

1 **Reference-based QUantification Of gene**

2 **Dispensability (QUOD)**

3

4Katharina Sielemann^{1,2}, Bernd Weisshaar^{1,*} and Boas Pucker^{1,3}

6¹Genetics and Genomics of Plants, Center for Biotechnology (CeBiTec) & Faculty of
7Biology, Bielefeld University, 33615 Bielefeld, Germany

8²Graduate School DILS, Bielefeld Institute for Bioinformatics Infrastructure (BIBI),
9Bielefeld University, 33615 Bielefeld, Germany

10³Evolution and Diversity, Department of Plant Sciences, University of Cambridge,
11Cambridge, UK

12*Correspondence: bernd.weisshaar@uni-bielefeld.de

13

14Email addresses:

15kfrey@cebitec.uni-bielefeld.de

16bpucker@cebitec.uni-bielefeld.de

17bernd.weisshaar@uni-bielefeld.de

18

19

20 **Abstract**

21 **Background**

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

30 **Results**

31 Here, we present a new method which overcomes this distinct classification by
32 quantifying gene dispensability and present a dedicated tool for reference-based
33 QUantification Of gene Dispensability (QUOD). As a proof of concept, sequence data
34 of 966 *Arabidopsis thaliana* accessions (Ath-966) were processed to calculate a
35 gene-specific dispensability score for each gene based on normalised coverage in
36 read mappings. We validated this score by comparison of highly conserved
37 Benchmarking Universal Single Copy Orthologs (BUSCOs) to all other genes. The
38 average scores of BUSCOs were significantly lower than the scores of non-BUSCOs.
39 Analysis of variation demonstrated lower variation values between replicates of a
40 single accession than between iteratively, randomly selected accessions from the
41 whole dataset Ath-966. Functional investigations revealed defense and antimicrobial
42 response genes among the genes with high-dispensability scores.

43 **Conclusions**

44 Instead of classifying a gene as core or dispensable, QUOD assigns a dispensability
45 score to each gene. Hence, QUOD facilitates the identification of candidate
46 dispensable genes, associated with high dispensability scores, which often underlie
47 lineage-specific adaptation to varying environmental conditions.

48

49

50 **Keywords**

51 pangenomics, genomics, dispensability, bioinformatics, bioinformatic tool,
52 presence/absence variations

53

54

55 **Background**

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

67 Genes of a pangenome are thought to be divided into a core and a dispensable gene
68 set, the latter is also often referred to as ‘accessory’ in the literature. Core genes
69 occur in all investigated genomes, whereas dispensable genes only occur in a single
70 or a few genomes (8). In eukaryotic pangenome studies, core and dispensable genes
71 are mostly identified based on sequence similarity e.g. using GET_HOMOLOGUES-
72 EST Markov clustering (9), OrthoMCL gene family clustering (10) or BLASTN (11).
73 Sometimes, a third category of ‘conditionally dispensable’ genes is invoked (12) or
74 genes might be classified as ‘cloud’, ‘shell’, ‘soft-core’ and ‘core’ (13) or even as
75 ‘core’, ‘softcore’, ‘dispensable’ and ‘private’ (14). However, this distinct classification
76 is not based on the biological dispensability of genes and relies on one or multiple
77 arbitrary cutoffs. Some studies consider genes as ‘core’ if these genes occur in at
78 least 90 % of the investigated genomes (11); in other studies, only genes which are
79 found in all genomes are part of the core genome (10). In addition, dependency
80 groups might influence the dispensability of certain genes. The possibility that two
81 genes might be ‘replaced’ by a specific number of other genes has to be considered.
82 Some genes, of e.g. a gene family, might be required in a specific proportion and
83 therefore are only conditionally dispensable (12). Further, assemblies of genomes or
84 transcriptomes might be incomplete leading to artificially missing genes (15). One
85 way to circumvent this is to rely only on high-quality reference genome sequences,
86 thus avoiding additional assemblies which are potential sources of errors.

87 Here, we present QUOD - a bioinformatic tool to quantify gene dispensability. An *A.*
88 *thaliana* dataset of about 1,000 accessions was used to calculate a per gene
89 dispensability score derived from the coverage of all genes in the given genomes.
90 This score was validated by comparison of scores of BUSCOs and the functional
91 investigation of genes with high-dispensability scores. Our tool is easy to use for all

92kinds of plant species. QUOD extends the distinct classification of genes as ‘core’
93and ‘dispensable’ based on an arbitrary threshold to a continuous dispensability
94score.

95

96

97**Methods**

98**Selection and preprocessing of datasets**

99Genomic reads (FASTQ format) of the investigated genomes were retrieved from the
100Sequence Read Archive (SRA) (16) via fastq-dump. BWA-MEM (v.0.7.13) (17) was
101applied to map all genomic paired-end Illumina reads to the corresponding reference
102genome sequence using default parameters as well as *-m* to discard secondary
103alignments. For *A. thaliana*, all available 1,135 datasets (18) (Additional file 1) were
104subjected to a mapping against the AthNd-1_v2c genome sequence (19). The
105resulting BAM files of these mappings were subjected to QUOD.

106

107**Calculation of gene dispensability scores – QUOD**

108QUOD calculates a reference-based gene dispensability score for each structurally
109annotated gene based on supplied mapping files (BAM) (one per investigated
110genome) and a structural annotation of the reference sequence (GFF)
111(<https://github.com/ksielemann/QUOD>). The tool is written in Python3 and consists of
112six different components (Additional file 2). During the first part of the analysis, the
113read coverage per position (I) as well as the read coverage per gene (II) are
114calculated. In the next step, genomes with an average coverage below a given cutoff

115(default=10) are discarded and excluded from further analyses (III). Finally, an input
116matrix is constructed (IV) and a dispensability score is determined for each gene (V).
117QUOD assigns high gene dispensability scores to more likely dispensable genes.
118Optionally, the results can be visualized as a colored histogram and a box plot (VI).
119The dispensability score ($ds(g)$) is calculated as follows (cov.=coverage):

$$120 \text{ dispensability score (gene } g) = 1 / \left[\frac{\sum_{n=1}^N \left(\frac{\text{average cov. of gene } g \text{ in genome } n}{\text{average cov. over all genes in genome } n} \right)}{\text{total number of genomes (N)}} \right]$$

121

122 **Comprehension of the dispensability score composition**

123For further investigation of the score composition of selected genes of interest, the
124script `'score_composition.py'` can be used
125(https://github.com/ksielemann/QUOD/blob/master/score_composition.py). As output,
126a table including (I) the dispensability score, (II) the average coverage of all
127investigated genome sequences, (III) the average coverage of the accessions with
128the highest and (IV) lowest 10 % of all coverage values, respectively, (V) the number
129of accessions with zero coverage and (VI) the coverage for each accession,
130separately, is provided. Further, the coverage distribution for each gene can be
131visualized in a box plot.

