input_gff <output of  
405 gff3_merge> --pfam_results <hmmscan output> --pfam_cutoff  
406 1e-10 --output_file <path to MAKER standard gene list>
```

407

```
408 get_subset_of_fastas.pl -l <path to MAKER standard gene  
409 list> -f <fasta_merge output transcript/protein fasta> -o  
410 <path MAKER standard transcript/protein fasta>
```

411

```
412 create_maker_standard_gff.pl --input_gff <output of  
413 gff3_merge> --output_gff <path to MAKER standard GFF3> --  
414 maker_standard_gene_list <path to MAKER standard gene list>  
415
```

416 Despite our use of a custom repeat library that was used for masking repeat
417 elements in the genome, some TE-related genes remain unmasked, and we performed
418 additional analyses to identify and remove any TE-related predictions from our MAKER
419 standard gene set. Predicted proteins were compared to a database of Gypsy transposable
420 elements (3.1.b2) [37]. Predicted proteins were also aligned with blastp to a database of
421 transposases [38,39] (
422 http://weatherby.genetics.utah.edu/MAKER/wiki/index.php/Repeat_Library_Construction-Advanced)
423). A GFF3 file of TE-related genes was derived from the MSU-RGAP gene
424 annotation GFF3 file
425 ([http://rice.plantbiology.msu.edu/pub/data/Eukaryotic_Projects/o_sativa/annotation_dbs/
426 pseudomolecules/version_7.0/all.dir/all.gff3](http://rice.plantbiology.msu.edu/pub/data/Eukaryotic_Projects/o_sativa/annotation_dbs/pseudomolecules/version_7.0/all.dir/all.gff3)) and was compared to the MAKER standard
427 GFF3 file using gffcompare [40].

```
428  
429 hmmscan --tblout <Gypsy HMM analysis output> -E 1e-5 --domE  
430 1e-5 <path to gypsy_db_3.1b2.hmm> <path to maker standard  
431 proteins fasta>  
432  
433 blastp -db <Tpases020812 database> -query <path to MAKER  
434 standard protein fasta> -out <path to Tpases blast output> -  
435 evaluate 1e-10 -outfmt 6  
436
```

```
437 gffcompare -o <TE comparison output file> -r <MSU RGAP TE
438 GFF3> <MAKER standard GFF3>
439
440 The create_no_TE_genelist.py script use the data derived above, the Pfam hmmscan
441 results file and a list of TE-related Pfam domains (TE_Pfam_domains.txt; available on
442 Childs Lab GitHub repository) to create a list of MAKER standard genes with no TE-
443 related predictions.
444
445 create_no_TE_genelist.py --input_file_TEpfam
446 <TE_Pfam_domains.txt> --input_file_maxPfam <MAKER max
447 predictions hmmscan output file> --
448 input_file_geneList_toKeep <path to MAKER standard gene
449 list> --input_file_TEhmm <Gypsy HMM analysis output> --
450 input_file_TEblast <path to Tpsases blast output> --
451 input_file_Terefmap <TE comparison output refmap file> --
452 output_file <path to TE filtered MAKER standard gene list>
453
454 create_maker_standard_gff.pl --input_gff <MAKER standard
455 GFF3> --output_gff <TE filtered MAKER standard GFF3> --
456 maker_standard_gene_list <path to TE filtered MAKER standard
457 gene list>
458
459 get_subset_of_fastas.pl -l <path to TE filtered MAKER
460 standard gene list> -f <fasta_merge output
461 transcript/protein fasta> -o <TE filtered MAKER standard
462 transcript/protein fasta>
```

463

464 This high-quality gene set without TE-related genes was used for all analyses
465 presented in the Results section. In addition to this standard MAKER annotation, two
466 additional annotations were created using either the SNAP HMM alone or the
467 AUGUSTUS HMM alone, and high-quality gene sets without TE-related genes were
468 identified for each of these annotations, which were used for comparisons to the final
469 GC-specific six HMMs annotation described below.

470

471 **Training GC-specific HMMs with high-GC and low-GC gene sequences**

472 In order to train high and low GC-specific HMMs for SNAP and AUGUSTUS, it
473 was necessary to use training data that consisted of gene models with CDS GC content
474 within specific ranges. The transcript-based gene predictions from the initial MAKER
475 run (when the `est2genome` parameter was used) served as the starting point for GC-
476 specific HMM training (Fig. 3). After generating the GFF3 file describing the transcript-
477 based gene models, the genome FASTA file was processed by the Perl script
478 `MAKER_GC_cutoff_determination.pl`.

479

```
480 MAKER_GC_cutoff_determination.pl --fasta <full path to file  
481 with genome FASTA sequences> --gff <full path to MAKER  
482 created GFF3 of est2genome> --name <BASE_NAME for output  
483 files> --peak <peak determination window, odd integer,  
484 default is 5> --smooth <smoothing window, odd integer,  
485 default is 7>
```

486

487 The MAKER_GC_cutoff_determination.pl script helps to identify the GC values of the
488 peaks in a bimodal grass gene GC content distribution. The script pulls out the CDS
489 FASTA sequences for the transcript-based gene predictions from the GFF3 and calculates
490 the GC content for each gene prediction. The script assigns the gene GC values to
491 integer bins based on the --smooth parameter, which helps to smooth the calculation of
492 the distribution by using a moving window average and writes the results to a file. This
493 output file can be used in R to plot the distribution of gene GC content. A FASTA file of
494 the CDS sequences and a GC content file (showing nucleotide composition and GC
495 content of each prediction) are also produced. In addition, a text file is created with the
496 high and low peak values of the bimodal gene GC distribution that were used here as set
497 points in creating the high and low GC HMM training sets. These peaks are determined
498 by taking each GC bin and looking at a set number of bins on either side (set by --
499 peak). A peak is identified when the $((\text{peak} - 1) / 2)$ bins on each side of a GC bin
500 have lower calculated GC values than the middle GC bin. However, users may pick their
501 own high-GC and low-GC cutoff values, and the gene GC content distribution graph may
502 aid in picking those cutoff values. The MAKER_GC_training_set_create.py script relies
503 on two files produced from the MAKER_GC_cutoff_determination.pl script:

504 BASE_NAME_gc_content.txt and BASE_NAME_cutoff.txt. The
505 MAKER_GC_training_set_create.py script will create high-GC and low-GC GFF3 files
506 that can be used for training SNAP and AUGUSTUS.

507

```
508 MAKER_GC_training_set_create.py --input_file_gff <path to  
509 MAKER GFF file> --input_file_GC_content  
510 <BASE_NAME_gc_content.txt file> --input_file_GC_cutoff
```

```
511 <BASE_NAME_cutoff.txt file> --output_file_low <path to the  
512 low GC GFF file> --output_file_high <path to the high GC GFF  
513 file> --genome_fasta <path to the genome fasta file>
```

514

515 As detailed in Figure 3, this script is run after est2genome transcript alignment to create
516 new GC-specific HMM training sets. All subsequent steps should only use the high or
517 low GC output files.

518

519 **MAKER six HMMs annotation**

520 After the creation of high and low GC SNAP and AUGUSTUS HMMs, a final
521 MAKER run is performed using the standard, high and low GC HMMs at the same time.
522 When using multiple SNAP and AUGUSTUS HMMs for this six HMMs annotation,
523 predictions from the different HMMs can be identified by providing the path to a specific
524 HMM, a colon, and an HMM-specific identifier (see below). Providing a comma-
525 separated list to the `snaphmm` and `augustus_species` parameters allows the
526 designation of multiple HMMs. To create this new six HMMs structural annotation, the
527 following parameters are set in the MAKER `maker_opts.ctl` file:

528

