

1 **Genome-wide association study in the pseudocereal quinoa**
2 **reveals selection pattern typical for crops with a short**
3 **breeding history**

4 Dilan S. R. Patiranage¹, Elodie Rey², Nazgol Emrani^{1,*}, Gordon Wellman², Karl Schmid³, Sandra
5 M. Schmöckel⁴, Mark Tester² and Christian Jung^{1,*}

6 ¹Plant Breeding Institute, Christian-Albrechts-University of Kiel, Am Botanischen Garten 1-9,
7 24118 Kiel, Germany

8 ²King Abdullah University of Science and Technology (KAUST), Biological and Environmental
9 Sciences & Engineering Division (BESE), Thuwal, 23955-6900, Saudi Arabia.

10 ³Institute of Plant Breeding, Seed Science and Population Genetics, University of Hohenheim,
11 Fruwirthstr. 21, 70599 Stuttgart, Germany.

12 ⁴Department of Physiology of Yield Stability, University of Hohenheim, Fruwirthstr. 21, 70599

13

14 *Corresponding authors email: c.jung@plantbreeding.uni-kiel.de

15 *Corresponding authors email: n.emrani@plantbreeding.uni-kiel.de

16 Correspondence to: Dr. Nazgol Emrani
17 Plant Breeding Institute
18 Christian-Albrechts-University of Kiel
19 Olshausenstrasse 40
20 D-24118 Kiel
21 Germany
22 Tel.: +49-4318807364
23 Fax: +49-4318802566

24 Email: n.emrani@plantbreeding.uni-kiel.de

25 Correspondence to: Prof. Dr. Christian Jung
26 Plant Breeding Institute
27 Christian-Albrechts-University of Kiel
28 Olshausenstrasse 40
29 D-24118 Kiel
30 Germany
31 Tel.: +49-4318807364
32 Fax: +49-4318802566

33 Email: c.jung@plantbreeding.uni-kiel.de

34

35

36 **Keywords:** *Chenopodium quinoa*, plant breeding, population structure, re-sequencing, adaptation,
37 domestication, genetic variation

38 Abstract

39 Quinoa germplasm preserves useful and substantial genetic variation, yet it remains untapped due to
40 a lack of implementation of modern breeding tools. We have integrated field and sequence data to
41 characterize a large diversity panel of quinoa. Whole-genome sequencing of 310 accessions
42 revealed 2.9 million polymorphic high confidence SNP loci. Highland and Lowland quinoa were
43 clustered into two main groups, with F_{ST} divergence of 0.36 and fast LD decay of 6.5 and 49.8 Kb,
44 respectively. A genome-wide association study uncovered 600 SNPs stably associated with 17
45 agronomic traits. Two candidate genes are associated with thousand seed weight, and a resistance
46 gene analog is associated with downy mildew resistance. We also identified pleiotropically acting
47 loci for four agronomic traits that are highly responding to photoperiod hence important for the
48 adaptation to different environments. This work demonstrates the use of re-sequencing data of an
49 orphan crop, which is partially domesticated to rapidly identify marker-trait association and
50 provides the underpinning elements for genomics-enabled quinoa breeding.

51 Introduction

52 Climate change poses a great threat to crop production worldwide. In temperate climates of the
53 world, higher temperatures and extended drought periods are expected. Moreover, crop production
54 in industrialized countries depends on only a few major crops resulting in narrow crop rotations.
55 Therefore, rapid transfer of wild species into crops using genetic modification and targeted
56 mutagenesis is currently discussed^{1,2}. Alternatively, orphan crops with a long tradition of
57 cultivation but low breeding intensity can be genetically improved by genomics assisted selection
58 methods. Quinoa (*Chenopodium quinoa* Willd.) is a pseudocereal crop species with a long history
59 of cultivation. It was first domesticated about 5000-7000 years ago in the Andean region. Quinoa
60 was a staple food during the pre-Columbian era, and the cultivation declined after the introduction
61 of crops like wheat and barley by the Spanish rulers. Owing to diversity, biotic and abiotic stress
62 tolerance, and ecological plasticity, quinoa can adapt to a broad range of agroecological regions^{3,4}.
63 Due to a high seed protein content and a favorable amino acid composition, its biological value is
64 even higher than beef, fish, and other major cereals^{5,6}. These favorable characteristics contributed
65 to the increasing worldwide popularity of quinoa among consumers and farmers.

66 A spontaneous hybridization event between two diploid species between 3.3 and 6.3 million years
67 ago gave rise to the allotetraploid species quinoa ($2n = 4x = 36$) with a genome size of 1.45-1.5 Gb
68 (nuclear DNA content $1C = 1.49$ pg)^{7,8}. A reference genome of the coastal Chilean quinoa
69 accession PI 614886 has been published with 44,776 predicted gene models together with whole-
70 genome re-sequencing of *C. pallidicaule* and *C. suecicum* species, close relatives of the A and B
71 subgenome donor species, respectively⁹. The organellar genomes are originated from the A-
72 genome ancestor¹⁰.

73 Quinoa belongs to the Amaranthaceae, together with some other economically important crops like
74 sugar beet, red beet, spinach, and amaranth. It reproduces sexually after self-pollination. Facultative
75 autogamy was reported for plants in close proximity with outcrossing rates in a range of 0.5 to
76 17.36 %^{11,12}. Thus, quinoa accessions are typically homozygous inbred lines. Nonetheless,
77 heterozygosity in some accessions has been reported, which indicates cross-pollination¹³. The
78 inflorescences are panicles, which are often highly branched. Florets are tiny, which is a significant
79 obstacle for hand-crossing. However, routine protocols for F₁ seed production in combination with
80 marker-assisted selection have been developed recently^{14,15}.

81 Systematic breeding of quinoa is still at its infancy compared to major crops. Until recently,
82 breeding has been mainly limited to Bolivia¹⁶ and Peru¹⁷, which are the major growing areas of
83 quinoa. Therefore, quinoa can be regarded as a partially domesticated crop. Many accessions suffer
84 from seed shattering, branching, and non-appropriate plant height, which are typical domestication
85 traits. Apart from these characters, grain yield and seed size, downy mildew resistance,
86 synchronized maturity, stalk strength, and low saponin content are major breeding objectives¹⁸. In
87 the past years, activities have been intensified to breed quinoa genotypes adapted to temperate
88 environments, for example, Europe, North America, and China¹⁹. Here, the major problem is the
89 adaptation to long-day conditions because quinoa is predominantly a short-day plant due to its
90 origin from regions near the equator.

91 There are only a few studies about the genetic diversity of quinoa. They were mainly based on
92 phenotypic observations^{16,20} and low throughput marker systems like random amplified
93 polymorphic DNA²¹, amplification fragment length polymorphisms²², and microsatellites²³. A
94 limited number of single nucleotide polymorphisms (SNP) based on expressed sequence tags were
95 published²⁴. Maughan, et al.²⁵ used five bi-parental populations to identify ca. 14,000 SNPs, from
96 which 511 KASP markers were developed. Genotyping 119 quinoa accessions gave the first insight
97 into the population structure of this species²⁵. Now, the availability of a reference genome enables
98 genome-wide genotyping (Jarvis et al. 2017). Jarvis, et al.⁹ re-sequenced 15 accessions and
99 identified ca. 7.8 million SNPs. In another study, 11 quinoa accessions were re-sequenced, and 8
100 million SNPs and ca. 842 thousand indels were identified²⁶.

101 Our study aimed to analyze the population structure of quinoa and patterns of variation by re-
102 sequencing a diversity panel encompassing germplasm from all over the world. Using millions of
103 markers, we performed a genome-wide association study using multiple-year field data. Here, we
104 identified QTLs that control agronomically important traits important for breeding cultivars to be
105 grown under long-day conditions. We are discussing the fundamental differences between an
106 underutilized crop and crops with a long breeding history. Our results provide useful information
107 for further understanding the genetic basis of agronomically important traits in quinoa and will be
108 instrumental for future breeding.

