

1                   **Reference-based QUantification Of gene**  
2                                   **Dispensability (QUOD)**

3

4 Katharina Frey<sup>1,2,\*</sup>, Bernd Weisshaar<sup>1</sup> and Boas Pucker<sup>1,3</sup>

5

6 <sup>1</sup>Genetics and Genomics of Plants, Center for Biotechnology (CeBiTec), Bielefeld  
7 University, 33615 Bielefeld, Germany

8 <sup>2</sup>Graduate School DILS, Bielefeld Institute for Bioinformatics Infrastructure (BIBI),  
9 Bielefeld University, 33615 Bielefeld, Germany

10 <sup>3</sup>Evolution and Diversity, Department of Plant Sciences, University of Cambridge,  
11 Cambridge, UK

12 \*Correspondence: [kfrey@cebitec.uni-bielefeld.de](mailto:kfrey@cebitec.uni-bielefeld.de)

13

14

15

16

17

18

19

20

21

22

23

24

25

## 26 **Abstract**

27 Dispensability of genes in a phylogenetic lineage, e.g. a species, genus, or higher-  
28 level clade, is gaining relevance as most genome sequencing projects move to a  
29 pangenome level. Most analyses classify genes as core genes, which are present in  
30 (almost) all investigated individual genomes, and dispensable genes, which only  
31 occur in a single or a few investigated genomes. The binary classification as 'core' or  
32 'dispensable' is often based on arbitrary cutoffs of presence/absence in the analysed  
33 genomes. Even when extended to 'conditionally dispensable', this concept still  
34 requires the assignment of genes to distinct groups.

35 Here, we present a new method which overcomes this distinct classification by  
36 quantifying gene dispensability and present a dedicated tool for reference-based  
37 QUantification Of gene Dispensability (QUOD). As a proof of concept, sequence data  
38 of 966 *Arabidopsis thaliana* accessions were processed to calculate a gene-specific  
39 dispensability score based on normalised coverage in read mappings. We validated  
40 this score by comparison of highly conserved Benchmarking Universal Single Copy  
41 Orthologs (BUSCOs) to all other genes. The average scores of BUSCOs were  
42 significantly higher than the scores of non-BUSCOs. Scores of genes involved in  
43 specialised metabolism were generally lower than scores of primary metabolite  
44 genes, indicating higher dispensability of genes involved in secondary metabolic  
45 processes. Analysis of variation demonstrated lower variation values between  
46 replicates of a single accession than between iteratively, randomly selected  
47 accessions from the whole dataset.

48 Instead of classifying a gene as core or dispensable, QUOD assigns a dispensability  
49 score to each gene. Hence, QUOD facilitates the identification of candidate  
50 dispensable genes which often underlie lineage-specific adaptation to varying  
51 environmental conditions.

52

53

## 54 **Introduction**

55 Genetic variation is not restricted to single nucleotide polymorphisms or small  
56 insertions and deletions but extends also to (large) structural variations. These  
57 structural variations include copy number variations (CNVs) and presence/absence  
58 variations (PAVs), which can cause substantial variation of the gene content among  
59 individual genomes<sup>1,2</sup>. The comparative analysis of multiple genomes of the same  
60 phylogenetic clade allows the identification of PAVs underlying phenotypic traits. In  
61 the case of crop species, the identification of PAVs underlying specific agronomic  
62 traits which only occur in a single or a few species is feasible<sup>3-5</sup>. As more highly  
63 contiguous genome sequences become available, pangenomes are suitable to  
64 describe and investigate the gene set diversity of a biological clade, e.g. species,  
65 genus or higher<sup>6,7</sup>.

66 Genes of a pangenome are thought to be divided into a core and a dispensable gene  
67 set, the latter is also often referred to as ‘accessory’ in the literature. Core genes  
68 occur in all or almost all investigated genomes, whereas dispensable genes only  
69 occur in a single or a few genomes<sup>8</sup>. Sometimes, a third category of ‘conditionally  
70 dispensable’ genes is invoked<sup>9</sup>. However, this distinct classification relies on one or  
71 multiple arbitrary cutoffs. Some studies consider genes as ‘core’ if these genes occur  
72 in at least 90 % of the investigated genomes<sup>10</sup>; in other studies, only genes which  
73 are found in all genomes are part of the core genome<sup>11</sup>. In addition, dependency  
74 groups might influence the dispensability of certain genes. The possibility that two  
75 genes might be ‘replaced’ by a specific number of other genes has to be considered.  
76 Some genes, of e.g. a gene family, might be required in a specific proportion and  
77 therefore are only conditionally dispensable<sup>9</sup>. Further, assemblies of genomes or  
78 transcriptomes might be incomplete leading to artificially missing genes<sup>12</sup>. One way  
79 to circumvent this is to rely on a high-quality reference genome sequence, thus  
80 avoiding additional assemblies which are potential sources of errors.