```
529 #-----Re-annotation Using MAKER Derived GFF3
```

```
530 maker_gff = path to MAKER alignment GFF3
```

```
531 est_pass=1
```

```
532 protein_pass=1
```

```
533 rm_pass=1
```

534

```
535 #-----Gene Prediction
```

536 snaphmm= path to standard SNAP HMM:orig_snap, path to high
537 GC HMM:high_snap, path to low GC HMM:low_snap
538 AUGUSTUS_species= path to standard AUGUSTUS
539 directory:orig_aug, path to high GC AUGUSTUS
540 directory:high_aug, path to low GC AUGUSTUS
541 directory:low_aug
542 keep_preds=1

543

544 Once the six HMMs MAKER annotation is finished, a final high-quality MAKER gene
545 set composed of gene models with transcript, protein or Pfam domain support was
546 created using the same protocol that was used above for the standard annotation.

547

548 **Creation of SNAP and AUGUSTUS HMMs trained with transcripts of randomized** 549 **GC content**

550 To assess the impact of GC specific HMM training on the structural annotation of
551 *O. sativa*, three MAKER annotations were created using HMMs trained with transcripts
552 with randomized GC content from the standard annotation. The following Perl scripts
553 were used, which create the training GFF3 files based on a random seed instead of
554 percentage GC content. The random_dataset_generate.pl script takes as input the
555 MAKER standard transcript FASTA and outputs three transcript FASTAs to be used for
556 downstream GFF3 creation and HMM training.

557

```
558 random_dataset_generate.pl --transcript <name of transcript  
559 FASTA file > --random1 <name of output random file1> --  
560 random2 <name of output random file2> --random3 <name of
```


561 output random file3>

562

563 The seq_name.pl script was run for each of the three random output FASTA files, and

564 generates a list of MAKER standard transcript names from each transcript FASTA.

565

566 seq_name.pl

567 --fastafile <path to an output file from

568 random_dataset_generate.pl>

569 --output <name of text file with ID names for each FASTA

570 sequence>

571

572 Finally, the random_gff3_create.pl script requires as inputs the MAKER

573 standard GFF3 with the genome FASTA included and the gene IDs from each of the

574 random FASTAs, and the script generates the final randomized GFF3s that were used for

575 SNAP and AUGUSTUS HMM training.

576

577 random_gff3_create.pl

578 --align_gff <path to MAKER GFF3 with FASTA included>

579 --rand_1 <path to random 1 IDs>

580 --rand_2 <path to random 2 IDs>

581 --rand_3 <path to random 3 IDs>

582

583 The outputs of these steps are three GFF3 files containing the coordinates of randomly

584 selected gene predictions. Each of the GFF3 files created by random_gff3_create.pl was

585 then used for SNAP and AUGUSTUS training.

586

587 **Isochore-specific AUGUSTUS training in *O. sativa***

588 To compare the MAKER GC specific HMM training protocol to the isochore-
589 specific AUGUSTUS method, we trained AUGUSTUS in its isochore-sensitive mode as
590 detailed below [6]. After isochore-specific AUGUSTUS training, the resulting HMM was
591 used in a MAKER run with all other parameters as had been used for the standard
592 annotation to create a MAKER structural annotation based only on isochore-specific
593 AUGUSTUS gene predictions. An additional annotation was also created with the
594 isochore-specific AUGUSTUS HMM and the standard SNAP HMM for comparison to
595 the six HMMs annotation method.

596 After one round of traditional AUGUSTUS training [6] which creates the
597 augustus.gb.train and augustus.gb.test genbank formatted gene files, change the
598 gc_range_min value to 0.32, gc_range_max value to 0.73 and the decomp_num_steps
599 value to 7 in the parameters.cfg file in the newly created AUGUSTUS species HMM
600 directory. The following three commands will then complete the isochore-specific
601 training of AUGUSTUS.

602