109 **Results**

110 **Re-sequencing 310 quinoa accessions reveals high sequence variation**

111 We assembled a diversity panel made of 310 quinoa accessions representing regions of major
112 geographical distributions of quinoa (Supplementary Fig. 1). The diversity panel comprises
113 accessions with different breeding history (Supplementary Table 1). We included 14 accessions
114 from a previous study, of which 7 are wild relatives⁹. The sequence coverage ranged from 4.07 to
115 14.55, with an average coverage of 7.78. We mapped sequence reads to the reference genome V2
116 (CoGe id53523). Using mapping reads, we identified 45,330,710 single nucleotide polymorphisms
117 (SNPs).

118

119

120

121

122

123

124 **Table 1:** Summary statistics of genome-wide single nucleotide polymorphisms identified in 303 quinoa accessions

Parameter	Type	All genotypes (quinoa only)	Highland population	Lowland population
SNP	Total	2,872,935	2,590,907	1,938,225
	Intergenic	2,452,347	2,227,952	1,649,310
	Introns	251,481	101,546	172,692
	Exons	114,654	214,945	78,248
Nucleotide diversity			5.78×10^{-4}	3.56×10^{-4}
Tajima's <i>D</i>			0.884	-0.384
Population divergences	<i>F_{ST}</i> (Weighted average)		0.36	

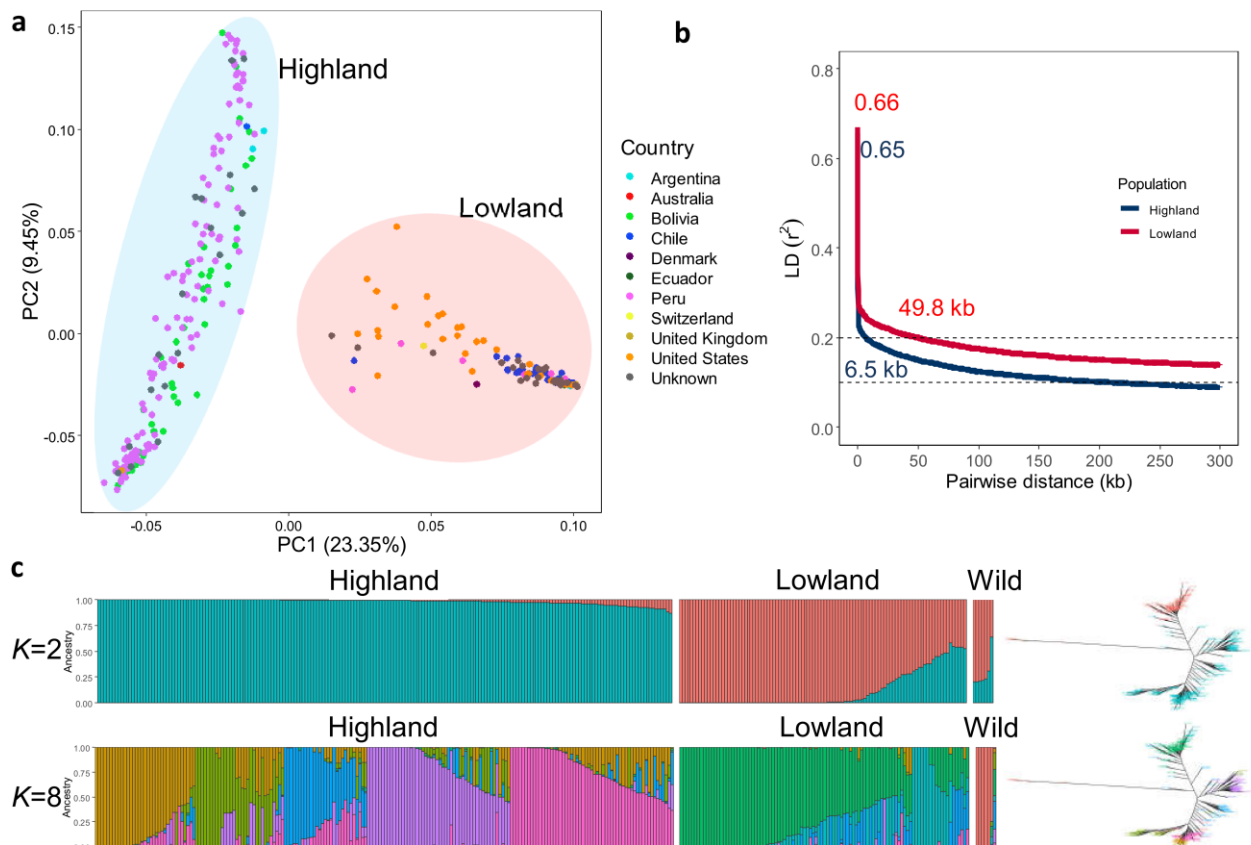
125

126 After filtering the initial set of SNPs, we identified 4.5 million SNPs in total for the base SNP set.
 127 We further filtered the SNPs for MAF >5 % (HCSNPs). We obtained 2.9 million high confident
 128 SNPs for subsequent analysis (Supplementary Table 2). Across the whole genome, SNP density
 129 was high, with an average of 2.39 SNPs/kb. However, SNP densities were highly variable between
 130 genomic regions and ranged from 0 to 122 SNPs/kb (Supplementary Fig. 3). We did not observe
 131 significant differences in SNP density between the two subgenomes (A subgenome 2.43 SNPs/kb;
 132 B subgenome 2.35 SNPs/kb). Then, we split the SNPs by their functional effects as determined by
 133 SnpEff²⁷. Among SNPs located in non-coding regions, 598,383 and 617,699 SNPs were located
 134 upstream (within 5kb from the transcript start site) and downstream (within 5kb from the stop site)
 135 of a gene, whereas 114,654 and 251,481 SNPs were located within exon and intron sequences,
 136 respectively (Table 1). We further searched for SNPs within coding regions. We found 70,604
 137 missense SNPs and 41,914 synonymous SNPs within coding regions of 53,042 predicted gene
 138 models.

139 Linkage disequilibrium and population structure of the quinoa diversity panel

140 Across the whole genome, LD decay between SNPs averaged 32.4 kb. We did not observe
 141 substantial LD differences between subgenome A (31.9kb) and subgenome B (30.7kb)
 142 (Supplementary Fig. 4C). The magnitude of LD decay among chromosomes did not vary drastically
 143 except for chromosome Cq6B, which exhibited a substantially slower LD decay (Supplementary
 144 Fig. 4 A and B).