132

133 **Identification of plastid sequences**

134Genes of Ath-966 with high similarity to plastid sequences were flagged via BLASTp
135(20) of the encoded peptides against all organelle peptide sequences obtained from
136the National Center for Biotechnology Information (NCBI). As a control, the

137sequences were also searched against themselves. Peptide sequences of Nd-1 with
138a score ratio ≥ 0.8 were considered plastid-like sequences when comparing BLAST
139hits against self-hits (19).

140

141**Score comparison between contrasting gene sets**

142Genes structurally annotated in AthNd-1_v2c were classified with BUSCO v3 (21)
143running in protein mode on the encoded peptide sequences using ‘brassicales
144odb10’ (order level) as reference (22). For comparison, BUSCO was additionally
145executed using ‘chlorophyta odb10’ (phylum level) and ‘embryophyta odb10’ (clade
146level) as reference. BUSCOs include single-copy genes and universal genes which
147are present in $> 90\%$ of all species in the reference dataset and are used to measure
148the completeness of assemblies and annotations (21). The scores of BUSCO and
149non-BUSCO genes were compared using matplotlib (23) for visualization (violin plot)
150and a Mann–Whitney U test implemented in the Python package dabest (24) for
151determination of the significance ([https://github.com/ksielemann/QUOD/blob/master/
152BUSCO_comparison.py](https://github.com/ksielemann/QUOD/blob/master/BUSCO_comparison.py)). Further, a Levene’s test, implemented in the Python
153package SciPy (25), was calculated to test for equal variances among BUSCO genes
154and non-BUSCO genes. The dispensability score of non-BUSCO genes might
155deviate more from the mean as non-BUSCO genes might be less conserved
156compared to BUSCO genes and might include multi-copy genes. Note that for all
157analyses performed within this study, the score of the size ‘infinity’ (detected for one
158gene) was set to the next highest score to enable calculations.

159A list of Nd-1 transposable element (TE) genes, which are Nd-1 gene structures
160overlapping with sequences annotated as TEs, was obtained from Pucker *et al.* (19).

161 First, the score distribution of TE and non-TE genes was determined using a Mann–
162 Whitney U test implemented in the Python package SciPy (25)
163 (https://github.com/ksielemann/QUOD/blob/master/analyse_TE_genes_and_scores.p
164 y). Next, the minimal distance of each gene to its closest TE gene was calculated
165 after extracting the gene positions from the Nd-1 annotation file. Mixed linear
166 modelling was performed using Statsmodels v0.12.0 (26) to determine the interaction
167 between the distance to the closest TE gene and the gene dispensability score
168 (https://github.com/ksielemann/QUOD/blob/master/mixed_linear_effects.py).

169

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

171 Length and number of exons per gene were extracted from the Nd-1 annotation file.
172 Linear mixed modelling was performed for gene length, exon number and the gene
173 dispensability score for the whole dataset Ath-966 as well as for three large *A.*
174 *thaliana* gene families (TAPscan (27)), namely MYBs (28), AP2/EREBP (29) and
175 WRKYs (30) using Statsmodels v0.12.0 (26)
176 (https://github.com/ksielemann/QUOD/blob/master/mixed_linear_effects.py).

177

178 **Variation between replicates**

179 A total of 14 genomic datasets of the *A. thaliana* accession Col-0 were received from
180 the SRA (Additional file 3) to assess the technical variation between replicates of the
181 same accession. Col-0 was selected for this analysis, because multiple independent
182 and high-quality datasets are only available for this accession. Each dataset was
183 mapped to the TAIR10 reference genome sequence using BWA-MEM because a
184 Col-0 read mapped against AthNd-1_v2c would result in multiple differences caused

185by accession-specific differences. The mappings were then subjected to QUOD,
186expecting a dispensability score close to one for each gene as there should be no
187variability between datasets of the same accession. As the distributions are different
188(Kolmogorov-Smirnov test, $p \approx 3e-27$) and the sample size (n) is high, the Levene's
189test was selected to test for equal variances, regarding the gene dispensability
190scores. The test was applied for (1) the dataset including replicates only and (2)
191iteratively (100x), randomly chosen subsets ($n=14$) of Ath-966
192(https://github.com/ksielemann/QUOD/blob/master/variance_in_repl_test.py).

193

194 **Functional annotation**

195All genes of the *A. thaliana* Nd-1 genome sequence were annotated via reciprocal
196best blast hits (RBHs) and best BLAST hits against Araport11 (19). Functional
197enrichment analyses (PANTHER protein classes and 'biological process' GO terms)
198were performed using the PANTHER Classification System of the Gene Ontology
199(31).

200

201 **Read mapper comparison**

202To evaluate the impact of the read mapping, the results of different mappers were
203compared. In addition to BWA-MEM (v.0.7.13; see above) (17), Bowtie2 (v2.4.1;
204default parameters) (32) and STAR (v2.5.1b) (33) were selected for this analysis.
205STAR parameters required alignments with a similarity of at least 95% over at least
20690% of the read pair length. The average coverage values per gene were
207investigated for correlation using the Spearman correlation coefficient implemented in
208the Python package SciPy (25).

209

210 **Data Availability**

211 The tool QUOD (QUOD.py) can be downloaded from GitHub
212 (<https://github.com/ksielemann/QUOD>; <http://doi.org/10.5281/zenodo.4066818>). A
213 data set to test QUOD is available on 'PUB - Publications at Bielefeld University'
214 (<http://doi.org/10.4119/unibi/2946079>).