81 Here, we present QUOD - a bioinformatic tool to quantify gene dispensability. An *A.*  
82 *thaliana* dataset of about 1,000 accessions was used to calculate a per gene  
83 dispensability score derived from the coverage of all genes in the given genomes.  
84 This score was validated by comparison of scores of BUSCOs and scores of  
85 secondary metabolite genes. Our tool is easy to use for all kinds of plant species and  
86 datasets of different levels of phylogenetic resolution. QUOD extends the distinct

87 classification of genes as ‘core’ and ‘dispensable’ based on an arbitrary threshold to  
88 a continuous dispensability score.

89

90

## 91 **Methods**

### 92 **Selection and preprocessing of datasets**

93 Genomic reads (FASTQ format) of the investigated genomes were retrieved from the  
94 Sequence Read Archive (SRA) <sup>13</sup> via fastq-dump. BWA-MEM <sup>14</sup> was applied to map  
95 all genomic paired-end Illumina reads to the corresponding reference genome  
96 sequence. For *A. thaliana*, all available 1,135 datasets <sup>15</sup> (File S1) were subjected to  
97 a mapping against the AthNd-1\_v2c genome sequence <sup>16</sup>. The resulting BAM files of  
98 these mappings were subjected to QUOD (<https://github.com/KatharinaFrey/QUOD>).

99

### 100 **Calculation of gene dispensability scores – QUOD**

101 QUOD calculates a reference-based gene dispensability score for each annotated  
102 gene based on a supplied mapping file (BAM) and annotation of the reference  
103 sequence (GFF) (<https://github.com/KatharinaFrey/QUOD>). The tool is written in  
104 Python3 and consists of six different modules. During the first part of the analysis, the  
105 read coverage per position (1) as well as the read coverage per gene (2) are  
106 calculated. In the next step, genomes with an average coverage below a given cutoff  
107 (default=10) are discarded and excluded from further analyses (3). Finally, an input  
108 matrix is constructed (4) and a dispensability score is determined for each gene (5).  
109 Optionally, the results can be visualized as a colored histogram and a violin plot (6).

110 The dispensability score ( $ds(g)$ ) is calculated as follows:

$$\text{dispensability score (gene } g) = \frac{\sum \left( \frac{\text{average coverage of gene } g}{\text{mean coverage over all genes in genome } x} \right)}{\text{total number of genomes}}$$

111

### 112 **Identification of plastid sequences**

113 Genes of the *A. thaliana* test set with high similarity to plastid sequences were  
114 flagged via BLASTp<sup>17</sup> of the encoded peptides against all organelle peptide  
115 sequences obtained from the National Center for Biotechnology Information (NCBI).  
116 As a control, the sequences were also searched against themselves. Peptide  
117 sequences of Nd-1 with a score ratio  $\geq 0.8$  were considered plastid-like sequences  
118 when comparing BLAST hits against self-hits<sup>16</sup>.

119

## 120 **Score comparison between contrasting gene sets**

121 Genes annotated in AthNd-1\_v2c were classified with Benchmarking Universal  
122 Single Copy Orthologs (BUSCO) v3<sup>18</sup> running in protein mode on the encoded  
123 peptide sequences using 'brassicales odb10' as reference<sup>19</sup>. BUSCOs include  
124 single-copy genes and universal genes which are present in > 90% of all species in  
125 the reference dataset and are used to measure the completeness of assemblies and  
126 annotations<sup>18</sup>. The scores ( $\leq 2$ ) of BUSCO and non-BUSCO genes were compared  
127 using matplotlib<sup>20</sup> for visualization (boxplot) and a Mann–Whitney U test  
128 implemented in the Python package SciPy<sup>21</sup> for determination of the significance.

129 Primary and secondary metabolite genes were identified using the KEGG Orthology  
130 for *A. thaliana* (ath00001.keg) downloaded from the KEGG database<sup>22</sup> and a list of  
131 primary and secondary metabolic pathways obtained from Mukherjee *et al.*<sup>23</sup>. In  
132 total, 1,887 genes of Nd-1 were assigned to the primary metabolism, whereas 440  
133 genes were assigned to the secondary metabolism. The average scores of these  
134 groups of genes were compared statistically by a Mann-Whitney U test within a  
135 custom Python script using the package SciPy<sup>21</sup> and a boxplot was created using  
136 matplotlib.

137 A list of Nd-1 transposable element (TE) genes was obtained from Pucker *et al.*<sup>16</sup>.  
138 First, the score distribution ( $ds > 2$  is set to  $ds = 2$ ) of TE and non-TE genes was  
139 determined using a Mann–Whitney U test implemented in the Python package SciPy  
140<sup>21</sup>. Next, the minimal distance of each gene to its closest TE gene was calculated  
141 after extracting the gene positions from the Nd-1 annotation file. Pearson's  
142 correlation coefficient between per gene distance to TEs and the gene dispensability  
143 score ( $ds > 2$  is set to  $ds = 2$ ) was determined.