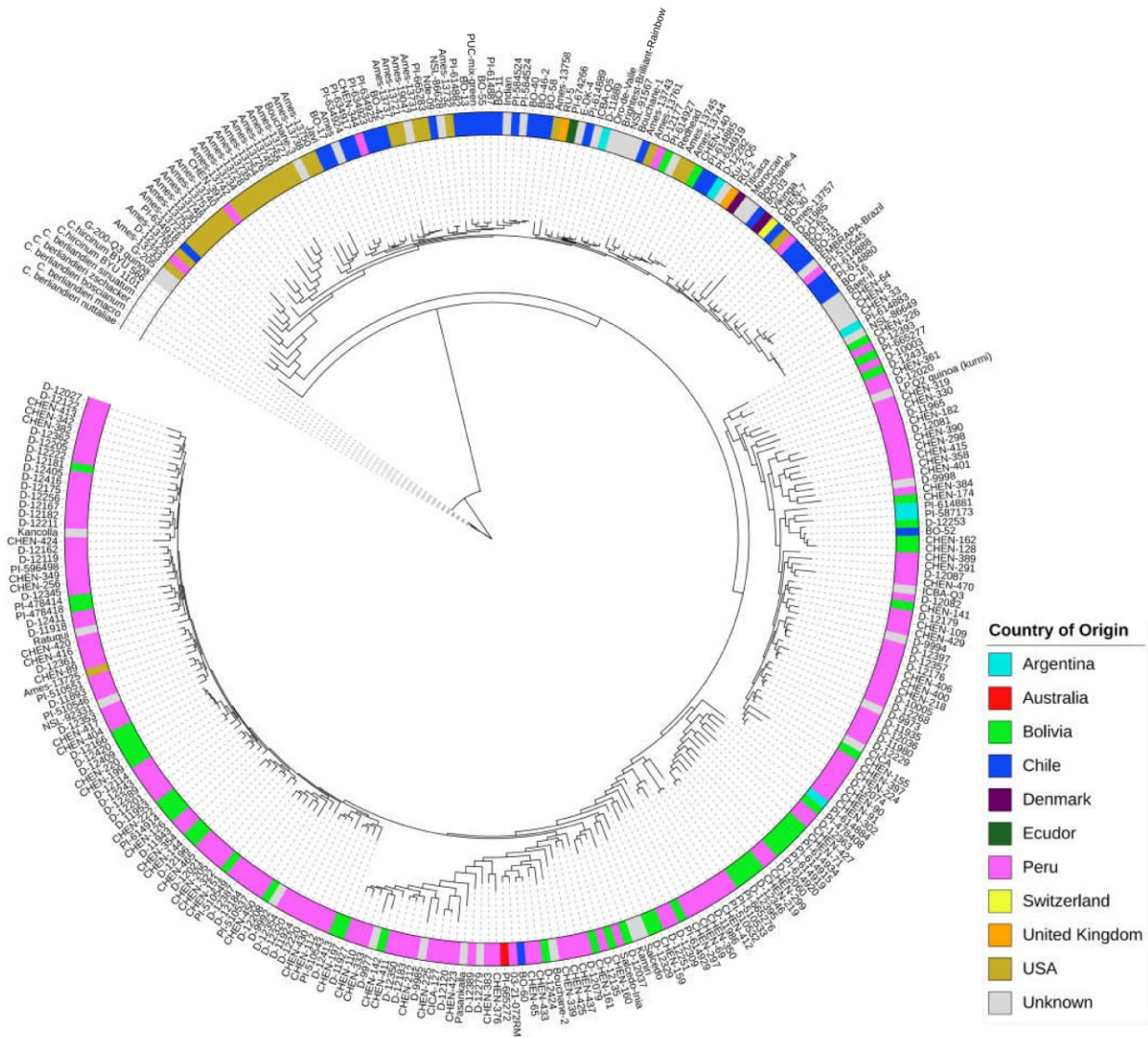
145 Then, we unraveled the population structure of the diversity panel. We performed principal
 146 component (PCA_(SNP)), population structure, and phylogenetic analyses. PCA_(SNP) showed two main
 147 clusters consistent with previous studies¹³. The first and second principal components (PC1_(SNP) and
 148 PC2_(SNP)) explained 23.35% and 9.45% of the variation, respectively (Fig. 1A). 202 (66.67%)
 149 accessions were assigned to subpopulation 1 (SP1) and 101 (33.33%) to subpopulation 2 (SP2). SP1
 150 comprised mostly Highland accessions, whereas Lowland accessions were found in SP2. PCA
 151 demonstrated a higher genetic diversity of the Highland population (Fig. 1A). We also calculated
 152 PCs for each chromosome separately. For 16 chromosomes, the same clustering as for the whole
 153 genome was calculated. Nevertheless, two chromosomes, Cq6B, and Cq8B showed three distinct
 154 clusters (Supplementary Fig. 5). This is due to the split of the Lowland population into two clusters.
 155 We reason that gene introgressions on these two chromosomes from another interfertile group
 156 might have caused these differences. This is also supported by a slower LD decay on chromosome
 157 Cq6B (Supplementary Fig. 4B).



158

159 **Fig. 1:** Genetic diversity and population structure of the quinoa diversity panel. (a) PCA of 303 quinoa accessions. PC1
 160 and PC2 represent the first two components of analysis, accounting for 23.35% and 9.45% of the total variation,
 161 respectively. The colors of dots represent the origin of accessions. Two populations are highlighted by different colors:
 162 Highland (light blue) and Lowland (pink). (b) Subpopulation wise LD decay in Highland (blue) and Lowland
 163 population (red). (c) Population structure is based on ten subsets of SNPs, each containing 50,000 SNPs from the
 164 whole-genome SNP data. Model-based clustering was done in ADMIXTURE with different numbers of ancestral
 165 kinships ($K=2$ and $K=8$). $K=8$ was identified as the optimum number of populations. Left: Each vertical bar represents
 166 an accession, and color proportions on the bar correspond to the genetic ancestry. Right: Unrooted phylogenetic tree of
 167 the diversity panel. Colors correspond to the subpopulation.

168 We also performed a population structure analysis with the ADMIXTURE software. We used cross-
 169 validation to estimate the most suitable number of populations. Cross-validation error decreased as
 170 the K value increased, and we observed that after $K = 5$, cross-validation error reached a plateau
 171 (Supplementary Fig. 6B). We observed allelic admixtures in some accessions, likely owing to their
 172 breeding history. The wild accessions were also clearly separated at the smallest cross-validation
 173 error of $K=8$, except two *C. hircinum* accessions (Fig. 1C). The reason for this could be that because
 174 *C. hircinum* is the closest crop wild relative, it also may have outcrossed with quinoa. The Highland
 175 population was structured into five groups, while the Lowland accessions were split into two
 176 subpopulations. The broad agro-climatic diversity of the Andean Highland germplasm might have
 177 caused a higher number of subpopulations.



178

179 **Fig. 2:** Maximum likelihood tree of 303 quinoa and seven wild *Chenopodium* accessions from the diversity panel.
 180 Colors are depicting the geographical origin of accessions.

181 We analyzed the phylogenetic relationships between quinoa accessions using 434,077 SNPs.
 182 Constructing a maximum likelihood tree gave rise to five clades (Fig. 2). We found that the
 183 placement of the wild quinoa species was concordant with the previous reports confirming that
 184 quinoa was domesticated from *C. hircinum*⁹. However, we found that the *C. hircinum* accession
 185 BYU 566 (from Chile) was placed at the base of both Lowland and Highland clades, which is in
 186 contrast to Jarvis, et al.⁹, where this accession was placed at the base of coastal quinoa. As
 187 expected, accessions from the USA and Chile are closely related because the USDA germplasm had
 188 been collected at these geographical regions.