215

216

217 **Results**

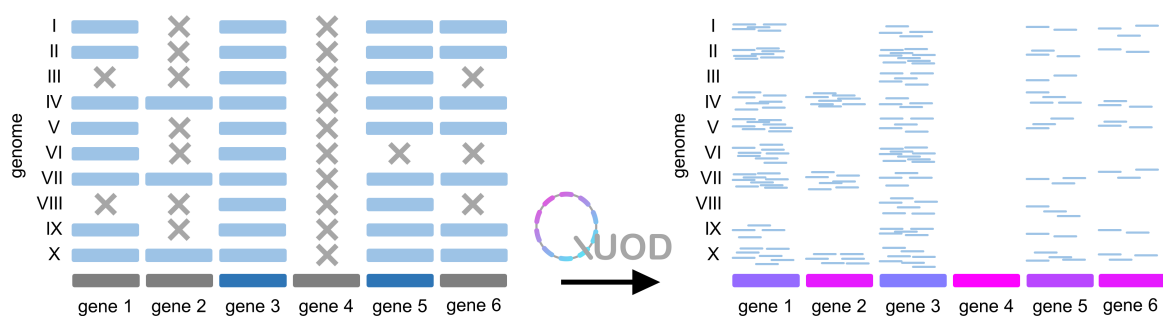
218 In this study, a bioinformatic tool was developed to calculate a gene-specific
219 dispensability score based on the normalised coverage in a read mapping. QUOD
220 allows the quantification of dispensability by calculation of a single score for each
221 gene (Figure 1). The binary classification of gene dispensability can be compared to
222 the original method of mRNA detection by endpoint RT-PCR providing only
223 qualitative results (34–36) which was replaced by quantitative analyses like RNA-
224 Seq.

225

226 **Gene dispensability scores**

227 The gene dispensability score would initially be dependent on the sequencing depth
228 per genome. By division of the average coverage of gene g in genome n (N = total
229 number of investigated genome sequences) by the average coverage over all genes
230 in genome n , the score is normalised for differences in the sequencing depth of the
231 investigated genomes. A high value indicates that a gene is likely to be missing in

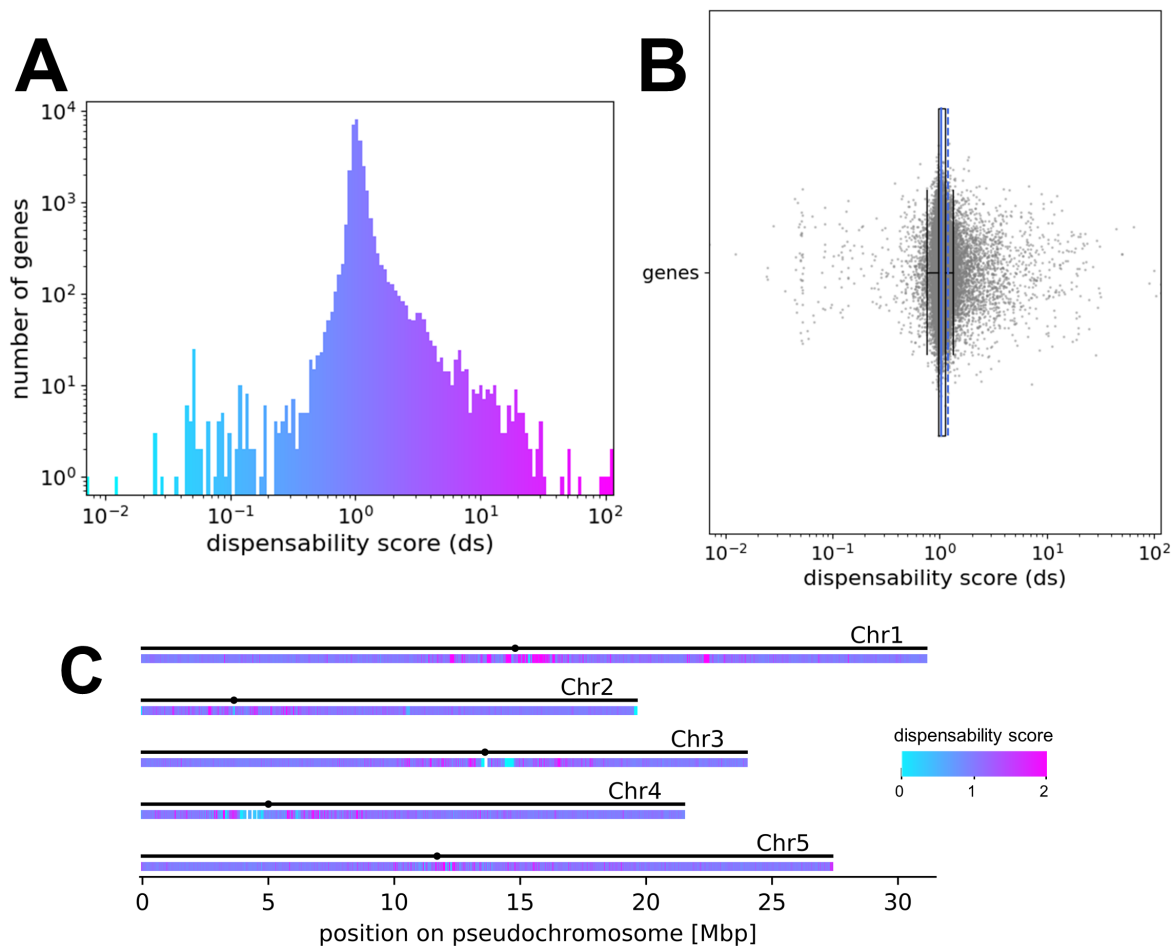
232some genomes and therefore more likely dispensable than a gene with a lower
233dispensability score. Due to this quantification approach, this method is not based on
234an arbitrary cutoff to determine the core genome and the dispensable genome of any
235given pangenome dataset. An example: Using a cutoff of 'gene n occurs in at least
23690 % of all genomes' to be considered a 'core' gene (dark blue), genes 1,2,4 and 6
237(dark grey) would be considered 'dispensable' (Figure 1). However, considering the
238coverage (right panel), it is not clear if e.g. gene 1 is truly biologically dispensable.
239QUOD does not rely on any thresholds for the classification of genes into 'core' and
240'dispensable', but provides a score based on the normalised coverage in a read
241mapping. The genes could theoretically be ranked as well using the percentage of
242presence/absence of a gene in the investigated genomes. However, this alternative
243approach would still rely on a threshold, e.g. the number of mapped reads for a gene
244to be considered present in a genome. This threshold is avoided using the QUOD
245method.



247Figure 1: Illustration of the QUOD method using a fictional dataset. On the left side,
248genes are classified as 'core' (dark blue) or 'dispensable' (dark grey) according to a
249cutoff. On the right side, gene dispensability is quantified according to a
250dispensability score based on the normalised coverage in a read mapping (I-X:

251investigated genomes). Coloring of genes (right side) indicates different
252dispensability scores. Extremely rare genes, which are absent from most genomes
253but present in the reference, can be easily detected using QUOD.