144

## 145 **Correlation of gene length and exon number with the dispensability score**

146 Length and number of exons per gene were extracted from the Nd-1 annotation file.  
147 The Pearson's correlation coefficients of gene length and exon number, respectively,  
148 with the gene dispensability score ( $ds > 2$  is set to  $ds = 2$ ) were calculated via SciPy  
149 <sup>21</sup> for the whole dataset as well as for three large *A. thaliana* gene families (TAPscan  
150 <sup>24</sup>), namely MYBs, AP2/EREBP and WRKYs .

151

## 152 **Variation between replicates**

153 A total of 14 genomic datasets of the *A. thaliana* accession Col-0 were received from  
154 the SRA (File S2) to assess the technical variation between replicates of the same  
155 accession. Each dataset was mapped to the TAIR10 reference genome sequence  
156 using BWA-MEM. The mappings were then subjected to QUOD expecting a  
157 dispensability score close to one for each gene. As the distributions are different  
158 (Kolmogorov-Smirnov test,  $p \approx 3e-27$ ) and the sample size ( $n$ ) is high, the Levene's  
159 test was selected to test for equal variances, regarding the gene dispensability  
160 scores. The test was applied for (1) the dataset including replicates only and (2)  
161 iteratively (100x), randomly chosen subsets ( $n=14$ ) of the dataset including all  
162 available *A. thaliana* accessions.

163

## 164 **Functional annotation**

165 All genes of the *A. thaliana* Nd-1 genome sequence were functionally annotated via  
166 reciprocal best blast hits (RBHs) and best BLAST hits against Araport11 <sup>16</sup>.

167

## 168 **Data Availability**

169 The tool QUOD (QUOD.py) can be downloaded from GitHub  
170 (<https://github.com/KatharinaFrey/QUOD>).

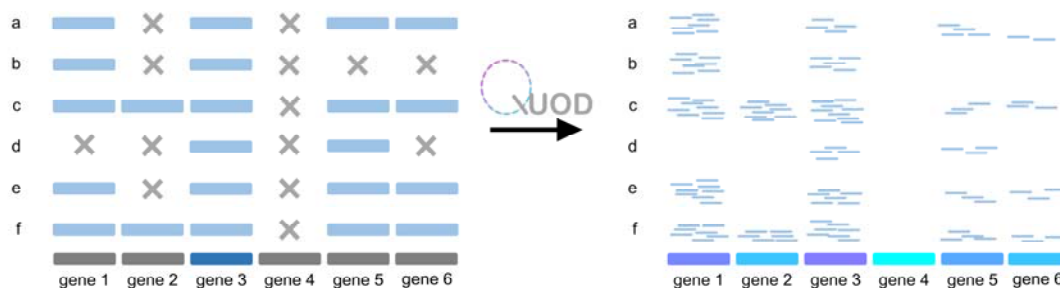
171

172

## 173 **Results**

174 In this study, a bioinformatic tool was developed to calculate a gene-specific  
175 dispensability score based on the normalised coverage in a read mapping. QUOD  
176 allows the quantification of dispensability by calculation of a score for each gene  
177 (Figure 1). The binary classification of gene dispensability can be compared to the  
178 original method of mRNA detection by endpoint RT-PCR providing only qualitative  
179 results<sup>25-27</sup> which was replaced by quantitative analyses like RNA-Seq.

180 The gene dispensability score would initially be dependent on the sequencing depth  
181 per genome. By division of the average coverage of gene *g* by the mean coverage  
182 over all genes in genome *x*, the score is normalised for differences in the sequencing  
183 read depth of the investigated genomes. A low value indicates that a gene is likely to  
184 be missing in some genomes and therefore more likely dispensable than a gene with  
185 a higher dispensability score. Due to this quantification approach, this method is not  
186 based on an arbitrary cutoff to determine the core genome and the dispensable  
187 genome of any given pangenome dataset. An example: Using a cutoff of 'gene *x*  
188 occurs in at least 90 % of all genomes' to be considered a 'core' gene, genes 1,2,4,5  
189 and 6 (dark grey) would be considered 'dispensable' (Figure 1). However,  
190 considering the coverage (right panel), is e.g. gene 1 truly dispensable?