```
603 [AUGUSTUS_installation_dir]/scripts/optimize_augustus.pl --
```

```
604 species=<species_name> augustus.gb.train
```

605

```
606 etraining --species=<species_name> augustus.gb.train
```

607

```
608 augustus --species=<species_name> augustus.gb.test
```

609

610 **Identification of novel high or low GC content gene predictions**

611 The BEDtools v 2.23.0 [41] intersect command was used to compare two GFF3
612 files containing MAKER gene coordinates to identify novel gene models that were
613 unique to the gene predictions created with high or low GC HMMs. Those predictions
614 that were only created by high GC HMMs but not standard or low GC HMMs were
615 considered novel high GC HMM predictions, while predictions created only by the low
616 GC HMMs were deemed novel low GC HMM predictions.

617

618 **Identification of orthologs of novel *O. sativa* gene predictions in other grass species**

619 Paralogs of novel high or low GC gene predictions in *O. sativa* and orthologs in
620 other grass species were identified using OrthoMCL (v1.4) [27] using default parameters
621 with the predicted proteins of the high or low GC unique genes. Predicted proteins of
622 *Brachypodium distachyon* (v 3.1;
623 https://phytozome.jgi.doe.gov/pz/portal.html#!info?alias=Org_Bdistachyon) *Zea mays*
624 (v5b+, Phytozome 11:
625 http://phytozome.jgi.doe.gov/pz/portal.html#!info?alias=Org_Zmays), *Oryza sativa* ssp.
626 Nipponbare (v7.0, <http://rice.plantbiology.msu.edu/>) and *Sorghum bicolor* (v3.1,
627 https://phytozome.jgi.doe.gov/pz/portal.html#!info?alias=Org_Sbicolor) were used for
628 comparison.

629

630 **Translating ribosome affinity purification (TRAP) sequencing analysis**

631 Paired end TRAP-seq and mRNA-seq reads were trimmed using Cutadapt v1.8.1
632 using default parameters. RNA quantification was conducted with the kallisto v 0.42.5

633 [42] pseudoalignment method using the six HMMs MAKER predicted transcripts with a
634 bootstrap value of 100. Transcripts per million (TPM) were calculated for both the
635 TRAP-seq and mRNA-seq reads for each tissue, and the translome enrichment index
636 (TEI) was calculated as the ratio of transcripts per million of TRAP-seq to mRNA-seq for
637 each high-quality six HMMs MAKER transcript. TRAP-seq and mRNA-seq data are
638 available in NCBI BioProject PRJNA298638.

639

640 **Figure Creation**

641 Figures were created in R (v. 3.1.1) [43] using the following packages: ggplot2 [44],
642 reshape2 [45] and NMF [46].

643

644 **Availability of Data**

645 Perl and python scripts for the GC-specific MAKER protocol are deposited in
646 Github: https://github.com/Childs-Lab/GC_specific_MAKER. Additional relevant data,
647 including FASTAs of MAKER predictions and GFF3 files can be found at the Dryad data
648 repository: (<http://www.datadryad.org>. Data can only be deposited after acceptance.
649 During the revision process the Dryad DOI will be added here.)

650

651 **Acknowledgments**

652 This research was supported by grant IOS-1126998 to KLC from the U.S. National
653 Science Foundation. We would like to thank John Hamilton for help creating the TE-
654 related Pfam domain list.

655

656 **Figure 1. Bimodal distribution and coding region GC content in *Oryza sativa*.** A)

657 Distribution of GC content of IRGSP v7 predicted coding regions. B) GC content of

658 IRGSP v7 predicted coding regions vs. genomic GC content 5Kb upstream and

659 downstream of predicted coding regions.

660

661 **Figure 2. Distribution of GC content of high-quality MAKER gene predictions.**