254As a proof of concept, *A. thaliana* sequence reads of 1,135 accessions were mapped
255to the *A. thaliana* Nd-1 genome sequence. All accessions with less than 10-fold read
256coverage were discarded. The remaining sequencing dataset Ath-966 was analysed
257with QUOD to calculate a dispensability score for each gene (Figure 2). Genes with
258high dispensability scores, colored in pink, are considered to be likely dispensable,
259whereas genes with dispensability scores close to one (dark purple/dark blue) are
260considered to be core genes.



262 Figure 2: Distribution of the gene dispensability scores for Ath-966. A) Histogram
263 coloured according to the dispensability score. The x-axis represents the
264 dispensability score and the y-axis shows the number of genes in each bin in
265 logarithmic scale. B) Box plot representing the dispensability score (x-axis) of all
266 genes (y-axis). The mean is represented by the dashed blue line, the other blue line
267 represents the median of the scores. C) Genome-wide distribution of genes with
268 different dispensability scores in *A. thaliana* Nd-1. The coloured heatmap shows the
269 respective gene dispensability scores. There are low (blue) and high (pink) scoring
270 genes clustered in repetitive regions, including centromeric and telomeric areas. The
271 x-axis represents the size (in Mbp) of each pseudochromosome in the assembly. The
272 black dots represent the position of the centromeres of the five chromosomes in the
273 AthNd1_v2c assembly (19).

274

275 **Genome-wide distribution of the gene dispensability scores**

276 Next, the genome-wide distribution of genes with specific gene dispensability scores
277 was investigated in *A. thaliana* (Figure 2C). A high plasticity between accessions,
278 which means a high number of genes with exceptionally high and low scores (pink
279 and blue), in the (peri-)centromeric regions is visible based on a heatmap (Figure
280 2C).

281 As high and low scoring genes cluster in repetitive regions (mainly centromeres), the
282 score distribution of TE genes was investigated (Additional file 4). Scores of TE
283 genes are evenly distributed across all dispensability scores. In total, the mean score
284 of TE genes (mean ds \approx 1.501) is significantly higher when compared to non-TE
285 genes (mean ds \approx 1.168) (Mann-Whitney U test, $p \approx 6E-8$), which are more frequent

286across scores close to one. Moreover, the minimal distance of each gene to its
287closest TE gene and the dispensability scores revealed no relation (Additional file 4).

288To test the hypothesis whether genes with higher dispensability scores/more likely
289dispensable genes are shorter and whether introns accumulate in core genes, the
290correlation of the gene dispensability score with gene length and exon number,
291respectively, were determined for the Ath-966 and for three selected gene families
292separately. However, no clear trend was detectable (Additional file 5).

293

294**Validation of the reliability**

295Validation of the reliability of the gene dispensability quantification was achieved by
296comparison of BUSCOs and non-BUSCOs (Additional file 6). BUSCO genes show on
297average slightly lower scores than non-BUSCO genes for all three reference datasets
298($p < 0.001$, Mann-Whitney U test). Levene's test was used to test for equal variances.
299The results show that the variances for all reference datasets differ significantly
300between BUSCO and non-BUSCO genes ($p < 0.001$, Levene's test). Thus, the
301deviation of the dispensability score from the respective mean is significantly higher
302for non-BUSCO genes in comparison to BUSCO genes.

303Further, functional annotation of BUSCO outliers, which are genes of the 'brassicales
304odb10' BUSCO gene set with dispensability scores below 0.75 or above 1.25,
305revealed, amongst others, several repeat proteins, transmembrane proteins, a 'stress
306induced protein', and multiple hypothetical proteins (Additional file 7).

307Genes with high and low gene dispensability scores were assessed in more detail.
308Among genes with high dispensability scores, several significantly enriched
309PANTHER protein classes were detected, e.g. defense/immunity and antimicrobial

310response proteins, small GTPases and G-proteins (Table 1). Among genes with
 311dispensability scores < 0.8, genes encoding proteins of the extracellular matrix were
 312significantly enriched (Table 1). 'Biological process' GO term enrichment revealed
 313several significantly enriched terms associated with the regulation of cellular
 314processes as well as associated with response to stimuli among genes with
 315dispensability scores > 2 (Table 1). Genes with low dispensability scores show
 316enrichment of primary metabolic processes (Table 1).

317

318Table 1: Closer investigation of genes with scores >2 and genes with scores < 0.8.
 319Significantly enriched PANTHER protein classes (padj < 0.05) as well as significantly
 320enriched GO biological process terms (padj < 0.05) are shown. Abbreviations: p =
 321process, mp = metabolic process.

PANTHER protein classes (padj < 0.05) of genes with scores >2	
small GTPase (PC00208)	4.21E-05
defense/immunity protein (PC00090)	4.24E-05
antimicrobial response protein (PC00051)	5.24E-05
G-protein (PC00020)	4.05E-04
protein class (PC00000)	2.04E-03
Unclassified	2.44E-03
protein-binding activity modulator (PC00095)	3.72E-02
PANTHER protein classes (padj < 0.05) of genes with scores <0.8	
extracellular matrix structural protein (PC00103)	5.40E-06
extracellular matrix protein (PC00102)	1.14E-05
Unclassified	2.68E-05
protein class (PC00000)	3.57E-05
metabolite interconversion enzyme (PC00262)	3.04E-02
GO biological process terms (padj < 0.05) of genes with scores >2	
cellular p (GO:0009987)	2.62E-08
mp (GO:0008152)	4.62E-07
cellular mp (GO:0044237)	2.85E-06
primary mp (GO:0044238)	2.37E-05
organic substance mp (GO:0071704)	3.02E-05
regulation of cellular mp (GO:0031323)	9.82E-04