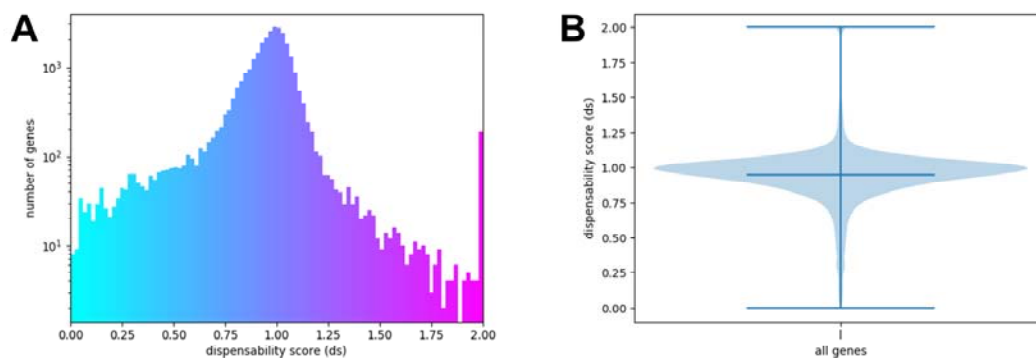


191

192 Figure 1: Illustration of the QUOD method using an artificial dataset. On the left side,  
193 genes are classified as 'core' or 'dispensable' according to a cutoff. On the right side,  
194 gene dispensability is quantified according to a dispensability score based on the  
195 normalised coverage in a read mapping (a-f: investigated genomes). Coloring of  
196 genes (right side) indicates different dispensability scores. Extremely rare genes can  
197 be easily detected using QUOD.

198 As a proof of concept, *A. thaliana* sequencing datasets of 1,135 accessions were  
199 downloaded and mapped to the Nd-1 genome sequence. All genomes with less than  
200 10-fold read coverage were discarded. The remaining sequencing datasets covering

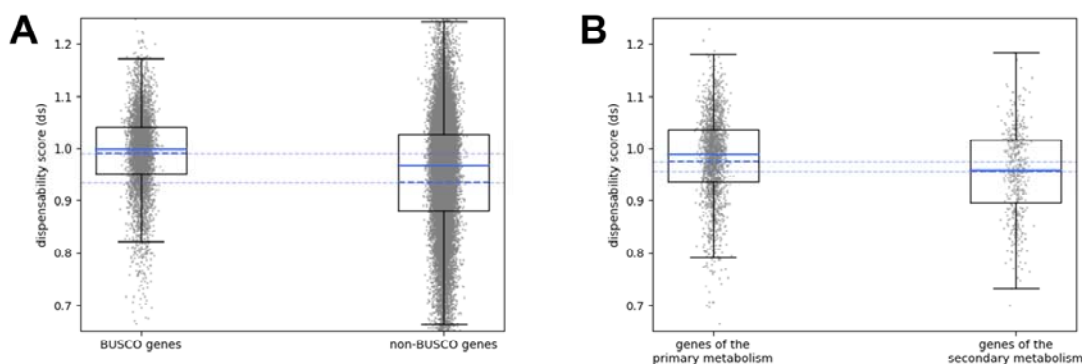
201 966 accessions were analysed with QUOD to calculate a dispensability score for  
202 each gene (Figure 2).



203

204 Figure 2: Distribution of the gene dispensability scores for the *A. thaliana* datasets. A)  
205 Histogram colored according to the dispensability score. The x-axis represents the  
206 dispensability score, whereas the y-axis shows the number of genes in each bin. B)  
207 Violin plot representing the dispensability score (y-axis) of all genes (x-axis).

208 Further validation of the reliability of the gene dispensability quantification was  
209 achieved by comparison of BUSCOs and non-BUSCOs (Figure 3A) as well as  
210 primary and secondary metabolite genes (Figure 3B). BUSCO genes show  
211 significantly higher scores than non-BUSCO genes (U test,  $p \approx 1e-122$ ). Secondary  
212 metabolite genes show significantly lower scores than primary metabolite genes (U  
213 test,  $p \approx 6e-11$ ).



214

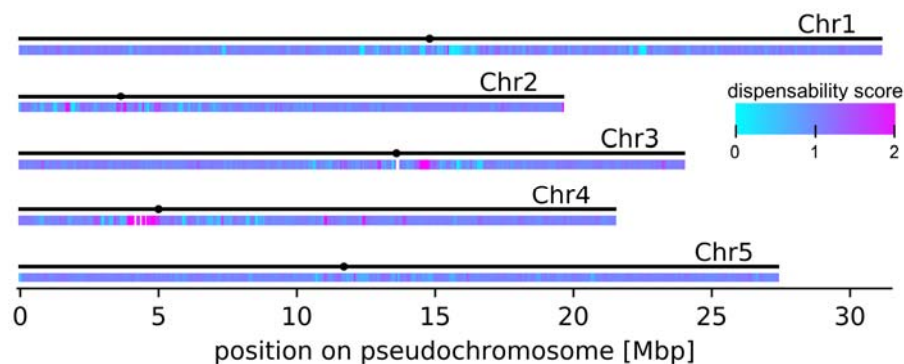
215 Figure 3: Box plot of the gene dispensability scores of (A) BUSCO and non-BUSCO  
216 genes and (B) primary and secondary metabolite genes. The mean is represented by  
217 the dashed, blue line, the other blue line represents the median of the scores. A)



218 BUSCO genes (n = 4586, mean ds  $\approx$  0.991) show significantly higher scores than  
219 non-BUSCO genes (n = 25354, mean ds  $\approx$  0.934) (U test,  $p \approx 1e-122$ ). B) Secondary  
220 metabolite genes (n = 440, mean ds  $\approx$  0.957) show significantly lower scores than  
221 primary metabolite genes (n = 1,887, mean ds  $\approx$  0.975) (U test,  $p \approx 6e-11$ ).