662 Distribution of GC content of various MAKER annotations created through the GC-

663 specific MAKER protocol. The high-quality standard and high GC MAKER genes retain

664 the bimodal distribution that is common to the Poaceae, while the high-quality low GC

665 MAKER genes and the novel high and low GC gene predictions have unimodal

666 distributions centered on GC content associated with the GC content of the HMM

667 training data.

668

669 **Figure 3. Six HMMs MAKER structural annotation method.** The center workflow

670 depicts the standard method for training hidden markov models for use in MAKER, while

671 the low GC (top) and high GC (bottom) training methods can be used after creating high

672 and low GC HMM training data sets. After separately training HMMs with the low and

673 high GC training data, all three SNAP HMMs and all three AUGUSTUS HMMs were

674 specified in the maker_opts.ctl file (see the Methods section), and MAKER was run to

675 create the six HMMs annotation, which incorporates gene predictions from the standard,

676 high and low GC MAKER runs.

677

678 **Figure 4. Heatmap visualization of annotated edit distance (AED) and GC content**
679 **of MAKER predicted gene models.** A) MAKER genes predicted using the high GC
680 HMMs. B) MAKER genes predicted using the low GC HMMs. C) Novel genes predicted
681 using the high GC HMMs. D) Novel genes predicted using the low GC HMMs. E) Gene
682 predictions from the high GC HMMs that improved gene predictions made by the
683 standard MAKER protocol. F) Gene predictions from the low GC HMMs that improved
684 gene predictions made by the standard MAKER protocol. G) Gene predictions from the
685 standard MAKER protocol. H) Gene predictions from the MAKER six HMMs
686 annotation.

687

688 **Figure 5. AED scores of high and low GC novel genes in *Oryza sativa*.** AED scores for
689 the novel A) high and B) low GC gene predictions generated through the MAKER
690 sixHMM annotation method.

691

692 **Figure 6. Translatome Enrichment Index (TEI) analysis of novel high and low GC**
693 **genes.** Heatmap of Translatome Enrichment Index (TEI) of the A) novel genes predicted
694 by low GC HMMs and B) novel genes predicted by high GC HMMs gene predictions,
695 which measures the ratio of TRAP-seq to mRNA seq for a specific transcript. Values are
696 scaled by row to a sum of one for visualization purposes.

697

698 Additional Files:

699

700 **Additional File 1: (.pdf) Venn diagram depicting the overlap between the rice GC-**
701 **specific sixHMM annotation and IGRSP v7 annotation.** Of the 7,004 genes that are
702 only present in the IGRSPv7 annotation, 1365 (19.5%) are designated as “hypothetical”,
703 while 4327 (61.8%) are designated as “expressed”.

704

705 **Additional File 2: (.xlsx) Transcript evidence used to reannotate *Oryza sativa* and**
706 **additional information from gene predictions generated from alternative MAKER**
707 **methods.** Table S1. Number of predictions, average transcript length and AED_{0.5} of gene
708 predictions generated by alternative MAKER approaches. Table S2. Number of
709 predictions, average transcript length and AED_{0.5} of the three randomly replicated
710 standard MAKER annotations. Table S3. RNA-seq transcript evidence used in the
711 reannotation of the *Oryza sativa* genome.

712

713 **Additional File 3: (.pdf) AED curves from various MAKER annotation methods.**
714 Figure S1. AED curves of MAKER annotations of *Oryza sativa* using various *ab initio*
715 prediction methods. Figure S2. Distribution of GC content of MAKER annotations of