189 Genomic patterns of variations between Highland and Lowland quinoa

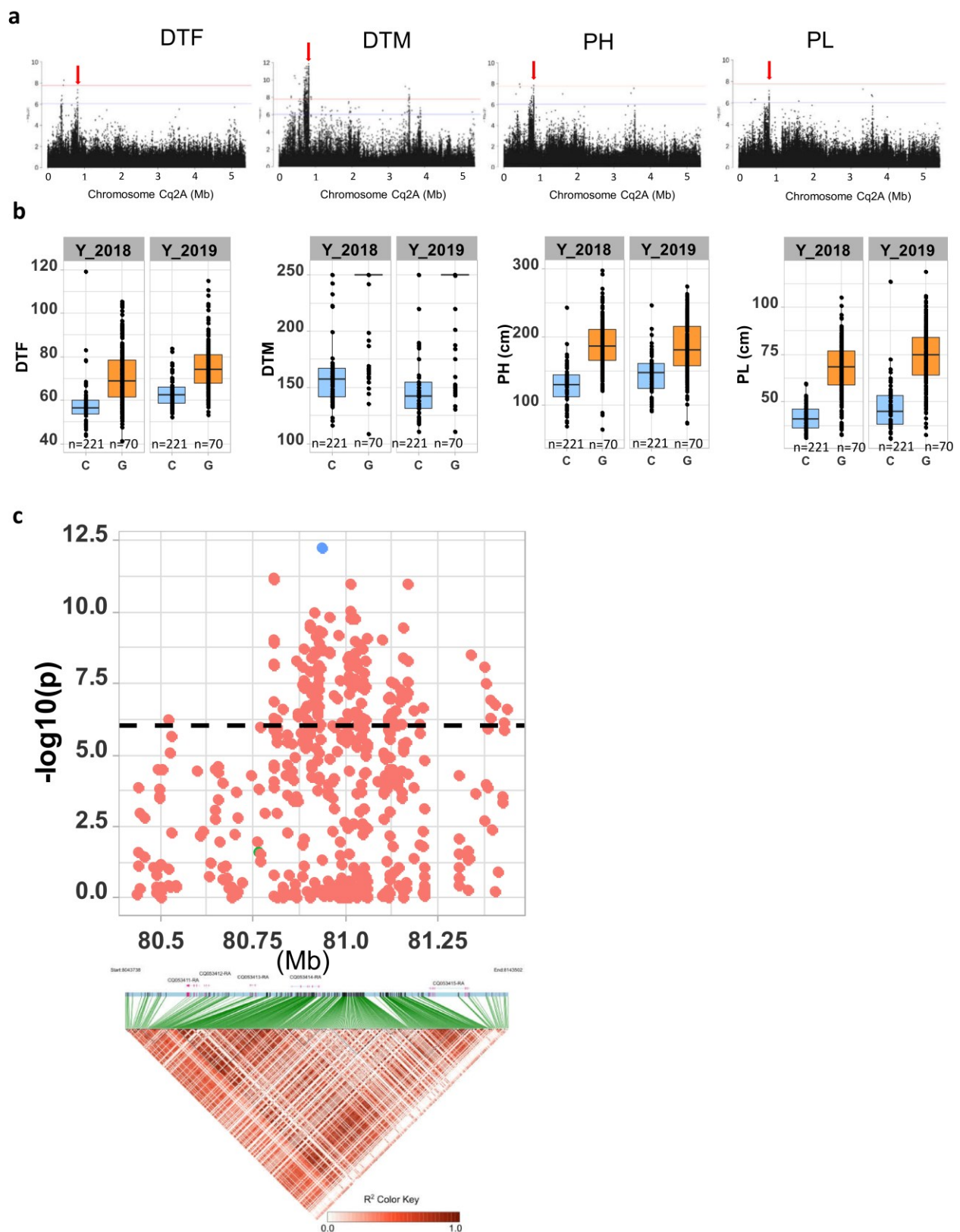
190 We were interested in patterns of variation in response to geographical diversification. We used
 191 principal component analysis derived clusters and phylogenetic analysis to define two diverged
 192 quinoa populations (namely Highland and Lowland). These divergent groups are highly correlated
 193 with Highland and Lowland geographical origin. We used the base SNP set to analyze diversity
 194 statistics. To detect genomic regions affected by the population differentiation, we measured the

195 level of nucleotide diversity using 10 kb non-overlapping windows²⁸. Then we calculated the
196 whole genome-wide LD decay across the two populations (Highland vs. Lowland); LD decayed
197 more rapidly in Highland quinoa (6.5 kb vs. 49.8 kb) (Fig. 1B). To measure nucleotide diversity, we
198 scanned the quinoa genome with non-overlapping windows of 10 kb in length in both populations
199 separately. The nucleotide diversity of the Highland population (5.78×10^{-4}) was 1.62 fold higher
200 compared to the Lowland population (3.56×10^{-4}) (Table 1 and Supplementary Fig. 7). We
201 observed left-skewed distribution and negative Tajima's *D* value (-0.3883) in the Lowland
202 populations indicating recent population growth (Table 1 and Supplementary Fig. 8). Genomic
203 regions favorable for adaptation to Highlands should have substantially lower diversity in the
204 Highland population than the Lowland population. Therefore, we calculated the nucleotide diversity
205 ratios between Highland and Lowland to identify major genomic regions that are underlying the
206 population differentiation. The F_{ST} value between populations was estimated to be 0.36, illustrating
207 strong population differentiation. Concerning the regions of variants, the number of exonic SNPs is
208 substantially higher in the Highland population (Table 1 and Supplementary Fig. 7).

209 Mapping agronomically important trait loci in the quinoa genome

210 We evaluated 13 qualitative and four dichotomous traits on 350 accessions across two different
211 environments. At the time of the final harvest, 254 accessions did not reach maturity (senescence).
212 All accessions produced seeds therefore used in seed analysis. For all traits, substantial phenotypic
213 variation among accessions was found. High heritabilities were calculated for all quantitative traits
214 except for number of branches (NoB) and stem lying (STL), which indicates that the phenotypic
215 variation between the accessions is mostly caused by genetic variation (Supplementary Table 3).
216 Trait correlations between years were also high (Supplementary Fig. 9), which is in accordance with
217 the heritability estimates. We found the strongest positive correlation between days to maturity
218 (DTM) and panicle length (PL), and plant height (PH) and PL, whereas the strongest negative
219 correlation was found between DTM and thousand seed weight (TSW) (Supplementary Fig. 10).
220 Then a principal component analysis was performed based on 12 quantitative traits ($PCA_{(PHEN)}$) to
221 explore the phenotypic relationship among quinoa accessions. The first two principal components
222 explained 62.12% of the phenotypic variation between the accessions. The score plot of the
223 principal components showed a similar clustering pattern as the SNP based PCA analysis
224 ($PCA_{(SNP)}$) (Fig. 1A and Supplementary Fig. 11A). $PCA_{(PHEN)}$ variables factor map indicated that
225 most Lowland accessions were high yielding with high TSW and dense panicles. Moreover, these
226 accessions are early flowering and early maturing, and they are short (Supplementary Fig. 11B).
227 Phenotype-based $PCA_{(PHEN)}$ also showed that the Lowland accessions are better adapted/selected for
228 cultivation in long-day photoperiods compared to the Highland accessions. These results are in
229 accordance with LD, nucleotide diversity, and Tajima's *D* estimations, implying the Lowland
230 accessions went through a stronger selection during breeding.

231 Then, we calculated the best linear unbiased estimates (BLUE) of the traits investigated. In total,
232 294 accessions shared the re-sequencing information and phenotypes out of 350 phenotypically
233 evaluated accessions. For GWAS analysis, we used ~2.9 million high-confidence SNPs. In total, we
234 identified 1480 significant ($P < 9.41e-7$) SNP-trait associations (MTA) for 17 traits (Supplementary
235 Fig. 12). The number of MTAs ranged from 4 (STL) to 674 (DTM) (Supplementary Table 4). In
236 agreement with previous reports, we defined an MTA as "consistent" when it was detected in both
237 years²⁹. We identified 600 consistent MTAs across eleven traits. TSW and DTM showed the highest
238 number of "consistent" associations. Among these, 143 MTAs are located within a gene, and 22
239 SNPs resulted in a missense mutation (Supplementary Table 5). MTA for the duration from bolting
240 to flowering (DTB to DTF), NoB, Seed yield, STL, and growth type (GT) were not "consistent"
241 between years (Supplementary Fig. 12). This is also reflected by the low heritability estimations of
242 these traits, indicating considerably higher genotype x environment interactions.



243

244 **Fig. 3:** Genomic regions associated with important agronomic traits (a) Significant marker-trait associations for days to
 245 flowering, days to maturity, plant height, and panicle density on chromosome Cq2A. Red color arrows indicate the SNP
 246 loci pleiotropically acting on all four traits. (b) Boxplots showing the average performance for four traits over two
 247 years, depending on single nucleotide variation (C or G allele) within locus Cq2A_8093547. (c) Local Manhattan plot
 248 from region 80.40 - 81.43 Mb on chromosome Cq2A associated with PC1 of the days to flowering (DTF), days to

249 maturity (DTM), plant height (PH), and panicle length (PL), and local LD heat map (bottom). The colors represent the
250 pairwise correlation between individual SNPs. Green color dots represent the strongest MTA (Cq2A_ 8093547).