222 Functional annotation of BUSCO outliers, which include genes with dispensability  
223 scores below 0.75 or above 1.25, revealed, amongst others, several repeat proteins,  
224 transmembrane proteins, a 'stress induced protein', and multiple hypothetical  
225 proteins (File S3).

226 Next, the genome-wide distribution of genes with specific gene dispensability scores  
227 was investigated in *A. thaliana*. A high dynamic across centromeric regions is visible  
228 based on the colored heatmap (Figure 4). High and low scoring genes cluster in  
229 centromeric and telomeric regions.



230

231 Figure 4: Genome-wide distribution of genes with different dispensability scores in *A.*  
232 *thaliana* Nd-1. The colored heatmap shows the respective gene dispensability  
233 scores. There are high (pink) and low (blue) scoring genes clustered in repetitive  
234 regions, including centromeric and telomeric areas. The x-axis represents the size (in  
235 Mbp) of each pseudochromosome in the assembly. The black dots represent the  
236 position of the centromeres of the five chromosomes in the *AthNd1\_v2c* assembly <sup>16</sup>.

237 As high and low scoring genes cluster in repetitive regions (mainly centromeres), the  
238 score distribution of TEs was investigated (File S4). Scores of TE genes are evenly  
239 distributed across all dispensability scores. In total, the mean score of TE genes  
240 (mean ds  $\approx$  0.919) is significantly lower when compared to non-TE genes (mean ds  $\approx$   
241 0.953) (U test,  $p \approx 6e-08$ ), which are more frequent across scores close to one.  
242 Moreover, the minimal distance of each gene to its closest TE gene and the  
243 dispensability scores revealed no correlation (File S4).

244 To test the hypothesis whether genes with lower dispensability scores/more likely  
245 dispensable genes are shorter and whether introns accumulate in core genes, the  
246 correlation of the gene dispensability score with gene length and exon number,  
247 respectively, were determined for the whole dataset and for three selected gene  
248 families separately. However, no clear trend was detectable (File S5).

249 Further, the variation between replicates of the same accession (Col-0) was  
250 determined (File S6). The variation of the gene dispensability score distribution of the  
251 replicate dataset (one accession) ( $\sigma^2 \approx 0.0226$ ) is significantly lower than the  
252 variation between all iteratively, randomly selected subsets of *A. thaliana* accessions  
253 ( $\sigma^2 \approx 0.0392$ ) (Levene's test,  $p \approx 4e-19$ ).

254 The function of the 100 genes with the lowest gene dispensability scores was  
255 examined in detail for the *A. thaliana* (File S7) dataset. Fourteen genes of this *A.*  
256 *thaliana* test set are annotated as “disease resistance proteins”, whereas seven  
257 genes are annotated as transposons/transposases. Four genes are described as  
258 hypothetical proteins and 24 genes have no functional annotation.

259

260

## 261 **Discussion**

262 QUOD was developed for the quantification of gene dispensability in plant  
263 pangenome datasets, but the application is not limited to the plant kingdom.  
264 However, an accurate determination of gene dispensability scores free of systematic  
265 biases might rely on a uniform selection of genomes from the respective taxonomic  
266 group and on uniform read coverage of genes. In addition, non-random  
267 fragmentation of DNA prior to sequencing<sup>28</sup> may cause biases. The variation among  
268 replicates of the same accession (Col-0;  $\sigma^2 \approx 0.0226$ ) might be attributed to technical  
269 biases, e.g. during library preparation. We validated the reliability of the gene  
270 dispensability score by showing that BUSCO genes are more conserved and can be  
271 more likely considered core genes than non-BUSCOs as BUSCOs are associated  
272 with significantly higher scores (Figure 3A). In addition, secondary metabolite genes  
273 are more likely dispensable as shown by the significantly lower dispensability scores  
274 in comparison to genes of primary metabolic processes (Figure 3B).