716 *Oryza sativa* using various *ab initio* prediction methods. Figure S3. AED curves of
717 MAKER annotations of *Oryza sativa* using HMMs trained with randomized training data.

718

719 **Additional File 4: (.pdf) Distribution of GC content, MAKER six HMMs gene**
720 **predictions and novel genes predicted by the high and low GC HMMs in the *Oryza***
721 ***sativa* genome.** A) Genomic GC content in 300Mb bins. Warmer colors indicate higher
722 than average GC content while cooler colors indicate lower than average GC content. B)

723 Heatmap visualization of the density of MAKER six HMMs gene models. C) Genomic
724 location of novel genes predicted by the high GC HMMs. D) Genomic location of novel
725 genes predicted by the low GC HMMs.

726

727 **Additional File 5: (.txt) OrthoMCL orthogroups containing novel high and low GC**
728 **gene predictions.** OrthoMCL output listing the number of genes, taxa and gene names
729 for each orthogroup that contains at least one novel high or low GC prediction. Novel
730 genes are indicated by bold text.

731

732 REFERENCES

- 733 1. Korf I. Gene finding in novel genomes. BMC Bioinformatics. 2004;5:59.
- 734 2. Majoros WH, Pertea M, Salzberg SL. TigrScan and GlimmerHMM: Two open source
735 ab initio eukaryotic gene-finders. Bioinformatics. 2004;20:2878–287910.
- 736 3. Lukashin A V, Borodovsky M. GeneMark.hmm: new solutions for gene finding.
737 Nucleic Acids Res. 1998;26:1107–15.
- 738 4. Carels N, Bernardi G. Two classes of genes in plants. Genetics. 2000;154:1819–25.
- 739 5. Costantini M, Clay O, Auletta F, Bernardi G. An isochore map of human
740 chromosomes. Genome Res. 2006;16:536–41.
- 741 6. Stanke M, Waack S. Gene prediction with a hidden markov model and a new intron
742 submodel. Bioinformatics. 2003;19 Suppl 2:ii215-ii225.
- 743 7. Wong GK-S, Wang J, Tao L, Tan J, Zhang J, Passey D a, et al. Compositional
744 gradients in Gramineae genes. Genome Res. [Internet]. 2002;12:851–6. Available from:
745 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1383739&tool=pmcentrez&r>

- 746 endertype=abstract
- 747 8. Wang H-C, Singer GAC, Hickey DA. Mutational bias affects protein evolution in
748 flowering plants. *Mol. Biol. Evol.* [Internet]. 2004;21:90–6. Available from:
749 <http://mbe.oupjournals.org/cgi/doi/10.1093/molbev/msh003>
- 750 9. Romiguier J, Ranwez V, Douzery EJP, Galtier N. Contrasting GC-content dynamics
751 across 33 mammalian genomes: Relationship with life-history traits and chromosome
752 sizes. *Genome Res.* 2010;20:1001–9.
- 753 10. Clément Y, Fustier M-A, Nabholz B, Glémin S. The bimodal distribution of genic
754 GC content is ancestral to monocot species. *Genome Biol. Evol.* [Internet]. 2015;7:336–
755 48. Available from: <http://gbe.oxfordjournals.org/content/7/1/336.abstract>
- 756 11. Law M, Childs KL, Campbell MS, Stein JC, Olson AJ, Holt C, et al. Automated
757 update, revision, and quality control of the maize genome annotations using MAKER-P
758 improves the B73 RefGen_v3 gene models and identifies new genes. *Plant Physiol.*
759 [Internet]. 2015;167:25–39. Available from:
760 <http://www.plantphysiol.org/lookup/doi/10.1104/pp.114.245027>
- 761 12. Tatarinova T V, Alexandrov NN, Bouck JB, Feldmann KA. GC3 biology in corn,
762 rice, sorghum and other grasses. *BMC Genomics.* 2010;11:308.
- 763 13. Glémin S, Clément Y, David J, Ressayre A. GC content evolution in coding regions
764 of angiosperm genomes: a unifying hypothesis. *Trends Genet.* [Internet]. 2014;30:263–
765 70. Available from: <http://linkinghub.elsevier.com/retrieve/pii/S0168952514000808>