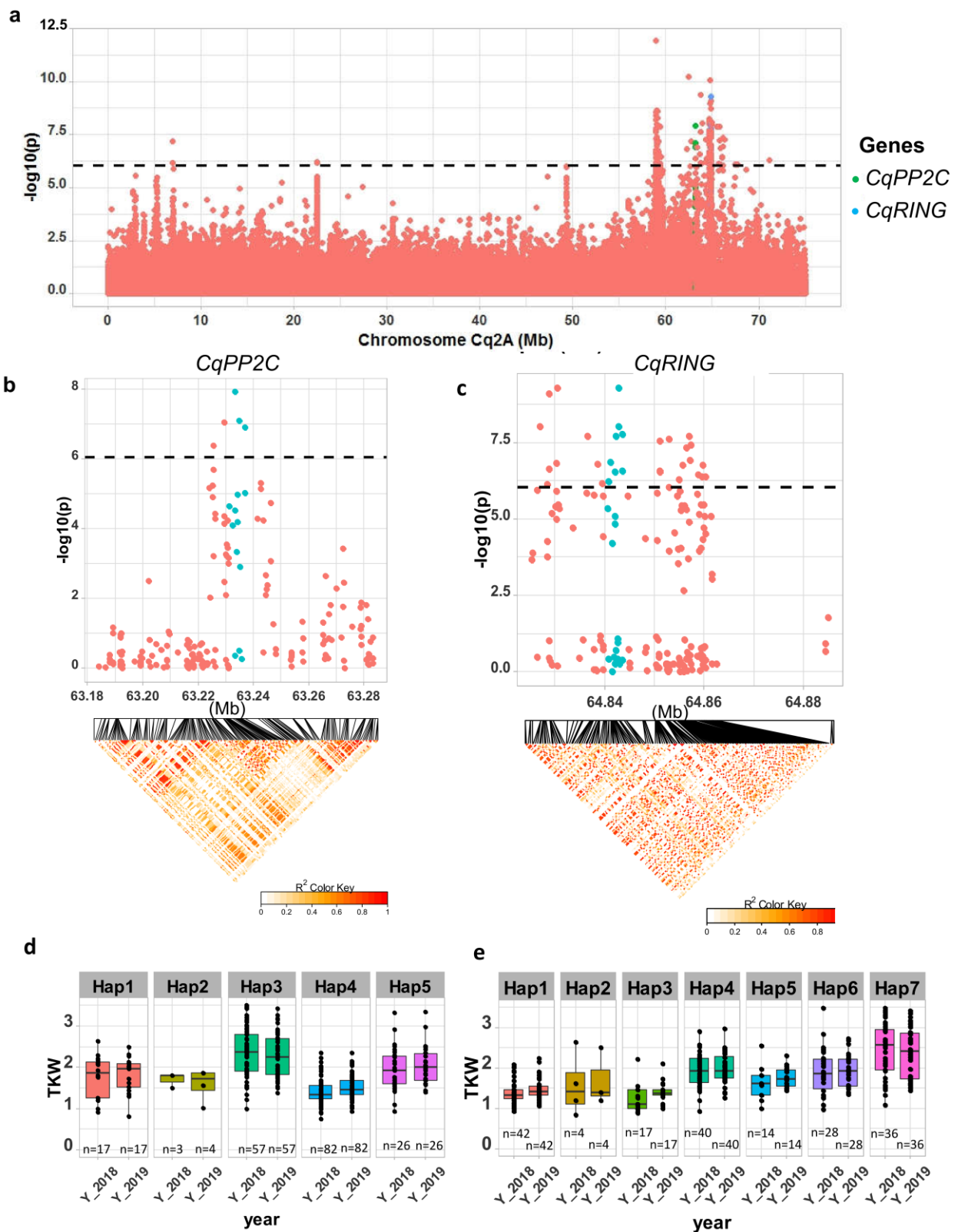
251 **Candidate genes for agronomically important traits**

252 First, we tested the resolution of our mapping study. We searched for major genes 50Kb down- and
253 upstream of significant SNPs for two qualitative traits in quinoa, flower color, and seed saponin
254 content. We identified highly significant MTAs for stem color on chromosome Cq1B (69.72-69.76
255 Mb). There are two genes (*CqCYP76AD1* and *CqDODAI1*) from the associated loci displaying high
256 homology to betalain synthesis pathway genes *BvCYP76AD1*³⁰ and *BvDODAI*³¹ from sugar beet
257 (Supplementary Fig. 14A and Supplementary Fig. 12). A significant MTA for saponin content on
258 chromosome Cq5B between 8.85 Mb to 9.2 Mb harbored the two *BHLH25* genes which have been
259 reported to control saponin content in quinoa⁹ (Supplementary Fig. 14B and Supplementary Fig.
260 12). This demonstrates that the marker density is high enough to narrow down to causative genes
261 underlying a trait.

262 Then, we examined four quantitative traits. We found seven MTA on chromosome Cq2A that are
263 associated with DTF, DTM, PH, and PL (cross-phenotype association), indicating evidence for
264 pleiotropic gene action (Fig. 3 and Supplementary Table 6). For further confirmation and to
265 investigate genes that are pleiotropically active on different traits, we followed a multivariate
266 approach³². First, we performed a PCA using the four phenotypes (cross-phenotypes). We found
267 89.94% of the variation could be explained by the first two principal components of the cross-
268 phenotypes (PCA_(CP)) (Supplementary Fig. 15). This indicates the adequate power of the PCA_(CP) to
269 reduce dimensions for the analysis of the cross-phenotypes association. We observed similar
270 clustering as in PCA_(SNP). Therefore, these results indicate that in quinoa, DTF, DTM, PH, and PL
271 are highly associated with population structure and thus, the adaptation to diverse environments.
272 Then, we performed a GWAS analysis using the first three PCs as traits (PC-GWAS)
273 (Supplementary Fig. 15C). We identified strong associations on chromosomes Cq2A, Cq7B (PC1),
274 and Cq8B (PC2) (Supplementary Fig. 16). Out of 468 MTAs (PC1:426 and PC2: 42) across the
275 whole genome, 222 (PC1:211 and PC2:11) are located within 95 annotated genes. We found 14
276 SNPs that changed the amino acid sequence in 12 predicted protein sequences of associated genes
277 (Supplementary Table 5). In the next step, we searched genes located within 50kb to an MTA.
278 Altogether, 605 genes were identified (PC1:520 and PC2:85) (Supplementary Table 7).

279 We found the region 80.50 -81.50 Mb on chromosome Cq2A to be of special interest because it
280 displays stable pleiotropic MTA for DTF, DTM, PH, and PL. The most significant SNP is located
281 within the *CqGLX2-2* gene, which encodes an enzyme of the glyoxalase family (Fig. 3). The
282 Arabidopsis *GLX2-1* has been shown to be essential for growth under abiotic stress³³. The allele
283 carrying a cytosine at the position with the most significant SNP resulted in early flowering,
284 maturing, and short panicles and plants (Fig. 3b). These traits are essential for the adaptation to
285 long-day conditions.

286 Thousand seed weight is an important yield component. We found a strong MTA between 63.2 –
287 64.87 Mb on chromosome Cq8B. Significantly associated SNPs were localized within two genes
288 (Fig. 4). One gene displays homology to *PP2C* encoding a member of the phosphatase-2C (*PP2C*)
289 protein family, which participates in Brassinosteroids signaling pathways and controls the
290 expression of the transcription factor *BZR1*³⁴. The second gene encodes a member of the RING-
291 type E3 ubiquitin ligase family. These genes are controlling seed size in soybean, maize, rice,
292 soybean, and Arabidopsis³⁵. We then checked haplotype variation and identified 5 and 7
293 haplotypes for *CqPP2C* and *CqRING* genes, respectively. Accessions carrying PP2C_hap3 and
294 RING_hap7 displayed larger seeds in both years (Fig. 4 and Supplementary Fig. 17)



295

296
297
298
299
300

Fig. 4: Identification of candidate genes for thousand seed weight. (a) Manhattan plot from chromosome Cq8B. Green and blue dots are depicting the *CqPP2C5* and the *CqRING* gene, respectively. (b) Top: Local Manhattan plot in the neighborhood of the *CqPP2C* gene. Bottom: LD heat map. (c) Top: Local Manhattan plot in the neighborhood of the *CqRING* gene. Bottom: LD heat map. Differences in thousand seed weight between five *CqPP2C* (d) and seven *CqRING* haplotypes (e).