regulation of biosynthetic p (GO:0009889)	9.92E-04
regulation of cellular biosynthetic p (GO:0031326)	1.04E-03
regulation of cellular macromolecule biosynthetic p (GO:2000112)	2.27E-03
regulation of macromolecule biosynthetic p (GO:0010556)	2.52E-03
regulation of primary mp (GO:0080090)	2.90E-03
macromolecule mp (GO:0043170)	2.93E-03
regulation of nitrogen compound mp (GO:0051171)	4.53E-03
regulation of RNA mp (GO:0051252)	4.69E-03
positive regulation of biological p (GO:0048518)	4.89E-03
response to organic substance (GO:0010033)	4.91E-03
positive regulation of cellular p (GO:0048522)	6.62E-03
regulation of RNA biosynthetic p (GO:2001141)	6.67E-03
regulation of mp (GO:0019222)	6.72E-03
regulation of nucleobase-containing compound mp (GO:0019219)	6.74E-03
regulation of nucleic acid-templated transcription (GO:1903506)	6.95E-03
developmental p (GO:0032502)	7.01E-03
response to hormone (GO:0009725)	7.25E-03
regulation of transcription, DNA-templated (GO:0006355)	7.27E-03
response to oxygen-containing compound (GO:1901700)	7.53E-03
anatomical structure development (GO:0048856)	7.62E-03
nitrogen compound mp (GO:0006807)	1.25E-02
response to endogenous stimulus (GO:0009719)	1.48E-02
regulation of gene expression (GO:0010468)	2.91E-02
system development (GO:0048731)	3.44E-02
regulation of macromolecule mp (GO:0060255)	3.45E-02
cellular lipid mp (GO:0044255)	4.10E-02
clathrin coat disassembly (GO:0072318)	4.14E-02
multicellular organismal p (GO:0032501)	4.19E-02
vesicle uncoating (GO:0072319)	4.26E-02
GO biological process terms (padj < 0.05) of genes with scores < 0.8	
cellular p (GO:0009987)	6.35E-07
mp (GO:0008152)	1.35E-06
organic substance mp (GO:0071704)	8.49E-06
cellular mp (GO:0044237)	2.92E-05
nitrogen compound mp (GO:0006807)	5.35E-04
primary mp (GO:0044238)	5.76E-04
macromolecule mp (GO:0043170)	3.82E-03
organonitrogen compound mp (GO:1901564)	9.67E-03
localization (GO:0051179)	4.87E-02

322

323The function of the 100 genes with the highest gene dispensability scores was
324examined in detail for Ath-966 (Additional file 8). Fourteen genes of Ath-966 are
325annotated as “disease resistance proteins”, whereas seven genes are annotated as
326transposons/transposases. Four genes are described as hypothetical proteins and 24
327genes have no functional annotation. In addition, an example for lineage specific
328adaptation is provided (Additional file 9). The gene NdCChr1.g3308 has a
329dispensability score of approx. 10. For 870 accessions, which account for approx. 90
330% of Ath-966, no coverage was detected. The gene is annotated as resistance gene
331mediating resistance against the bacterial pathogen *Pseudomonas syringae*.

332Next, the variation between replicates of the same accession (Col-0) was determined
333(Additional file 10). The variation of the gene dispensability score distribution of the
334replicate dataset (one accession) ($\sigma^2 \approx 0.0226$) is significantly lower than the variation
335between all iteratively, randomly selected subsets of *A. thaliana* accessions ($\sigma^2 \approx$
3360.0392) (Levene’s test, $p \approx 4e-19$). The average coverage per gene using different
337read mappers revealed strong correlations in all comparisons (Additional file 11). The
338coverage correlations, calculated using Spearman correlation coefficient, between
339BWA-MEM and bowtie2 ($r \approx 0.810$, $p \approx 0.0$), BWA-MEM and STAR ($r \approx 0.814$, $p \approx$
3400.0) as well as bowtie2 and STAR ($r \approx 0.760$, $p \approx 0.0$) are similar.

341

342

343Discussion

344QUOD was developed for the quantification of gene dispensability in plant
345pangenome datasets. Multiple accessions of several plant species have been
346sequenced and pose potential use cases for QUOD (Additional file 12). Dropping

347sequencing costs will lead to an increasing availability of comprehensive sequence
348datasets which would permit the application of QUOD. Additionally, QUOD is not
349restricted to plants, but could be applied to other species (e.g. pig (37)). However,
350an accurate determination of gene dispensability scores free of systematic biases
351might rely on a uniform selection of genomes from the respective taxonomic group
352and on uniform read coverage of genes. In addition, non-random fragmentation of
353DNA prior to sequencing (38) may cause biases. The variation among replicates of
354the same accession (Col-0; $\sigma^2 \approx 0.0226$) might be attributed to technical biases, e.g.
355during sequencing library preparation. The comparison of different read mappers
356revealed a significant correlation for the average coverage per gene. Outlier samples,
357detected by the investigation of the average coverage per gene using different read
358mappers, might indicate technical issues. Even though the correlations are strong,
359the same tool with the same parameter settings needs to be used for the read
360mapping of all compared datasets within one single QUOD run.

361Most genes show dispensability scores close to one as the majority of genes are
362widespread across species. The aim of QUOD is mainly the identification of the
363‘outliers’ and therefore the more dispensable genes, which are genes not present in
364all genomes. These dispensable genes represent a smaller fraction of the genome
365than the core genes. Genome level patterns are expected to be similar for all
366species. Further, QUOD is not an alternative to PAV detection methods as groups of
367genes can still always be defined using PAV methods, but QUOD provides a
368quantitative measurement for these cases.

369As already stated in the Introduction, genome assemblies might be incomplete
370leading to artificially missing genes (15). One way to circumvent this is to rely on a
371high-quality reference genome sequence, thus avoiding additional assemblies which

372 are potential sources of errors. Recently released telomere-to-telomere assemblies
373 indicate that these resources will be available for many plant species in the near
374 future (39). Further, the usage of QUOD with a synthetic reference derived from
375 multiple assemblies is possible and can be implemented in the future. A graph-based
376 assembly of a pangenome comprising multiple accessions is already feasible for
377 bacteria (40–42). However, for large plant genome sequences graph-based
378 pangenome assembly is computationally expensive and not yet robust for complex
379 structural variants like inversions(43). Even though there are still several
380 shortcomings, like loss of the sample information (44), improved methods might be
381 available in the near future and could be used for the improved quantification of gene
382 dispensability.

383

384 **Genome-wide distribution of the gene dispensability scores**

385 The genome-wide distribution of all gene dispensability scores (not only BUSCO
386 genes) of the *A. thaliana* genomes reveals the origin of exceptionally low
387 dispensability scores (Figure 2). Low scoring genes, which are colored in light blue
388 in Figure 2, might be TEs and other repeat genes associated with collapsed
389 sequences in the assembly. An accurate determination of the dispensability scores of
390 these genes might be possible using ideal genome sequences without any collapsed
391 regions and with specific read mappings e.g. using high quality long reads. However,
392 low scoring genes could still be useful to determine amplified TEs and other repeat
393 genes. Moreover, the genome-wide distribution plot (Figure 2C) shows that high and
394 low scoring genes cluster in repetitive regions, like centromeres or telomeres. Very
395 similar sequences, e.g. members of a gene family or close paralogs, might cause
396 read mapping errors confounding biases in the dispensability scores of these genes.