- 766 14. Holt C, Yandell M. MAKER2: an annotation pipeline and genome-database
767 management tool for second-generation genome projects. *BMC Bioinformatics* [Internet].
768 BioMed Central Ltd; 2011;12:491. Available from: <http://www.biomedcentral.com/1471->

- 769 2105/12/491
- 770 15. Kellner F, Kim J, Clavijo BJ, Hamilton JP, Childs KL, Vaillancourt B, et al. Genome-
771 guided investigation of plant natural product biosynthesis. *Plant J.* [Internet].
772 2015;82:680–92. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/25759247>
- 773 16. Neale DB, Wegrzyn JL, Stevens KA, Zimin A V, Puiu D, Crepeau MW, et al.
774 Decoding the massive genome of loblolly pine using haploid DNA and novel assembly
775 strategies. *Genome Biol.* [Internet]. 2014;15:R59. Available from:
776 <http://genomebiology.com/2014/15/3/R59>
- 777 17. Wegrzyn JL, Liechty JD, Stevens KA, Wu L-S, Loopstra CA, Vasquez-Gross HA, et
778 al. Unique features of the loblolly pine (*Pinus taeda* L.) megagenome revealed through
779 sequence annotation. *Genetics* [Internet]. 2014;196:891–909. Available from:
780 <http://www.genetics.org/cgi/doi/10.1534/genetics.113.159996>
- 781 18. Campbell MS, Law M, Holt C, Stein JC, Moghe GD, Hufnagel DE, et al. MAKER-P:
782 A tool kit for the rapid creation, management, and quality control of plant genome
783 annotations. *Plant Physiol.* [Internet]. 2014;164:513–24. Available from:
784 <http://www.plantphysiol.org/cgi/doi/10.1104/pp.113.230144>
- 785 19. Campbell MS, Holt C, Moore B, Yandell M. Genome annotation and curation using
786 MAKER and MAKER-P. *Curr. Protoc. Bioinforma.* 2014;48.
- 787 20. Kawahara Y, de la Bastide M, Hamilton JP, Kanamori H, McCombie WR, Ouyang S,
788 et al. Improvement of the *Oryza sativa* Nipponbare reference genome using next
789 generation sequence and optical map data. *Rice* [Internet]. 2013;6:4. Available from:
790 <http://www.thericejournal.com/content/6/1/4>
- 791 21. Zhu L, Zhang Y, Zhang W, Yang S, Chen J-Q, Tian D. Patterns of exon-intron

- 792 architecture variation of genes in eukaryotic genomes. *BMC Genomics* [Internet].
793 BioMed Central; 2009 [cited 2017 Feb 10];10:47. Available from:
794 <http://bmcgenomics.biomedcentral.com/articles/10.1186/1471-2164-10-47>
- 795 22. Eilbeck K, Moore B, Holt C, Yandell M. Quantitative measures for the management
796 and comparison of annotated genomes. *BMC Bioinformatics*. 2009;10:67.
- 797 23. Pruitt K, Brown G, Tatusova T, Maglott D. The Reference Sequence (RefSeq)
798 Database. National Center for Biotechnology Information (US); 2012;
- 799 24. International T, Initiative B. Genome sequencing and analysis of the model grass
800 *Brachypodium distachyon*. *Nature* [Internet]. 2010 [cited 2014 Feb 19];463:763–8.
801 Available from: <http://www.ncbi.nlm.nih.gov/pubmed/20148030>
- 802 25. Schnable PS, Ware D, Fulton RS, Stein JC, Wei F, Pasternak S, et al. The B73 maize
803 genome: complexity, diversity, and dynamics. *Science* (80-.). [Internet]. 2009;326:1112–
804 5. Available from: <http://www.sciencemag.org/cgi/doi/10.1126/science.1178534>
- 805 26. Paterson AH, Bowers JE, Bruggmann R, Dubchak I, Grimwood J, Gundlach H, et al.
806 The *Sorghum bicolor* genome and the diversification of grasses. *Nature* [Internet].
807 2009;457:551–6. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/19189423>
- 808 27. Li L, Stoeckert CJJ, Roos DS. OrthoMCL: Identification of ortholog groups for
809 eukaryotic genomes. *Genome Res*. [Internet]. 2003;13:2178–89. Available from:
810 <http://genome.cshlp.org/cgi/content/full/13/9/2178>