275 Based on the distribution of the scores in the boxplot (Figure 3), the difference  
276 between BUSCOs and non-BUSCOs appears small, even though the difference is  
277 significant (U test,  $p \approx 1e-122$ ). The expected difference in the dispensability scores  
278 of BUSCOs and non-BUSCOs might not be exceptionally high as multiple-copy  
279 genes are not included in the BUSCO gene set <sup>18</sup>. Further, functional annotation of  
280 BUSCO outliers revealed several repeat proteins and transmembrane proteins.  
281 Repeat proteins might lead to read mapping errors and consequently artificial  
282 variations in coverage and dispensability scores. Transmembrane proteins could be  
283 involved in biotic stress response and therefore might not be essential for some  
284 accessions <sup>29</sup>. This could explain the absence in some genomes resulting in low  
285 dispensability scores of these genes. Therefore, many important, high-scoring genes  
286 might lie outside of the BUSCO sampling space. The difference between scores of  
287 primary and secondary metabolite genes might be higher as well, as strongly variable  
288 secondary metabolite genes might not be detected in our analyses due to the RBH-  
289 based annotation.

290 The genome-wide distribution of all gene dispensability scores of the *A. thaliana*  
291 genomes reveals the origin of exceptionally high dispensability scores (Figure 4).  
292 High and low scoring genes cluster in repetitive regions, like e.g. centromeres or  
293 telomeres, due to variations in the recombination rate <sup>30</sup> and active transposable  
294 elements in these regions.

295 It was previously proposed, that likely dispensable genes are located closer to  
296 transposable elements which are important factors in pangenome evolution <sup>31</sup>. In our  
297 study, TE genes are widely distributed across all dispensability scores as TEs can  
298 occur in a high copy number in genomes leading to high scores and can as well be  
299 dispensable. Other studies detected a high number of TEs in the dispensable  
300 genome <sup>32</sup>, however, only certain TE families might be truly dispensable. A high-  
301 quality annotation of transposons and a following exclusion of these genes from the  
302 analysis or improved read mapping to the consensus sequence might improve the  
303 results. Moreover, heterochromatin or genome-purging mechanisms <sup>33</sup> could  
304 influence the gene dispensability scores in these regions. Additionally some of the  
305 high scoring genes were flagged as plastid-like sequences as original sequencing  
306 data from plants contain a high amount reads originating from plastid sequences <sup>34,35</sup>.  
307 Biases due to this plastid read contamination inflate the coverage of sequences with

308 high similarity to plastid sequences, resulting in an exceptionally high gene  
309 dispensability score. Therefore, especially genes with dispensability scores above  
310 two must be investigated carefully. Very similar sequences, e.g. members of a gene  
311 family or close paralogs, might cause read mapping errors confounding biases in the  
312 dispensability scores of these genes.

313 Even though it was previously proposed that dispensable genes are located close to  
314 transposable elements<sup>31</sup>, this result was not observed in our study. One limitation is  
315 the accurate assignment of reads to repetitive sections of the reference sequence  
316 during the read mapping<sup>12</sup>. Further, only a fraction of transposons might be correctly  
317 assembled and annotated due to several computational challenges in highly  
318 repetitive and peri-centromeric regions<sup>36</sup>. Therefore, a different strategy might be  
319 needed to accurately quantify dispensability of transposable elements.

320 Functional annotation of the 100 most likely dispensable genes revealed a high  
321 number of uncharacterised proteins, disease resistance proteins as well as  
322 transposons and transposases in the *A. thaliana* genomes. Transposable elements  
323 were detected in other studies as contributors to large structural variations between  
324 species and individuals and substantial part of the dispensable genome<sup>32</sup>. Previous  
325 pangenome analyses also revealed that the dispensable genome comprises  
326 functions like e.g. 'defense response', 'diseases resistance', 'flowering time' and  
327 'adaptation to biotic and abiotic stress'<sup>10,31,37</sup>. Therefore, in depth investigation of  
328 genes with low dispensability scores can result in the identification and  
329 characterisation of phenotypic variation<sup>38</sup> and important agronomic traits<sup>37</sup>. We  
330 envision several applications for the gene dispensability score generated by QUOD:  
331 (1) more accurate prediction if a PAV is associated with a specific trait, (2)  
332 development of dependency gene networks, and (3) improved modeling of the  
333 evolutionary value of genes.

334

## 335 **Conclusion**

336 QUOD (reference-based QUantification Of gene Dispensability) overcomes the  
337 problem of labeling genes as 'core' or 'dispensable' through implementation of a  
338 quantification approach. Instead of a distinct classification, QUOD provides a ranking

339 of all genes based on assigned gene-specific dispensability scores and therefore  
340 does not rely on any thresholds.

341

342

### 343 **Author Contributions**

344 KF, BW and BP designed the study, performed the experiments, analysed the data,  
345 and wrote the manuscript. All authors read and approved the final version of this  
346 manuscript.