397 Additionally, this can be explained by variation in the recombination rate (45) and
398 active TEs in these regions. It was previously proposed, that dispensable genes are
399 likely located closer to TEs which are important factors in genome evolution (9).
400 However, in the results of our study, TE genes are widely distributed across all
401 dispensability scores as TEs can occur with variable copy numbers in genomes
402 leading to low scores and can as well be dispensable. Other studies detected a high
403 number of TEs in the dispensable genome (46). However, it is possible that only
404 certain TE families might be truly dispensable. One limitation is the accurate
405 assignment of reads to repetitive sections of the reference sequence during the read
406 mapping (15). Further, only a fraction of transposons might be correctly assembled
407 and annotated due to several computational challenges in highly repetitive and peri-
408 centromeric regions (47). Therefore, a different strategy might be needed to
409 accurately quantify dispensability of TEs. A high quality annotation of transposons
410 and a following exclusion of these genes from the analysis or improved read mapping
411 to the consensus sequence might improve the results. Again, long reads could be an
412 alternative solution to handle regions which might be ambiguous in read mappings.
413 Moreover, heterochromatin or genome-purging mechanisms (48) could influence the
414 gene dispensability scores in these regions.

415 Additionally some of the low scoring genes were flagged as plastid-like sequences as
416 original sequencing data from plants contain high amounts of reads originating from
417 plastid sequences (49,50). Biases due to this plastid read contamination inflate the
418 coverage of sequences with high similarity to plastid sequences, resulting in an
419 exceptionally low gene dispensability score.

420

421 **Validation of the reliability**

422 We validated the reliability of the gene dispensability score by showing that more
423 conserved BUSCO genes get significantly lower dispensability scores than non-
424 BUSCO genes (Additional file 6). Based on the distribution of the scores in the violin
425 plot (Additional file 6), the difference between BUSCOs and non-BUSCOs appears
426 small, even though the difference is significant (U test, $p \approx 4E-113$, brassicales
427 reference). It is important to note that non-BUSCO genes can be highly conserved.
428 Consequently, the difference is only visible at the group level. The difference in the
429 dispensability scores of BUSCOs and non-BUSCOs is low as expected, because
430 conserved multiple-copy genes are not included in the BUSCO gene set (21).
431 Therefore, the variance of the dispensability scores of non-BUSCO genes is
432 significantly larger than the variance among BUSCO genes: non-BUSCO genes
433 comprise highly conserved multi-copy genes as well as less conserved genes.
434 Further, functional annotation of BUSCO outliers revealed several repeat proteins
435 and transmembrane proteins. Repeat proteins might lead to read mapping errors and
436 consequently artificial variations in coverage and dispensability scores.
437 Transmembrane proteins are thought to be involved in biotic stress response and
438 might not be essential for some accessions and therefore dispensable (51). This
439 could explain the absence in some genomes resulting in high dispensability scores of
440 these genes. Therefore, many important, lower-scoring genes might lie outside of the
441 BUSCO reference set.

442 Functional annotation of the 100 most likely dispensable genes revealed a high
443 number of uncharacterised proteins, disease resistance proteins as well as
444 transposons and transposases in the *A. thaliana* genomes. It is possible that these
445 genes are undergoing pseudogenization and have not been functionally annotated
446 due to the lack of a visible phenotype when mutated. TEs were detected in other

447 studies as contributors to large structural variations between species and individuals
448 and considered as a substantial part of the dispensable genome (46). Previous
449 pangenome analyses also revealed that the dispensable genome comprises
450 functions like ‘defense response’, ‘diseases resistance’, ‘flowering time’ and
451 ‘adaptation to biotic and abiotic stress’ (9,11,13). Comparable results were detected
452 for the enriched protein classes and ‘biological process’ GO terms (Table 1), even
453 though very general terms, like ‘protein class’, give little evidence about the function
454 of genes. Moreover, we provide a specific example for lineage specific adaptation
455 associated with a high dispensability score (Additional file 9): a gene mediating
456 resistance against the bacterial pathogen *Pseudomonas syringae*. Therefore, in
457 depth investigation of genes with high dispensability scores can result in the
458 identification and characterization of phenotypic variation (52) and important
459 agronomic traits (13). We envision several applications for the gene dispensability
460 score generated by QUOD: (1) more accurate prediction if a gene is associated with
461 a specific trait, (2) development of dependency gene networks, and (3) improved
462 modeling of the evolutionary value of genes.

463

464

465 **Conclusions**

466 QUOD (reference-based QUantification Of gene Dispensability) overcomes the
467 problem of labeling genes as ‘core’ or ‘dispensable’ through implementation of a
468 quantification approach. Instead of a distinct classification, QUOD provides a ranking
469 of all genes based on assigned gene-specific dispensability scores and therefore
470 does not rely on any thresholds.

471

472

473 **Declarations**

474 **Ethics approval and consent to participate**

475 Not applicable.

476 **Consent for publication**

477 Not applicable.

478 **Availability of data and materials**

479 The tool QUOD for the reference-based QUantification Of gene Dispensability

480 (QUOD.py) can be downloaded from GitHub (<https://github.com/ksielemann/QUOD>;

481 <http://doi.org/10.5281/zenodo.4066818>).

482 **Competing interests**

483 The authors declare that they have no competing interests.

484 **Funding**

485 KS is funded by Bielefeld University.

486 **Authors' contributions**

487 KS, BW and BP designed the study, performed the experiments, analysed the data,

488 and wrote the manuscript. All authors read and approved the final version of this

489 manuscript.

490 **Acknowledgements**

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

497

498

499 **References**

5001. Springer NM, Ying K, Fu Y, Ji T, Yeh C-T, Jia Y, et al. Maize Inbreds Exhibit
501 High Levels of Copy Number Variation (CNV) and Presence/Absence Variation
502 (PAV) in Genome Content. Ecker JR, editor. PLoS Genet. 2009 Nov
503 20;5(11):e1000734.
5042. Scherer SW, Lee C, Birney E, Altshuler DM, Eichler EE, Carter NP, et al.
505 Challenges and standards in integrating surveys of structural variation. Nat
506 Genet. 2007 Jul;39(S7):S7–15.
5073. Tao Y, Zhao X, Mace E, Henry R, Jordan D. Exploring and Exploiting Pan-
508 genomics for Crop Improvement. Molecular Plant. 2019 Feb;12(2):156–69.
5094. Lu F, Romay MC, Glaubitz JC, Bradbury PJ, Elshire RJ, Wang T, et al. High-
510 resolution genetic mapping of maize pan-genome sequence anchors. Nat
511 Commun. 2015 Nov;6(1):6914.