- 811 28. Mustroph A, Juntawong P, Bailey-Serres J. Isolation of plant polysomal mRNA by
812 differential centrifugation and ribosome immunopurification methods. In: Belostotsky
813 DA, editor. *Methods Mol. Biol.* [Internet]. 1st ed. Humana Press; 2009 [cited 2016 Jun
814 15]. p. 109–26. Available from: http://link.springer.com/10.1007/978-1-60327-563-7_6

- 815 29. Reynoso MA, Juntawong P, Lancia M, Blanco FA, Bailey-Serres J, Zanetti ME.
816 Translating ribosome affinity purification (TRAP) followed by RNA sequencing
817 technology (TRAP-SEQ) for quantitative assessment of plant translomes. In: Alonso
818 JM, Stepanova AN, editors. *Plant Funct. Genomics* [Internet]. Springer New York; 2015
819 [cited 2016 Jun 13]. p. 185–207. Available from: [http://link.springer.com/10.1007/978-1-](http://link.springer.com/10.1007/978-1-4939-2444-8_9)
820 [4939-2444-8_9](http://link.springer.com/10.1007/978-1-4939-2444-8_9)
- 821 30. Zhao D, Hamilton JP, Hardigan M, Yin D, He T, Vaillancourt B, et al. Analysis of
822 ribosome-associated mRNAs in rice reveals the importance of transcript size and GC
823 content in translation. *G3 (Bethesda)*. 2016;7.
- 824 31. Zanetti ME, Chang I-F, Gong F, Galbraith DW, Bailey-Serres J. Immunopurification
825 of polyribosomal complexes of Arabidopsis for global analysis of gene expression. *Plant*
826 *Physiol.* [Internet]. 2005 [cited 2016 Jun 13];138:624–35. Available from:
827 <http://www.plantphysiol.org/cgi/doi/10.1104/pp.105.059477>
- 828 32. Salamov AA, Solovyev V V. Ab initio Gene Finding in Drosophila Genomic DNA.
829 *Genome Res.* [Internet]. Cold Spring Harbor Laboratory Press; 2000 [cited 2016 Aug
830 2];10:516–22. Available from: <http://www.genome.org/cgi/doi/10.1101/gr.10.4.516>
- 831 33. Haas BJ, Delcher AL, Mount SM, Wortman JR, Jr RKS, Hannick LI, et al. Improving
832 the Arabidopsis genome annotation using maximal transcript alignment assemblies.
833 *Nucleic Acids Res.* [Internet]. Oxford University Press; 2003 [cited 2016 Aug
834 2];31:5654–66. Available from:
835 <http://nar.oxfordjournals.org/lookup/doi/10.1093/nar/gkg770>
- 836 34. Kodama Y, Shumway M, Leinonen R. The sequence read archive: explosive growth
837 of sequencing data. *Nucleic Acids Res.* [Internet]. 2012;40:D54–6. Available from:

- 838 <http://nar.oxfordjournals.org/lookup/doi/10.1093/nar/gkr854>
- 839 35. Cantarel BL, Korf I, Robb SMC, Parra G, Ross E, Moore B, et al. MAKER: An easy-
840 to-use annotation pipeline designed for emerging model organism genomes. *Genome*
841 *Res.* [Internet]. 2008;18:188–96. Available from:
842 <http://genome.cshlp.org/content/18/1/188.short>
- 843 36. Quevillon E, Silventoinen V, Pillai S, Harte N, Mulder N, Apweiler R, et al.
844 InterProScan: protein domains identifier. *Nucleic Acids Res.* [Internet]. 2005 [cited 2013
845 Mar 3];33:W116-20. Available from:
846 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1160203&tool=pmcentrez&r](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1160203&tool=pmcentrez&rendertype=abstract)
847 [endertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1160203&tool=pmcentrez&rendertype=abstract)
- 848 37. Llorens C, Munoz-Pomer A, Futami R, Moya A. The GyDB Collection of Viral and
849 Mobile Genetic Element Models. *Biotechvana* [Internet]. 2008 [cited 2017 Mar 6];
850 Available from:
851 http://biotechvana.uv.es/bioinformatics/article_files/31/pdf/gydb_collection2.pdf
- 852 38. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped
853 BLAST and PSI-BLAST: A new generation of protein database search programs. *Nucleic*