301 Downy mildew is one of the major diseases in quinoa, which causes massive yield damage.
302 Notably, our GWAS identified strong MTA for resistance against this disease. The most significant
303 SNPs are located in subgenome A (Supplementary Fig. 12). Thus, the A-genome progenitor seems
304 to be the donor of downy mildew resistance. We identified a candidate gene within a region 38.99 -
305 39.03 Mb on chromosome Cq2A, which showed the highest significant association (Supplementary
306 Fig. 14C). This gene encodes a protein with an NBS-LRR (nucleotide-binding site leucine-rich
307 repeat) domain often found in resistance gene analogs with a function against mildew infection ³⁶.

308 Discussion

309 We assembled a diversity set of 303 quinoa accessions and seven accessions from wild relatives.
310 Plants were grown under northern European conditions, and agronomically important traits were
311 studied. In total, 2.9 million SNPs were found after re-sequencing. We found substantial phenotypic
312 and genetic variation. Our diversity set was structured into two highly diverged populations, and
313 genomic regions associated for this diversity were localized. Due to a high marker density,
314 candidate genes controlling qualitative and quantitative traits were identified. The high genetic
315 diversity and rapid LD breakdown are reflecting the short breeding history of this crop.

316 We were aiming to assemble the first diversity set, which represents the genetic variation of this
317 species. Therefore, we established a permanent resource that is genotypically and phenotypically
318 characterized. We believe that this collection is important for future studies due to the following
319 reasons: We observed substantial phenotypic variation for all traits and high homogeneity within
320 accessions. Moreover, low or absent phenotypic variation within accessions demonstrates
321 homogeneity as expected for a self-pollinating species. Therefore, the sequence of one plant is
322 representative of the whole accession, which is important for the power of the GWAS.

323 Today, over sixteen thousand accessions of quinoa are stored *ex-situ* in seed banks in more than 30
324 countries ³⁷. Despite the enormous diversity, only a few accessions have been genotyped with
325 molecular markers. We found a clear differentiation into Highland and Lowland quinoa. In previous
326 studies, five ecotypes had been distinguished: Valley type, Altiplano type, Salar type, Sea level
327 type, and Subtropical type ¹⁹. Adaptation to different altitudes, tolerance to abiotic stresses such as
328 drought and salt, and photoperiodic responses are the major factors determining ecotypes ¹⁸. In our
329 study, we could further allocate the quinoa accessions to five Highland and two Lowland
330 subpopulations. This demonstrates the power of high-density SNP mapping to identity finer
331 divisions at higher *K*. The origin of accessions and ecotype differentiation could be meaningfully
332 interpreted by combining the information from phylogenetic data and population structure. As we
333 expected, North American accessions (accessions obtained from USDA) were clustering with
334 Chilean accessions, suggesting sequence-based characterization of ecotypes would be more
335 informative and reproducible. Moreover, high-density SNP genotyping unveiled the origin of
336 unknown or falsely labeled gene bank accessions, as recently proposed by Milner, et al. ³⁸. The
337 geographical origin of 52 accessions from our panel was unknown. We suggest using phylogenetic
338 data and admixture results to complement the available passport data. For instance, two accessions
339 with origin recorded as Chile are closely related to Peruvian and Bolivian accessions, which
340 suggests that they are also originating from Highland quinoa.

341 What can we learn about the domestication of quinoa and its breeding history by comparing our
342 results with data from other crops? LD decay is one parameter reflecting the intensity of breeding.
343 LD decay in quinoa (32.4 kb) is faster than in most studies with major crop species, e.g. rapeseed
344 (465.5 kb) ³⁹, foxtail millet (*Setaria italica*, 100 kb) ⁴⁰, pigeonpea (*Cajanus cajan*, 70 kb) ⁴¹,
345 soybean (150 kb) ⁴² and rice (200 kb) ⁴³. Although comparisons must be regarded with care due to
346 different numbers of markers and accessions, different types of reproduction, and the selection

347 intensity, the rapid LD decay in quinoa reflects its short breeding history and low selection
348 intensity. Moreover, quinoa is a self-pollinating species where larger linkage blocks could be
349 expected. However, cross-pollination rates in some accessions can be up to 17.36 %¹², which is
350 exploited by small Andean farmers who grow mixed quinoa accessions to ensure harvest under
351 different biotic and abiotic stresses. This may facilitate a certain degree of cross-pollination and
352 admixture.

353 Interestingly, the LD structure between Highland and Lowland populations is highly contrasting
354 (6.5 vs. 49.8 kb), indicating larger LD blocks in the Lowland population. Low nucleotide diversity
355 and negative Tajima's *D* were also observed in the Lowland population compared to Highland
356 quinoa. The population differentiation index and LD differences have been used to test the
357 hypothesis of multiple domestication events. As an example, different domestication bottlenecks
358 have been reported for japonica (LD decay: 65 kb) and indica rice (LD decay: 200 kb)⁴⁴. The
359 estimated *F_{ST}* value from this study (0.36) is in the similar range of *F_{ST}* estimates in rice subspecies
360 *indica* and *japonica* (0.55)⁴⁵ and melon (*Cucumis melo*) subspecies *melo* and *agrestis* (0.46)⁴⁶.
361 Two hypotheses have been proposed for the domestication of quinoa from *C. hircinum*; (1) one
362 event that gave rise to Highland quinoa and subsequently to Lowland quinoa and (2) two separate
363 domestication events giving rise to Highland and Lowland quinoa independently⁹. However, our
364 study is not strictly following the second hypothesis because *C. hircinum* accession BYU 566 was
365 basal to both clades of the phylogenetic tree (Highland and Lowland). Moreover, our wild
366 Chenopodium germplasm does not represent enough diversity for in-depth analysis of
367 domestication events. Therefore, we propose three possible scenarios to explain strong differences
368 in LD structure, nucleotide diversity, Tajima's *D* and *F_{ST}*, (1) two independent domestication events
369 with a strong bottleneck on lowland populations, (2) a single domestication but strong population
370 growth after adaptation of lowland quinoa or (3) strong adaptive selection after domestication. To
371 understand the history and genetics of domestication, it will be necessary to sequence a large
372 representative set of outgroup species such as *C. berlandieri*, *C. hircinum*, *C. pallidicaule*, and *C.*
373 *suecicum*.

374 Apart from marker density and sample size, the power of GWAS depends on the quality of the
375 phenotypic data. Plants were grown in Northern Europe. Therefore, the MTAs are, first of all,
376 relevant for temperate long-day climates. The share of genetic variances and thus, the heritabilities
377 were high across environments. We expect higher genotype x environment interaction for flowering
378 time, days to maturity, plant height, and panicle length if short-day environments will be included
379 because many accessions have a strong day-length response (data not shown). Furthermore, the
380 positions of genes controlling Mendelian traits were precisely coinciding with significant SNP
381 positions, as exemplified by the genes associated with saponin content and flower color. Hence, the
382 diversity panel provides sufficient power to identify SNP-trait associations for important agronomic
383 traits such as TSW and downy mildew tolerance. In different plant species, seed size is controlled
384 by six different pathways³⁵. We found two important genes controlling seed size from the
385 Brassinosteroid (*CqPP2C*) and the ubiquitin-proteasome (*CqRING*) pathway. The non-functional
386 allele of soybean *PP2C1* resulted in small seeds³⁴. We detected a superior haplotype (PP2C_hap3),
387 which results in larger seeds. *CqRING* encodes an E3 ubiquitin ligase protein. There are two RING-
388 type E3 ubiquitins known as *DA1* and *DA2*, which are involved in seed size controlling pathway.
389 They were found in Arabidopsis rice, maize, and wheat. Downy mildew is the most acute disease
390 for quinoa, caused by the fungus *Peronospora variabilis*⁴⁷. A recent study attempted identification
391 of genes based on a GWAS analysis. However, no significant associations were found, probably
392 due to the lack of power because of the small number of accessions used (61 and 88)⁴⁸. In our
393 study, a strong MTA suggests that the NBS-LRR gene on chromosome Cq2A contributes to downy
394 mildew resistance in quinoa. We propose using this sequence for marker-assisted selection in
395 segregating F₂ populations produced during pedigree breeding of quinoa.