5125. Swanson-Wagner RA, Eichten SR, Kumari S, Tiffin P, Stein JC, Ware D, et al.
513 Pervasive gene content variation and copy number variation in maize and its
514 undomesticated progenitor. *Genome Research*. 2010 Dec 1;20(12):1689–99.
5156. Tettelin H, Massignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, et al.
516 Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*:
517 Implications for the microbial ‘pan-genome’. *Proceedings of the National*
518 *Academy of Sciences*. 2005 Sep 27;102(39):13950–5.
5197. Vernikos G, Medini D, Riley DR, Tettelin H. Ten years of pan-genome analyses.
520 *Current Opinion in Microbiology*. 2015 Feb;23:148–54.
5218. Golicz AA, Batley J, Edwards D. Towards plant pangenomics. *Plant Biotechnol*
522 *J*. 2016 Apr;14(4):1099–105.
5239. Gordon SP, Contreras-Moreira B, Woods DP, Des Marais DL, Burgess D, Shu
524 S, et al. Extensive gene content variation in the *Brachypodium distachyon* pan-
525 genome correlates with population structure. *Nat Commun*. 2017 Dec;8(1):2184.
52610. Li Y, Zhou G, Ma J, Jiang W, Jin L, Zhang Z, et al. De novo assembly of
527 soybean wild relatives for pan-genome analysis of diversity and agronomic traits.
528 *Nat Biotechnol*. 2014 Oct;32(10):1045–52.
52911. Zhao Q, Feng Q, Lu H, Li Y, Wang A, Tian Q, et al. Pan-genome analysis
530 highlights the extent of genomic variation in cultivated and wild rice. *Nat Genet*.
531 2018 Feb;50(2):278–84.
53212. Marroni F, Pinosio S, Morgante M. Structural variation and genome complexity:
533 is dispensable really dispensable? *Current Opinion in Plant Biology*. 2014
534 Apr;18:31–6.

53513. Golicz AA, Bayer PE, Barker GC, Edger PP, Kim H, Martinez PA, et al. The
536 pangenome of an agronomically important crop plant *Brassica oleracea*. *Nat*
537 *Commun.* 2016 Dec;7(1):13390.
53814. Liu Y, Du H, Li P, Shen Y, Peng H, Liu S, et al. Pan-Genome of Wild and
539 Cultivated Soybeans. *Cell.* 2020 Jul;182(1):162-176.e13.
54015. Alkan C, Sajjadian S, Eichler EE. Limitations of next-generation genome
541 sequence assembly. *Nat Methods.* 2011 Jan;8(1):61–5.
54216. Leinonen, Rasko and Sugawara, Hideaki and Shumway, Martin and
543 International Nucleotide Sequence Database Collaboration. The sequence read
544 archive. *Nucleic Acids Research.* 2010;39(suppl_1):D19–D21.
54517. Li H. Aligning sequence reads, clone sequences and assembly contigs with
546 BWA-MEM. arXiv preprint arXiv:13033997. 2013;
54718. Alonso-Blanco C, Andrade J, Becker C, Bemm F, Bergelson J, Borgwardt KM, et
548 al. 1,135 Genomes Reveal the Global Pattern of Polymorphism in *Arabidopsis*
549 *thaliana*. *Cell.* 2016 Jul;166(2):481–91.
55019. Pucker B, Holtgräwe D, Stadermann KB, Frey K, Huettel B, Reinhardt R, et al. A
551 chromosome-level sequence assembly reveals the structure of the *Arabidopsis*
552 *thaliana* Nd-1 genome and its gene set. Feltus FA, editor. *PLoS ONE.* 2019 May
553 21;14(5):e0216233.
55420. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment
555 search tool. *Journal of Molecular Biology.* 1990 Oct;215(3):403–10.

55621. Simão, Felipe A and Waterhouse, Robert M and Ioannidis, Panagiotis and
557 Kriventseva, Evgenia V and Zdobnov, Evgeny M. BUSCO: assessing genome
558 assembly and annotation completeness with single-copy orthologs.
559 *Bioinformatics*. 2015;31(19):3210–3212.
56022. Zdobnov EM, Tegenfeldt F, Kuznetsov D, Waterhouse RM, Simão FA, Ioannidis
561 P, et al. OrthoDB v9.1: cataloging evolutionary and functional annotations for
562 animal, fungal, plant, archaeal, bacterial and viral orthologs. *Nucleic Acids Res*.
563 2017 Jan 4;45(D1):D744–9.
56423. Hunter JD. Matplotlib: A 2D Graphics Environment. *Comput Sci Eng*.
565 2007;9(3):90–5.
56624. Ho J, Tumkaya T, Aryal S, Choi H, Claridge-Chang A. Moving beyond P values:
567 data analysis with estimation graphics. *Nat Methods*. 2019 Jul;16(7):565–6.
56825. Jones E, Oliphant T, Peterson P, others. SciPy: Open source scientific tools for
569 Python. 2001;
57026. Seabold S, Perktold J. Statsmodels: Econometric and statistical modeling with
571 python. *Proceedings of the 9th Python in Science Conference*. 2010;57.
57227. Wilhelmsson PKI, Mühlich C, Ullrich KK, Rensing SA. Comprehensive Genome-
573 Wide Classification Reveals That Many Plant-Specific Transcription Factors
574 Evolved in Streptophyte Algae. *Genome Biology and Evolution*. 2017 Dec
575 1;9(12):3384–97.
57628. Stracke R, Werber M, Weisshaar B. The R2R3-MYB gene family in *Arabidopsis*
577 *thaliana*. *Current Opinion in Plant Biology*. 2001 Oct;4(5):447–56.

57829. Feng J-X, Liu D, Pan Y, Gong W, Ma L-G, Luo J-C, et al. An annotation update
579 via cDNA sequence analysis and comprehensive profiling of developmental,
580 hormonal or environmental responsiveness of the Arabidopsis AP2/EREBP
581 transcription factor gene family. *Plant Molecular Biology*. 2005 Dec;59(6):853–
582 68.
58330. Eulgem T, Rushton PJ, Robatzek S, Somssich IE. The WRKY superfamily of
584 plant transcription factors. *Trends in Plant Science*. 2000 May;5(5):199–206.
58531. Carbon, Seth, Mungall, Chris. Gene Ontology Data Archive [Internet]. Zenodo;
586 2018 [cited 2020 Sep 29]. Available from: <https://zenodo.org/record/3980761>
58732. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat*
588 *Methods*. 2012 Apr;9(4):357–9.
58933. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR:
590 ultrafast universal RNA-seq aligner. *Bioinformatics*. 2013;29(1):15–21.
59134. Wang AM, Doyle MV, Mark DF. Quantitation of mRNA by the polymerase chain
592 reaction. *Proceedings of the National Academy of Sciences*. 1989 Dec
593 1;86(24):9717–21.
59435. Gilliland G, Perrin S, Blanchard K, Bunn HF. Analysis of cytokine mRNA and
595 DNA: detection and quantitation by competitive polymerase chain reaction.
596 *Proceedings of the National Academy of Sciences*. 1990;87(7):2725–9.
59736. Chiang PW, Song WJ, Wu KY, Korenberg JR, Fogel EJ, Van Keuren ML, et al.
598 Use of a fluorescent-PCR reaction to detect genomic sequence copy number
599 and transcriptional abundance. *Genome Research*. 1996 Oct 1;6(10):1013–26.