854 *Acids Res.* 1997;25:3389–402.
- 855 39. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al.
856 BLAST+: architecture and applications. *BMC Bioinformatics* [Internet]. BioMed Central;
857 2009 [cited 2016 May 2];10:421. Available from: [http://www.biomedcentral.com/1471-](http://www.biomedcentral.com/1471-2105/10/421)
858 [2105/10/421](http://www.biomedcentral.com/1471-2105/10/421)
- 859 40. Pertea M, Pertea GM, Antonescu CM, Chang T-C, Mendell JT, Salzberg SL.
860 StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat.*

861 Biotechnol. [Internet]. Nature Research; 2015 [cited 2017 Feb 10];33:290–5. Available
862 from: <http://www.nature.com/doi/10.1038/nbt.3122>

863 41. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic
864 features. Bioinformatics [Internet]. 2010;26:841–2. Available from:
865 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2832824&tool=pmcentrez&r](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2832824&tool=pmcentrez&rendertype=abstract)
866 [endertype=abstract\http://bioinformatics.oxfordjournals.org/content/26/6/841.short](http://bioinformatics.oxfordjournals.org/content/26/6/841.short)

867 42. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq
868 quantification. Nat. Biotechnol. [Internet]. Nature Publishing Group, a division of
869 Macmillan Publishers Limited. All Rights Reserved.; 2016;34:525–7. Available from:
870 <http://dx.doi.org/10.1038/nbt.3519>

871 43. R Core Team. R: A Language and Environment for Statistical Computing [Internet].
872 Vienna, Austria; 2014. Available from: <http://www.r-project.org/>

873 44. Wickham H. ggplot2: elegant graphics for data analysis [Internet]. Springer New
874 York; 2009. Available from: <http://had.co.nz/ggplot2/book>

875 45. Wickham H. Reshaping Data with the {reshape} Package. J. Stat. Softw. [Internet].
876 2007;21:1–20. Available from: <http://www.jstatsoft.org/v21/i12/>

877 46. Gaujoux R, Seoighe C, Paatero P, Tapper U, Lee D, Seung H, et al. A flexible R
878 package for nonnegative matrix factorization. BMC Bioinformatics [Internet]. BioMed
879 Central; 2010 [cited 2016 Jun 10];11:367. Available from:
880 <http://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-11-367>

881

882

883

884 **Table 1. Numbers of high quality rice genes predicted by different MAKER**
885 **protocols**

886

Annotation	Number of Predictions	Average Transcript Length	AED _{0.5}	Percentage (%)
Standard Protocol	29133	1920	26809	92.0
High GC	26063	1439	22600	86.7
Low GC	26559	2046	25091	94.5
Six HMMs	29942	1947	27395	91.5

887

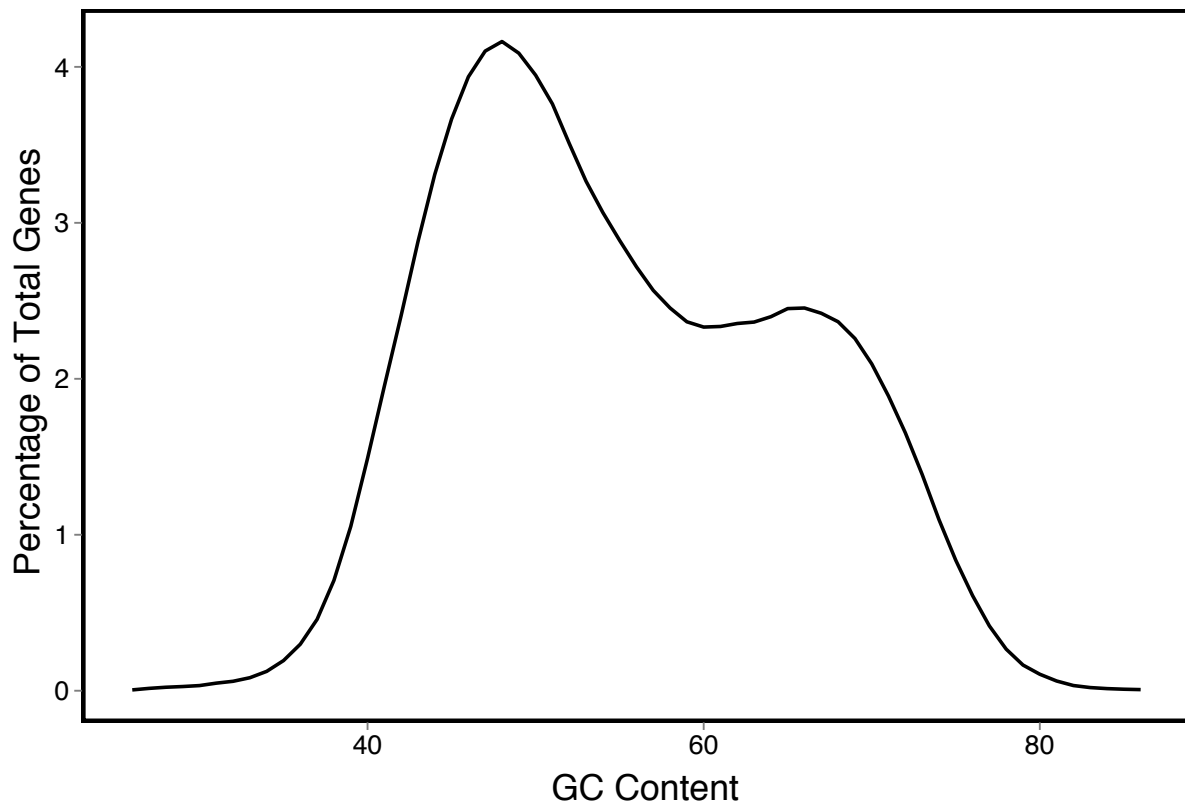
888

889

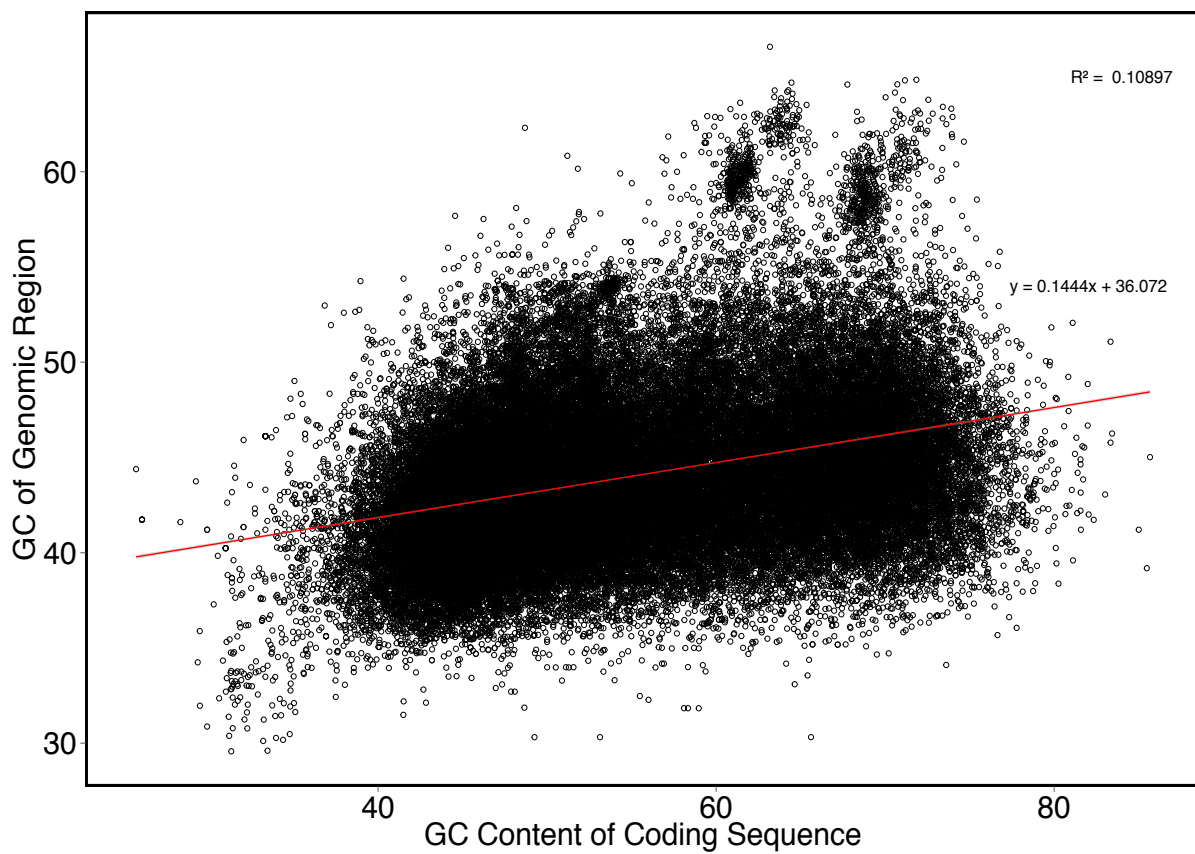
890 **Table 2. Distribution across the genome of rice of novel genes predicted by SNAP**
891 **and AUGUSTUS HMMs trained genes with high and low GC content.**

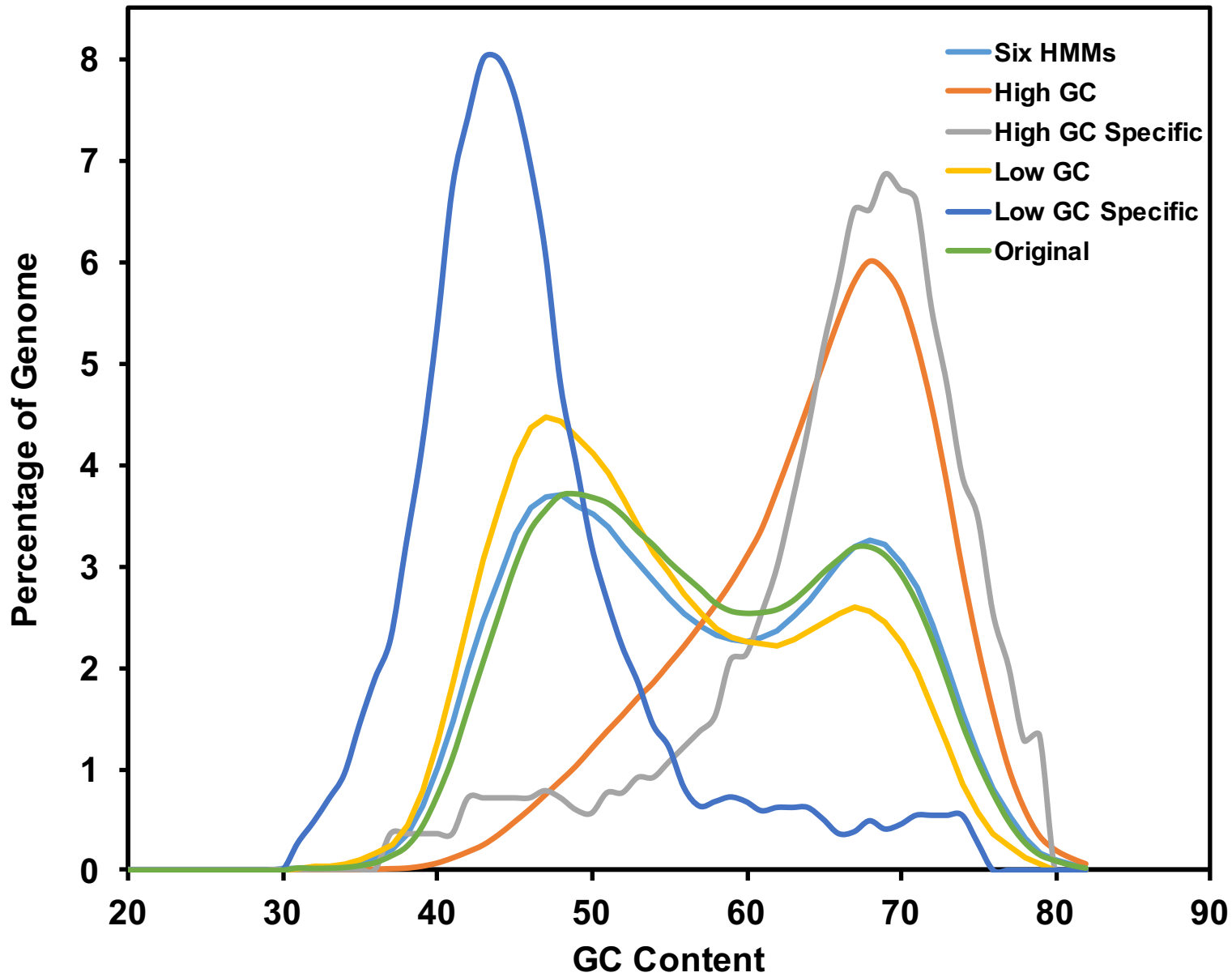
	Novel Low GC	Novel High GC ⁸⁹²
	HMM Predictions	HMM Predictions
Chr1	28	33
Chr2	39	23
Chr3	32	26
Chr4	42	35
Chr5	35	20
Chr6	31	30
Chr7	29	18
Chr8	26	22
Chr9	20	23
Chr10	29	21
Chr11	31	20
Chr12	31	26

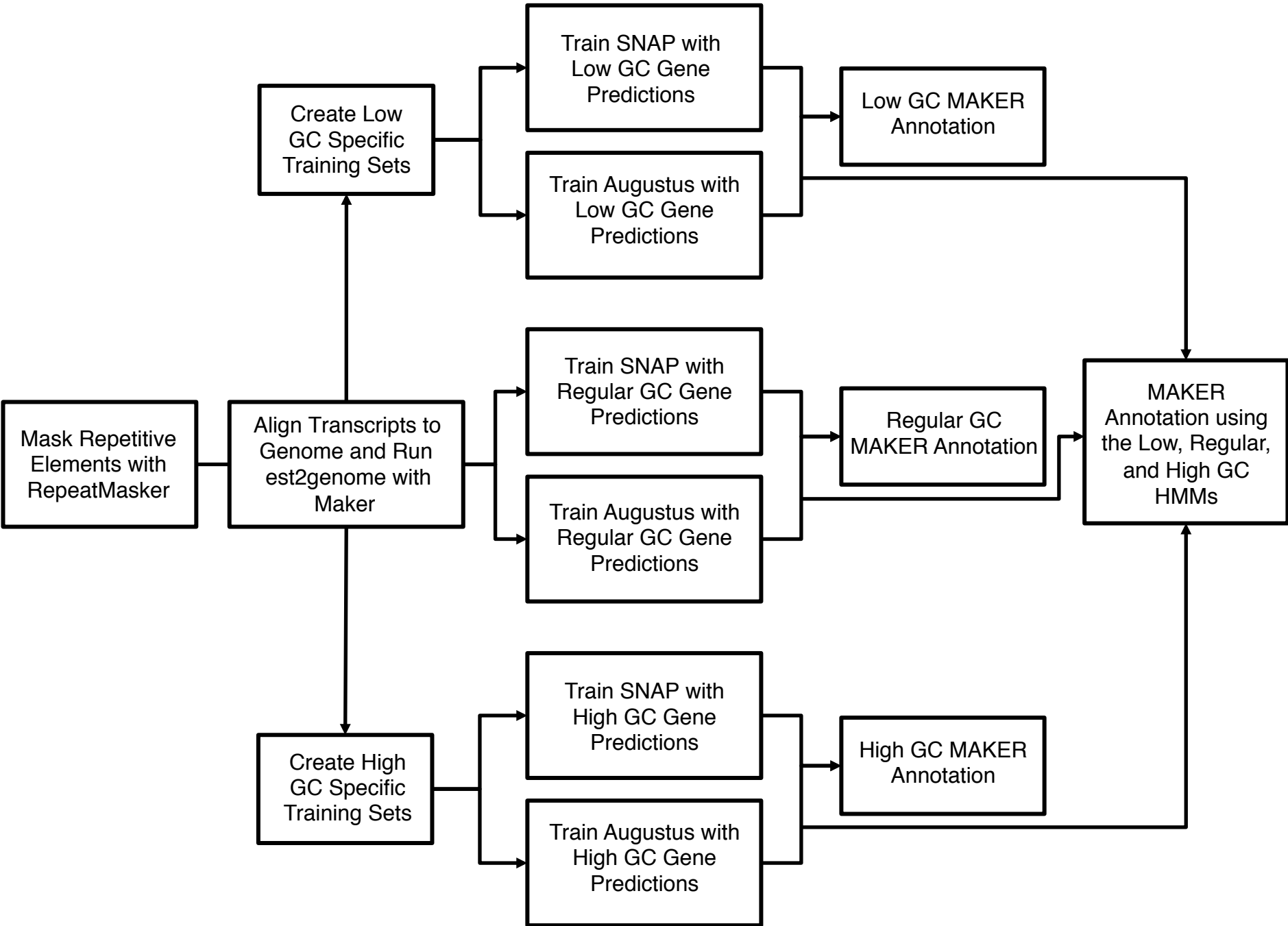
A

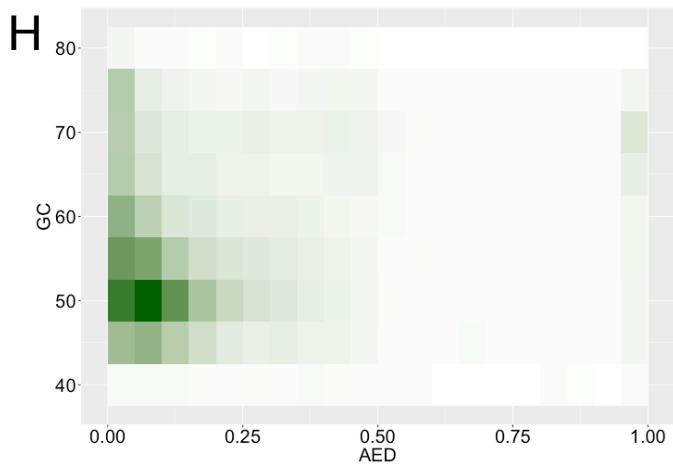
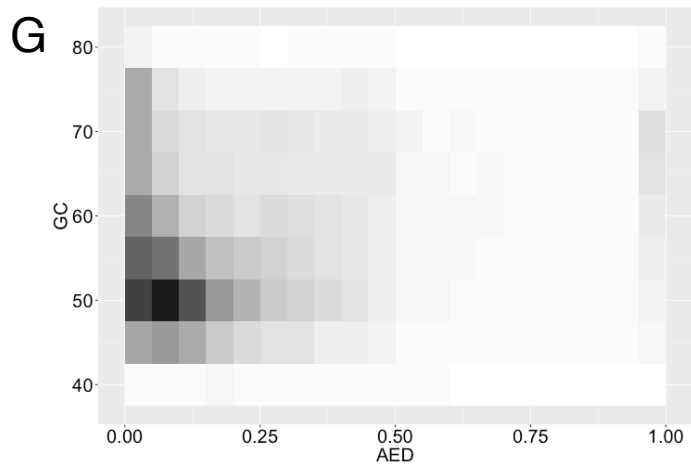
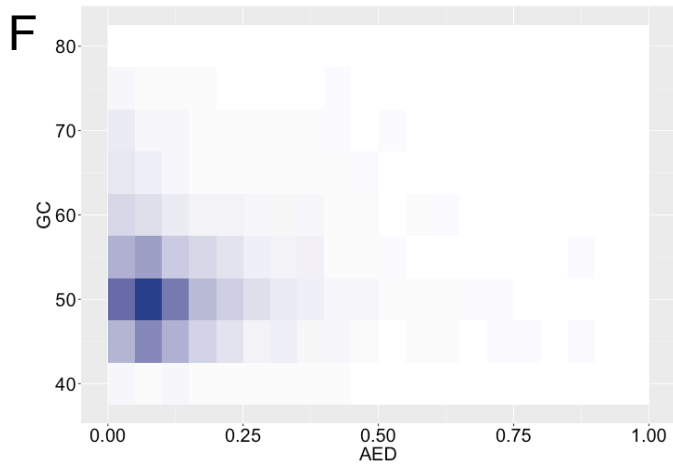
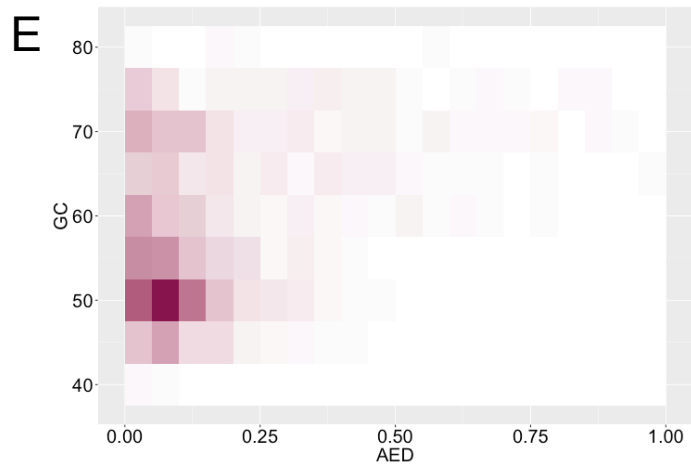
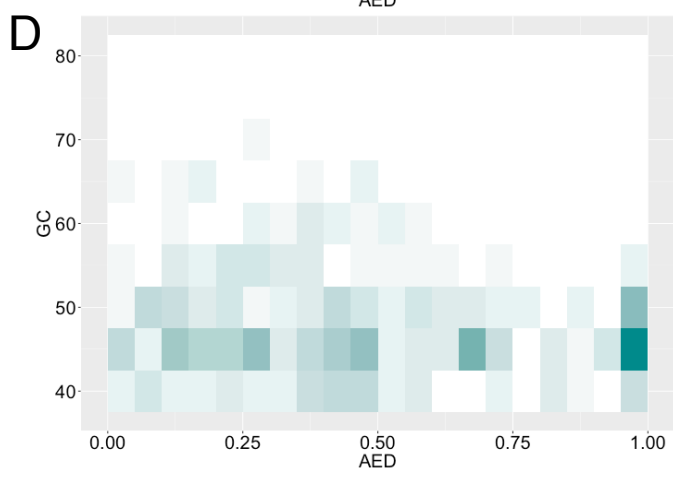
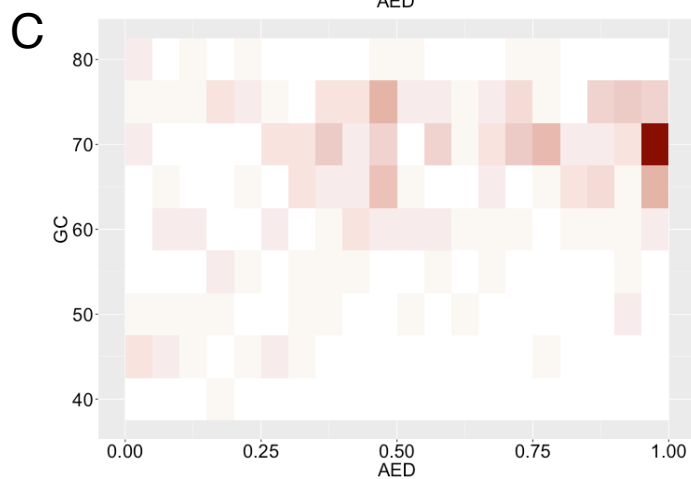
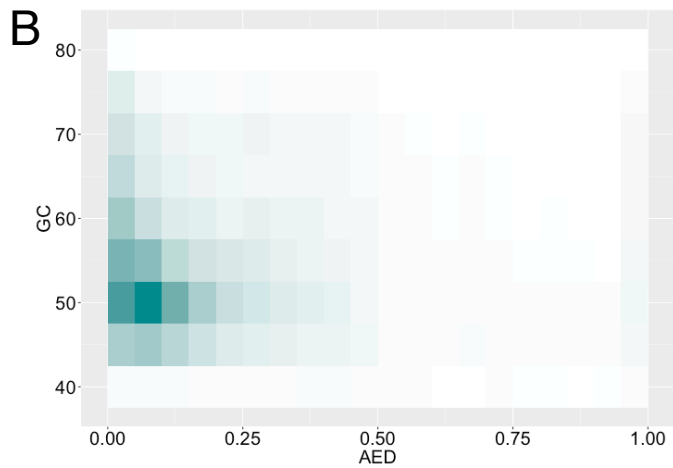
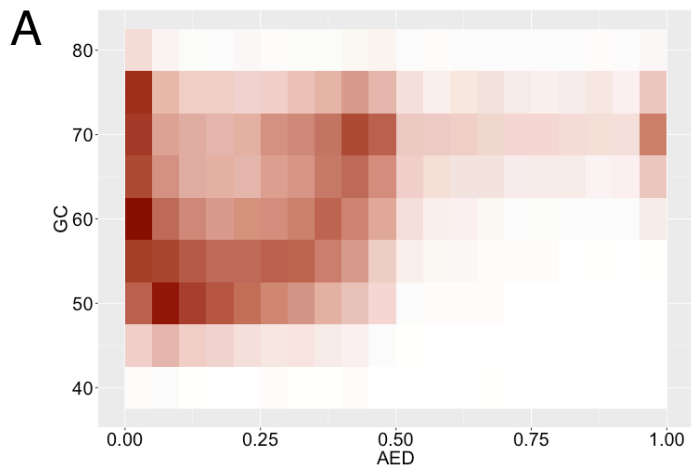


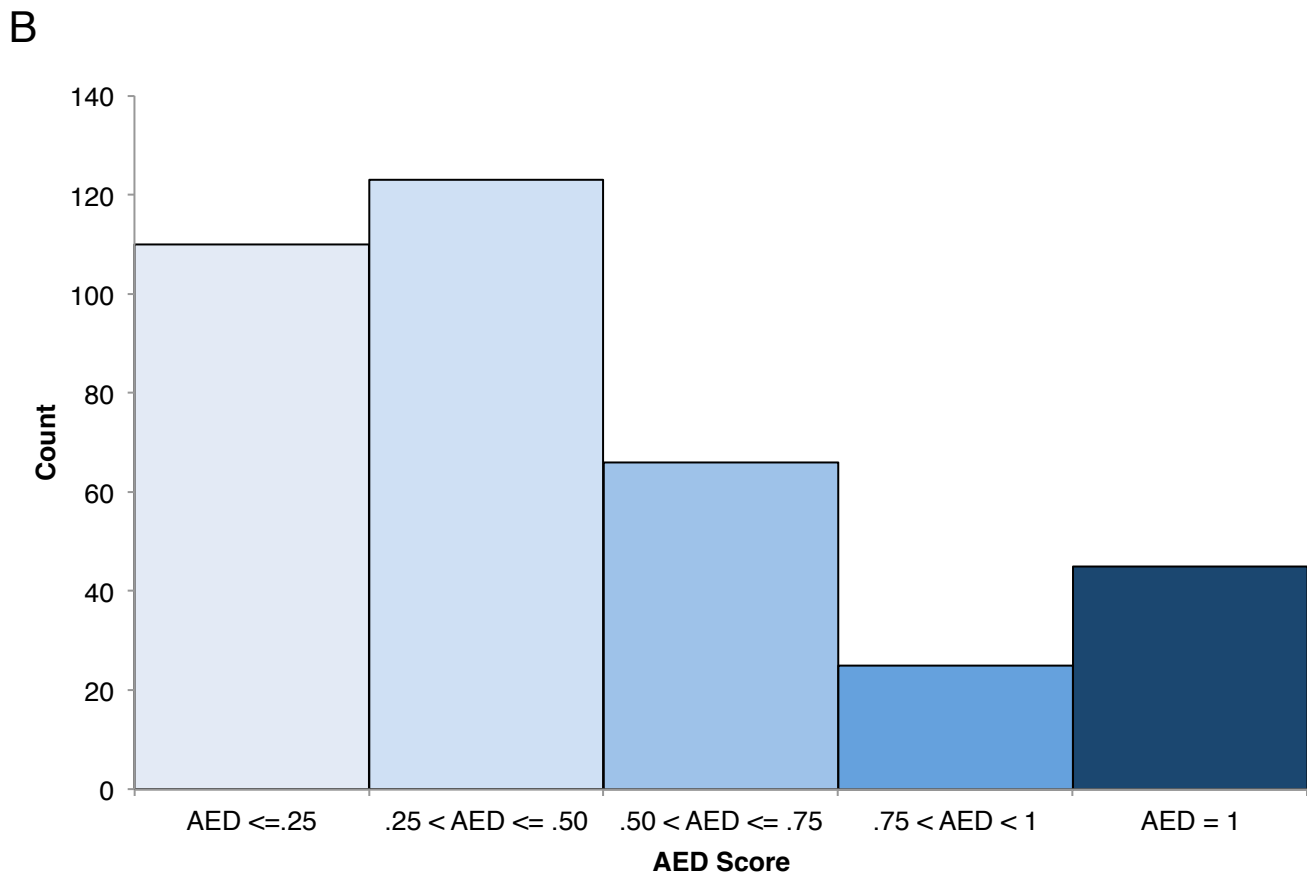
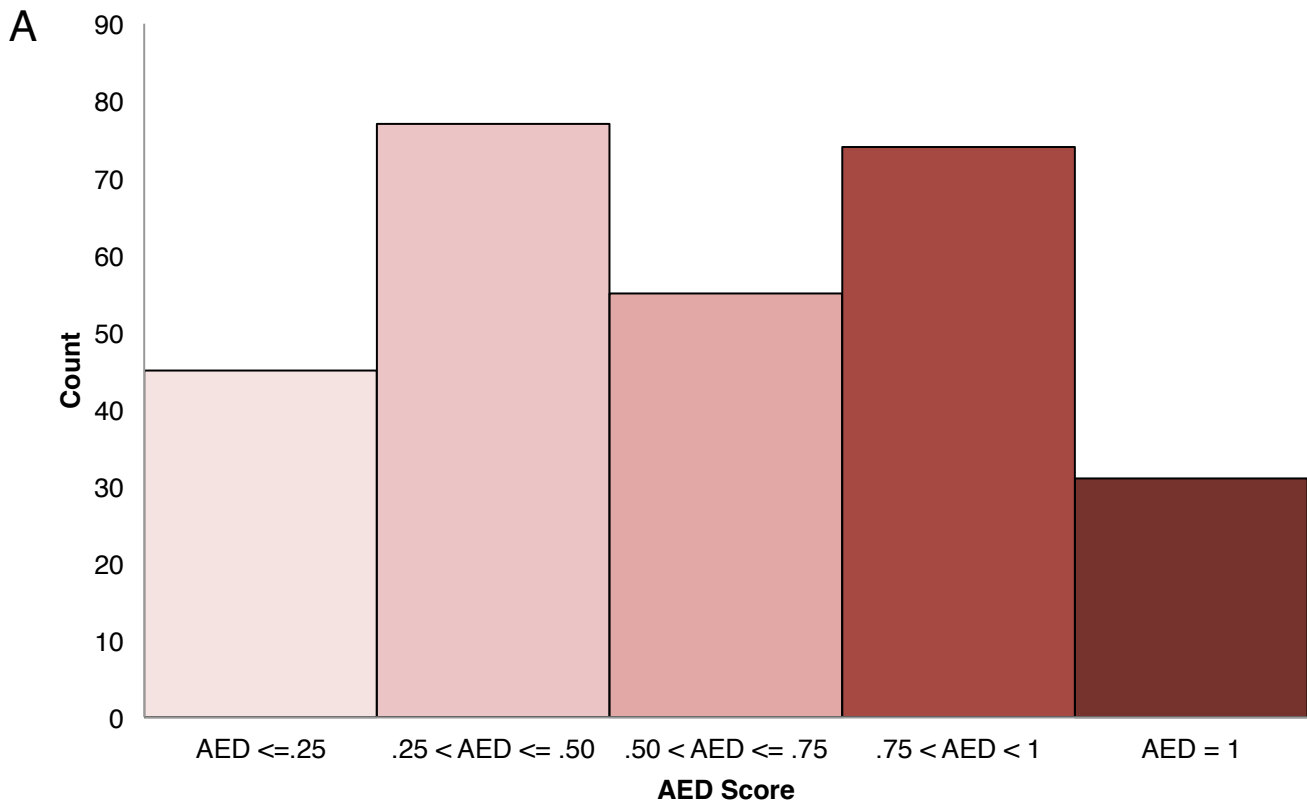
B



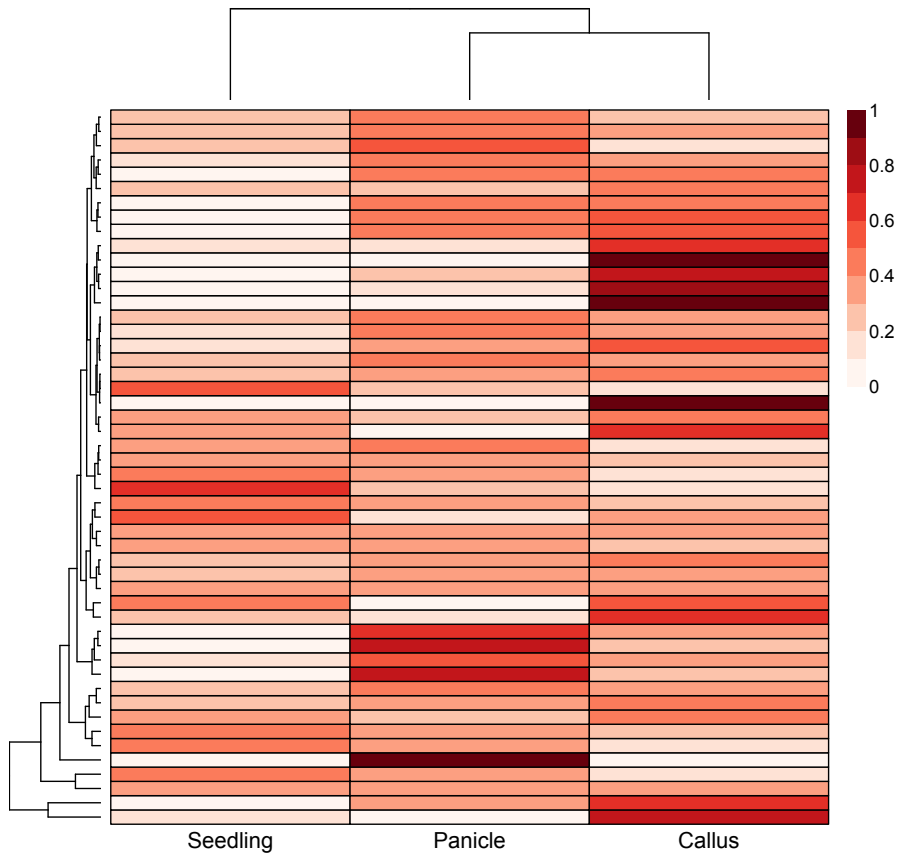








A



B