### 347 **Funding**

348 This research received no external funding.

### 349 **Acknowledgments**

350 We thank members of Genetics and Genomics of Plants for discussion of preliminary  
351 results. We are very grateful to Janik Sielemann and Nathanael Walker-Hale for  
352 helpful comments on the manuscript. We acknowledge support for the Article  
353 Processing Charge by the Deutsche Forschungsgemeinschaft and the Open Access  
354 Publication Fund of Bielefeld University. We thank the CeBiTec Bioinformatic  
355 Resource Facility team for great technical support.

### 356 **Conflicts of Interest**

357 The authors declare no conflict of interest.

### 358 **Supplementary Files**

359 File S1: SRA IDs of datasets downloaded to conduct the QUOD analysis of the *A.*  
360 *thaliana* genomes.

361 File S2: SRA IDs of datasets downloaded to conduct the analysis of replicates (Col-  
362 0).

363 File S3: Functional annotation of BUSCO outliers with a dispensability score smaller  
364 than 0.75 or greater than 1.25.

365 File S4: Distribution of scores of TE genes and non-TE genes and correlation of the  
366 distance to the closest TE gene with the gene dispensability score of the *A. thaliana*  
367 genomes.

368 File S5: Correlation of gene length and exon number with the dispensability scores of  
369 the *A. thaliana* genomes.

370 File S6: Analysis of variance of the gene dispensability score calculated for replicates  
371 of the *A. thaliana* Col-0 accession and iteratively, randomly chosen subsets of the  
372 whole dataset.

373 File S7: Functional annotation of the 100 most likely dispensable genes of the *A.*  
374 *thaliana* genomes.

375

376

## 377 **References**

378 1. Springer, N. M. *et al.* Maize Inbreds Exhibit High Levels of Copy Number Variation  
379 (CNV) and Presence/Absence Variation (PAV) in Genome Content. *PLoS Genet* **5**,  
380 e1000734 (2009).

381 2. Scherer, S. W. *et al.* Challenges and standards in integrating surveys of structural  
382 variation. *Nat Genet* **39**, S7–S15 (2007).

383 3. Tao, Y., Zhao, X., Mace, E., Henry, R. & Jordan, D. Exploring and Exploiting Pan-  
384 genomics for Crop Improvement. *Molecular Plant* **12**, 156–169 (2019).

385 4. Lu, F. *et al.* High-resolution genetic mapping of maize pan-genome sequence  
386 anchors. *Nat Commun* **6**, 6914 (2015).

387 5. Swanson-Wagner, R. A. *et al.* Pervasive gene content variation and copy number  
388 variation in maize and its undomesticated progenitor. *Genome Research* **20**,  
389 1689–1699 (2010).

- 390 6. Tettelin, H. *et al.* Genome analysis of multiple pathogenic isolates of *Streptococcus*  
391 *agalactiae*: Implications for the microbial ‘pan-genome’. *Proceedings of the*  
392 *National Academy of Sciences* **102**, 13950–13955 (2005).
- 393 7. Vernikos, G., Medini, D., Riley, D. R. & Tettelin, H. Ten years of pan-genome  
394 analyses. *Current Opinion in Microbiology* **23**, 148–154 (2015).
- 395 8. Golicz, A. A., Batley, J. & Edwards, D. Towards plant pangenomics. *Plant*  
396 *Biotechnol J* **14**, 1099–1105 (2016).
- 397 9. Marroni, F., Pinosio, S. & Morgante, M. Structural variation and genome  
398 complexity: is dispensable really dispensable? *Current Opinion in Plant Biology* **18**,  
399 31–36 (2014).
- 400 10. Zhao, Q. *et al.* Pan-genome analysis highlights the extent of genomic variation  
401 in cultivated and wild rice. *Nat Genet* **50**, 278–284 (2018).
- 402 11. Li, Y. *et al.* De novo assembly of soybean wild relatives for pan-genome  
403 analysis of diversity and agronomic traits. *Nat Biotechnol* **32**, 1045–1052 (2014).
- 404 12. Alkan, C., Sajjadian, S. & Eichler, E. E. Limitations of next-generation genome  
405 sequence assembly. *Nat Methods* **8**, 61–65 (2011).
- 406 13. Leinonen, Rasko and Sugawara, Hideaki and Shumway, Martin and  
407 International Nucleotide Sequence Database Collaboration. The sequence read  
408 archive. *Nucleic Acids Research* **39**, D19–D21 (2010).
- 409 14. Li, H. Aligning sequence reads, clone sequences and assembly contigs with  
410 BWA-MEM. *arXiv preprint arXiv:1303.3997* (2013).
- 411 15. Alonso-Blanco, C. *et al.* 1,135 Genomes Reveal the Global Pattern of  
412 Polymorphism in *Arabidopsis thaliana*. *Cell* **166**, 481–491 (2016).
- 413 16. Pucker, B. *et al.* A chromosome-level sequence assembly reveals the  
414 structure of the *Arabidopsis thaliana* Nd-1 genome and its gene set. *PLoS ONE*  
415 **14**, e0216233 (2019).