60037. Tian X, Li R, Fu W, Li Y, Wang X, Li M, et al. Building a sequence map of the pig
601 pan-genome from multiple de novo assemblies and Hi-C data. *Sci China Life*
602 *Sci.* 2020 May;63(5):750–63.
60338. Poptsova MS, Il'icheva IA, Nechipurenko DYU, Panchenko LA, Khodikov MV,
604 Oparina NY, et al. Non-random DNA fragmentation in next-generation
605 sequencing. *Sci Rep.* 2015 May;4(1):4532.
60639. Miga KH, Koren S, Rhie A, Vollger MR, Gershman A, Bzikadze A, et al.
607 Telomere-to-telomere assembly of a complete human X chromosome. *Nature.*
608 2020 Sep 3;585(7823):79–84.
60940. Gautreau G, Bazin A, Gachet M, Planel R, Burlot L, Dubois M, et al.
610 PPanGGOLiN: Depicting microbial diversity via a partitioned pangenome graph.
611 Ouzounis CA, editor. *PLoS Comput Biol.* 2020 Mar 19;16(3):e1007732.
61241. Ding W, Baumdicker F, Neher RA. panX: pan-genome analysis and exploration.
613 *Nucleic Acids Research.* 2018 Jan 9;46(1):e5–e5.
61442. Chan AP, Sutton G, DePew J, Krishnakumar R, Choi Y, Huang X-Z, et al. A
615 novel method of consensus pan-chromosome assembly and large-scale
616 comparative analysis reveal the highly flexible pan-genome of *Acinetobacter*
617 *baumannii*. *Genome Biol.* 2015 Dec;16(1):143.
61843. Chen N-C, Solomon B, Mun T, Iyer S, Langmead B. Reducing reference bias
619 using multiple population reference genomes [Internet]. *Bioinformatics*; 2020
620 Mar [cited 2020 Sep 29]. Available from:
621 <http://biorxiv.org/lookup/doi/10.1101/2020.03.03.975219>

62244. Li H, Feng X, Chu C. The design and construction of reference pangenome
623 graphs. arXiv:200306079 [q-bio] [Internet]. 2020 Mar 12 [cited 2020 Sep 29];
624 Available from: <http://arxiv.org/abs/2003.06079>
62545. Nachman M. Variation in recombination rate across the genome: evidence and
626 implications. *Current Opinion in Genetics & Development*. 2002 Dec
627 1;12(6):657–63.
62846. Morgante M, Depaoli E, Radovic S. Transposable elements and the plant pan-
629 genomes. *Current Opinion in Plant Biology*. 2007 Apr;10(2):149–55.
63047. Platt RN, Blanco-Berdugo L, Ray DA. Accurate Transposable Element
631 Annotation Is Vital When Analyzing New Genome Assemblies. *Genome Biol
632 Evol*. 2016 Feb;8(2):403–10.
63348. Lee S-I, Kim N-S. Transposable Elements and Genome Size Variations in
634 Plants. *Genomics Inform*. 2014;12(3):87.
63549. Yu J, Hu S, Wang J, Wong GK-S, Li S, Liu B, et al. A draft sequence of the rice
636 genome (*Oryza sativa* L. ssp. *indica*). *Science*. 2002 Apr 5;296(5565):79–92.
63750. Tang J, Xia H, Cao M, Zhang X, Zeng W, Hu S, et al. A Comparison of Rice
638 Chloroplast Genomes. *Plant Physiol*. 2004 May;135(1):412–20.
63951. Dodds PN, Rathjen JP. Plant immunity: towards an integrated view of plant–
640 pathogen interactions. *Nat Rev Genet*. 2010 Aug;11(8):539–48.
64152. Hirsch CN, Foerster JM, Johnson JM, Sekhon RS, Muttoni G, Vaillancourt B, et
642 al. Insights into the Maize Pan-Genome and Pan-Transcriptome. *Plant Cell*.
643 2014 Jan;26(1):121–35.

644

645

646

647 **Additional files**

648 Additional file 1 (.tsv): SRA IDs of datasets downloaded to conduct the QUOD
649 analysis of the *A. thaliana* genomes.

650 Additional file 2 (.pdf): Illustration of the different components of QUOD.

651 Additional file 3 (.tsv): SRA/ENA IDs of datasets downloaded to conduct the analysis
652 of replicates (Col-0).

653 Additional file 4 (.pdf): Distribution of scores of TE genes and non-TE genes and
654 correlation of the distance to the closest TE gene with the gene dispensability score
655 of the *A. thaliana* genomes.

656 Additional file 5 (.pdf): Correlation of gene length and exon number with the
657 dispensability scores of the *A. thaliana* genomes.

658 Additional file 6 (.pdf): Comparison of BUSCO analyses for 'chlorophyta',
659 'brassicales' and 'embryophyta' as reference.

660 Additional file 7 (.tsv): Functional annotation of BUSCO outliers (using 'brassicales
661 odb10' as reference) with a dispensability score smaller than 0.75 or greater than
662 1.25.

663 Additional file 8 (.tsv): Functional annotation of the 100 most likely dispensable genes
664 of the *A. thaliana* genomes.

665 Additional file 9 (.pdf): Example for lineage specific adaptation.

666Additional file 10 (.pdf): Analysis of variance of the gene dispensability score
667calculated for replicates of the *A. thaliana* Col-0 accession and iteratively, randomly
668chosen subsets of the whole dataset Ath-966.

669Additional file 11 (.pdf): Correlation of the average coverage per gene using three
670different read mappers: BWA-MEM, bowtie2 and STAR.

671Additional file 12 (.tsv): Examples of diploid species where multiple cultivars were
672already sequenced.