396 In this study, the advantage of multivariate analysis of cross-phenotype association became obvious.
397 We could identify candidate genes with a pleiotropic effect on days to flowering, days to maturity,
398 plant height, and panicle length. Interestingly, the most significant SNP was residing within a
399 putative *GLX-2* ortholog. *GLX* genes, among other functions, have been shown to impact cell
400 division and proliferation in *Amaranthus paniculatus*⁴⁹. Therefore, the *CqGLX-2* gene is one
401 candidate for controlling day length response.

402 This study also has a major breeding perspective. We aimed to elucidate the potential of quinoa for
403 cultivation in temperate climates. Evidently, many accessions are not adapted to northern European
404 climate and photoperiod conditions because they flowered too late and did not reach maturity before
405 October. Nevertheless, 48 accessions are attractive as crossing partners for breeding programs
406 because they are insensitive to photoperiod or long-day responsive. Moreover, they are attractive
407 due to their short plant height, low tillering capacity, favorable inflorescence architecture, and high
408 TSW. These are important characters for mechanical crop cultivation and combine harvesting. The
409 MTA found in this study offers a perspective to use parents with superior phenotypes in crossing
410 programs. We suggest a genotype building strategy by pyramiding favorable alleles (haplotypes). In
411 this way, also accessions from our diversity set, which are not adapted to long-day conditions but
412 with favorable agronomic characters, will be considered. Then, favorable genotypes will be
413 identified from offspring generations by marker-assisted selection using markers in LD with
414 significant SNPs. Furthermore, the MTA from this study will be useful for allele mining in quinoa
415 germplasm collections to identify yet unexploited genetic variation.

416 **Materials and Methods**

417 **Plant materials and growth conditions**

418 We selected 350 quinoa accessions for phenotyping, and of these, 296 were re-sequenced in this
419 study. Re-sequencing data of 14 additional accessions that had already been published⁹ were also
420 included in the study, together with the wild relatives (*C. belandieri* and *C. hircinum*)⁹. These
421 accessions represent different geographical regions of quinoa cultivation (Supplementary Table 1).
422 Plants were grown in the field in Kiel, Northern Germany, in 2018 and 2019. Seeds were sown in
423 the second week of April in 35x multi-tray pots. Then plants were transplanted to the field in the
424 first week of May as single-row plots in a randomized complete block design with three blocks. The
425 distances between rows and between plants were set to 60 cm and 20 cm, respectively. Each row
426 plot contained seven plants per accession.

427 We recorded days to bolting (DTB) as BBCH51 and days to flowering (DTF) as BBCH60 twice a
428 week during the growth period. Days to maturity (DTM) was determined when plants reached
429 complete senescence (BBSHC94). If plants did not reach this stage, DTM was set as 250 days. In
430 both years, plants were harvested in the second week of October. Plant height (PH), panicle length
431 (PL), and the number of branches (NoB) were phenotyped at harvest. Stem lying (STL)
432 (Supplementary Fig. 2) was scored on a scale from one to five, where score one indicates no stem
433 lying. Similarly, panicle density was recorded on a scale from one to seven, where density one
434 represents lax panicles, and panicle density seven represents highly dense panicles. Flower color
435 and stem color were determined by visual observation. Pigmented and non-pigmented plants were
436 scored as 1 and 0, respectively. Growth type was classified into two categories and analyzed as a
437 dichotomous trait as well. We observed severe mildew infection in 2019. Therefore, we scored
438 mildew infection on a scale from 1 to 3, where 1 equals no infection, and 3 equals severe infection.

439 **Statistical analysis**

440 We calculated the best linear unbiased estimates of the traits across years by fitting a linear mixed
441 model using the lme4 R package ⁵⁰. We used the following model:

$$442 Y_{ijk} = \mu + \text{Accession}_i + \text{Block}_j + \text{Year}_k + (\text{Accession} \times \text{Block})_{ij} + (\text{Accession} \times \text{Year})_{ik} + \text{Error}_{ijk}$$

443 Where μ is the mean, Accession_i is the genotype effect of the i -th accession, Block_j is the effect of
444 the j -th Block, Year_k is the effect of the k -th year, $(\text{Accession} \times \text{Block})_{ij}$ is the Accession-Block
445 interaction effect, $\text{Accession} \times \text{Year}_{ik}$ is the accession-year interaction effect, Error_{ijk} is the error of
446 the j -th block in the k -th year. We treated all items as random effects for heritability estimation, and
447 for best linear unbiased estimates (BLUE), accessions were treated as fixed effects. We analyzed the
448 principle components of phenotypes using the R package FactoMineR ⁵¹.

449 **Genome sequencing and identification of genomic variations**

450 For DNA extraction, two plants per genotype were grown in a greenhouse at the University of
451 Hohenheim, and two leaves from a single two-months old plant were collected and frozen
452 immediately. DNA was subsequently extracted using the AX Gravity DNA extraction kit (A\&A
453 Biotechnology, Gdynia, Poland) following the manufacturer's instructions. Purity and quality of
454 DNA were controlled by agarose gel electrophoresis and the concentration determined with a Qubit
455 instrument using SYBR green staining. Whole-genome sequencing was performed for 312
456 accessions at Novogene (China) using short-reads Illumina NovaSeq S4 Flowcell technology and
457 yielded an average of 10 Gb of paired-end (PE) 2 x 150 bp reads with quality $Q > 30$ Phred score per
458 sample, which is equivalent to $\sim 7X$ coverage of the haploid quinoa genome (~ 1.45 Gb). We then
459 used an automated pipeline ([https://github.com/IBEXCluster/IBEX-](https://github.com/IBEXCluster/IBEX-SNPcaller/blob/master/workflow.sh)
460 [SNPcaller/blob/master/workflow.sh](https://github.com/IBEXCluster/IBEX-SNPcaller/blob/master/workflow.sh)) compiled based on the Genome Analysis Toolkit. Raw
461 sequence reads were filtered with trimmomatic-v0.38 ⁵² using the following criteria: LEADING:20;
462 TRAILING:20; SLIDINGWINDOW:5:20; MINLEN:50. The filtered paired-end reads were then
463 individually mapped for each sample against an improved version of the QQ74 quinoa reference
464 genome (CoGe id53523) using BWA-MEM (v-0.7.17) ⁵³ followed by sorting and indexing using
465 samtools (v1.8) ⁵⁴. Duplicated reads were marked, and read groups were assigned using the Picard
466 tools (<http://broadinstitute.github.io/picard/>). Variants were identified with GATK (v4.0.1.1) ^{55 56}
467 using the "--emitRefConfidence" function of the HaplotypeCaller algorithm and "--heterozygosity"
468 value set at 0.005 to call SNPs and InDels for each accession. Individual g.vcf files for each sample
469 were then compressed and indexed with tabix (v-0.2.6) ⁵⁷ and combined into chromosome g.vcf
470 using GenomicsDBImport function of GATK. Joint genotyping was then performed for each
471 chromosome using the function GenotypeGVCFs of GATK. To obtain high confidence variants, we
472 excluded SNPs with the VariantFiltration function of GATK with the criteria: $QD < 2.0$; $FS > 60.0$;
473 $MQ < 40.0$; $MQRankSum < -12.5$; $ReadPosRankSum < -8.0$ and $SOR > 3.0$. Then, SNP loci
474 which contained more than 70% missing data, were filtered by VCFtools ⁵⁸ (v0.1.5), which resulted
475 in our initial set of $\sim 45M$ SNPs for all the 332 accessions, including 20 previously re-sequenced
476 accessions ⁹. All resequencing data are submitted to SRA under project id BioProject
477 PRJNA673789.