- 416 17. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local  
417 alignment search tool. *Journal of Molecular Biology* **215**, 403–410 (1990).
- 418 18. Simão, Felipe A and Waterhouse, Robert M and Ioannidis, Panagiotis and  
419 Kriventseva, Evgenia V and Zdobnov, Evgeny M. BUSCO: assessing genome  
420 assembly and annotation completeness with single-copy orthologs. *Bioinformatics*  
421 **31**, 3210–3212 (2015).
- 422 19. Zdobnov, E. M. *et al.* OrthoDB v9.1: cataloging evolutionary and functional  
423 annotations for animal, fungal, plant, archaeal, bacterial and viral orthologs.  
424 *Nucleic Acids Res* **45**, D744–D749 (2017).
- 425 20. Hunter, J. D. Matplotlib: A 2D Graphics Environment. *Comput. Sci. Eng.* **9**, 90–  
426 95 (2007).
- 427 21. Jones, E., Oliphant, T., Peterson, P. & others. SciPy: Open source scientific  
428 tools for Python. (2001).
- 429 22. Kanehisa, M. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic*  
430 *Acids Research* **28**, 27–30 (2000).
- 431 23. Mukherjee, D., Mukherjee, A. & Ghosh, T. C. Evolutionary Rate Heterogeneity  
432 of Primary and Secondary Metabolic Pathway Genes in *Arabidopsis thaliana*.  
433 *Genome Biol Evol* **8**, 17–28 (2016).
- 434 24. Wilhelmsson, P. K. I., Mühlich, C., Ullrich, K. K. & Rensing, S. A.  
435 Comprehensive Genome-Wide Classification Reveals That Many Plant-Specific  
436 Transcription Factors Evolved in Streptophyte Algae. *Genome Biology and*  
437 *Evolution* **9**, 3384–3397 (2017).
- 438 25. Wang, A. M., Doyle, M. V. & Mark, D. F. Quantitation of mRNA by the  
439 polymerase chain reaction. *Proceedings of the National Academy of Sciences* **86**,  
440 9717–9721 (1989).



- 441 26. Gilliland, G., Perrin, S., Blanchard, K. & Bunn, H. F. Analysis of cytokine  
442 mRNA and DNA: detection and quantitation by competitive polymerase chain  
443 reaction. *Proceedings of the National Academy of Sciences* **87**, 2725–2729 (1990).
- 444 27. Chiang, P. W. *et al.* Use of a fluorescent-PCR reaction to detect genomic  
445 sequence copy number and transcriptional abundance. *Genome Research* **6**,  
446 1013–1026 (1996).
- 447 28. Poptsova, M. S. *et al.* Non-random DNA fragmentation in next-generation  
448 sequencing. *Sci Rep* **4**, 4532 (2015).
- 449 29. Dodds, P. N. & Rathjen, J. P. Plant immunity: towards an integrated view of  
450 plant–pathogen interactions. *Nat Rev Genet* **11**, 539–548 (2010).
- 451 30. Nachman, M. Variation in recombination rate across the genome: evidence  
452 and implications. *Current Opinion in Genetics & Development* **12**, 657–663 (2002).
- 453 31. Gordon, S. P. *et al.* Extensive gene content variation in the Brachypodium  
454 distachyon pan-genome correlates with population structure. *Nat Commun* **8**, 2184  
455 (2017).
- 456 32. Morgante, M., Depaoli, E. & Radovic, S. Transposable elements and the plant  
457 pan-genomes. *Current Opinion in Plant Biology* **10**, 149–155 (2007).
- 458 33. Lee, S.-I. & Kim, N.-S. Transposable Elements and Genome Size Variations in  
459 Plants. *Genomics Inform* **12**, 87 (2014).
- 460 34. Yu, J. *et al.* A draft sequence of the rice genome (*Oryza sativa* L. ssp. indica).  
461 *Science* **296**, 79–92 (2002).
- 462 35. Tang, J. *et al.* A Comparison of Rice Chloroplast Genomes. *Plant Physiol.*  
463 **135**, 412–420 (2004).
- 464 36. Platt, R. N., Blanco-Berdugo, L. & Ray, D. A. Accurate Transposable Element  
465 Annotation Is Vital When Analyzing New Genome Assemblies. *Genome Biol Evol*  
466 **8**, 403–410 (2016).

- 467 37. Golicz, A. A. *et al.* The pangenome of an agronomically important crop plant  
468 *Brassica oleracea*. *Nat Commun* **7**, 13390 (2016).
- 469 38. Hirsch, C. N. *et al.* Insights into the Maize Pan-Genome and Pan-  
470 Transcriptome. *Plant Cell* **26**, 121–135 (2014).
- 471