478 In our panel, we had three triplicates for quality checking and nine duplicates between Jarvis et al.
479 2017 and 312 newly re-sequenced accessions. In order to remove duplicates, as a preliminary
480 analysis, we removed SNP loci with a minimum mean-depth < 5 across samples and SNP loci with
481 more than 5% missing data. Then, we filtered SNPs with a minor allele frequency lower than 0.05
482 ($MAF < 0.05$). After these filtering steps, we obtained a VCF file that contained 229,017 SNPs.
483 Then, we construct a maximum likelihood (ML) tree. First, we used the modelFinder ⁵⁹ in IQ-TREE

484 v1.6.619 (Nguyen et al. 2015) to determine the best model for ML tree construction. We selected
485 GTR+F+R8 (GTR: General time-reversible, F: Empirical base frequencies, R8: FreeRate model) as
486 the best fitting model according to the Bayesian Information Criterion (BIC) estimated by the
487 software. We used 1000 replicates with ultrafast bootstrapping (UFboots)⁶⁰ to check the reliability
488 of the phylogenetic tree. To visualize the phylogenetic tree, we used the Interactive Tree Of Life
489 tool (<https://itol.embl.de/>)⁶¹. Then, based on the phylogenetic tree, we removed duplicate accessions
490 and accessions with unclear identity. After the quality control, we retained 310 accessions (303
491 quinoa accessions and 7 wild *Chenopodium* accessions).

492 Then we used the initial SNP set and defined two subsets using the following criteria: (1) A base
493 SNP set of 5,817,159 biallelic SNPs obtained by removing SNPs with more than 50% missing
494 genotype data, minimum mean depth less than five, and minor allele frequency less than 1%. (2) A
495 high confidence (HCSNP) set of 2,872,935 SNPs from the base SNP set by removing SNPs with a
496 minor allele frequency of less than 5%. The base SNP set was used for the diversity statistics, and
497 the HCSNPs set was used for GWAS analysis.

498 We annotated the HCSNP using SnpEff 4.3T²⁷ and a custom database²⁷ based on the QQ74
499 reference genome and annotation (CoGe id53523). Afterward, we extracted the SNP annotations
500 using SnpSift⁶². Based on the annotations, SNPs were mainly categorized into five groups, (1)
501 upstream of the transcript start site (5kb), (2) downstream of the transcript stop site (5kb), (3)
502 coding sequence (CDS), (4) intergenic, and (5) intronic. We used SnpEff to categorize SNPs in
503 coding regions based on their effects such as synonymous, missense, splice acceptor, splice donor,
504 splice region, start lost, start gained, stop lost, and spot retained.

505 **Phylogenetic analysis and population structure analysis**

506 For population structure analysis, we employed SNP subsets, as demonstrated in previous studies,
507 to reduce the computational time⁶³. We created ten randomized SNP sets, each containing 50,000
508 SNPs. To create subsets, first, the base SNP set was split into 5000 subsets of an equal number of
509 SNPs. Then, 10 SNPs from each subset were randomly selected, providing a total of 50,000 SNPs
510 in a randomized set (randomized 50k set). We then repeated this procedure for nine more times and
511 finally obtained ten randomized 50k sets. Population structure analysis was conducted using
512 ADMIXTURE (Version: 1.3)⁶⁴. We ran ADMIXTURE for each subset separately with a
513 predefined number of genetic clusters K from 2 to 10 and varying random seeds with 1000
514 bootstraps. Also, we performed the cross-validation (CV) procedure for each run. Obtained Q
515 matrices were aligned using the greedy algorithm in the CLUMPP software⁶⁵. Population structure
516 plots were created using custom R scripts. We then combined SNP from the ten subsets to create a
517 single SNP set of 434,077 unique SNPs for the phylogenetic analysis. We used the same method
518 mentioned above to create the phylogenetic tree. Here we selected the model GTR+F+R6 based on
519 the BIC estimates. For the principal component analysis (PCA) we used the HCSNP set and
520 analysis was done in R package SNPrelate⁶⁶. We estimated the top 10 principal components. The
521 first (PC1) and second (PC2) were plotted using custom R scripts.

522 **Genomic patterns of variations**

523 Using the base SNP set, we calculated nucleotide diversity (π) for subpopulations and π ratios for
524 Highland and Lowland population regions with the top 1% ratios of $\pi_{\text{Highland}} / \pi_{\text{Lowland}}$ candidate
525 regions for population divergence. We also estimated Tajima's *D* values for both populations to
526 check the influence of selection on populations. F_{ST} values were calculated between Highland and
527 Lowland populations using the 10kb non-overlapping window approach. Nucleotide diversity,
528 Tajima's *D*, and F_{ST} calculations were carried out in VCFtools (v0.1.5)⁵⁸.

529 Linkage disequilibrium analysis

530 First, we calculated linkage disequilibrium in each population separately (Highland and Lowland).
531 Then, LD was calculated in the whole population, excluding wild accessions. For LD calculations,
532 we further filtered the HCSNP set by removing SNPs with >80% missing data²⁹. Using a set of
533 2,513,717 SNPs, we calculated the correlation coefficient (r^2) between SNPs up to 300kb apart by
534 setting -MaxDist 300 and default parameters in the PopLDdecay software⁶⁷. LD decay was plotted
535 using custom R scripts based on the ggplot2 package.

536 Genome-wide association study

537 We used the best linear unbiased estimates (BLUE) of traits and HCSNPs for the GWAS analysis.
538 Morphological traits were treated as dichotomous traits and analyzed using generalized mixed linear
539 models with the lme4 R software package⁵⁰. We used population structure and genetic relationships
540 among accessions to minimize false-positive associations. Population structure represented by the
541 PC was estimated with the SNPrelate software⁶⁶. Genetic relationships between accessions were
542 represented by a kinship matrix calculated with the efficient mixed-model association expedited
543 (EMMAX) software⁶⁸ using HCSNPs. Then, we performed an association analysis using the mixed
544 linear model, including K and P matrices in EMMAX. We estimated the effective number of SNPs
545 ($n=1,062,716$) using the Genetic type I Error Calculator (GEC)⁶⁹. We set the significant *P*-value
546 threshold (Bonferroni correction, $0.05/n$, $-\log_{10}(4.7e-08)=7.32$) and suggestive significant threshold
547 ($1/n$, $-\log_{10}(9.41e-7)=6.02$) to identify significant loci underlying traits. We plotted SNP *P*-values
548 on Manhattan plots using the qqman R package⁷⁰.

549 Acknowledgments

550 We thank David Jarvis for providing the updated version of the quinoa reference genome. We thank
551 Monika Bruisch, Brigitte Neidhardt-Olf, Elisabeth Kokai-Kota, Verena Kowalewski, and Gabriele
552 Fiene for technical assistance. The financial support of this work was provided by the Competitive
553 Research Grant (Grant No. OSR-2016-CRG5- 466 2966-02) of the King Abdullah University of
554 Science and Technology, Saudi Arabia and baseline funding from KAUST to Mark Tester.

555 Author contributions

556 C.J, M.T, and N.E directed the project and conceived the research. D.S.R.P conducted genomic data
557 analysis and GWAS analysis. D.S.R.P and N.E performed field experiments and phenotyping. E.R
558 conducted SNP identifications. G.W and S.M.S selected and assembled the diversity panel. K.S
559 contributed to DNA isolation, library preparation for genome sequencing. D.S.R.P, together with all
560 authors, wrote and finalized the manuscript.

561 Competing interests

562 The authors declare no competing interests.

563 Data availability

564 The raw sequencing data have been submitted to the NCBI Sequence Read Archive (SRA) under
565 the BioProject PRJNA673789. Quinoa reference genome version 2 is available at CoGe database
566 under genome id 53523. Source data are provided with the paper.

567 Code availability

568 Custom scripts used for SNP calling are available on GitHub:
569 <https://github.com/IBEXCluster/IBEX-SNPcaller/blob/master/workflow.sh>. Additional information
570 on other custom scripts will be available upon request.

571 References

- 572 1. Stetter, M.G., Gates, D.J., Mei, W. & Ross-Ibarra, J. How to make a domesticate. *Current*
573 *Biology* **27**, R896-R900 (2017).
- 574 2. Li, T. *et al.* Domestication of wild tomato is accelerated by genome editing. *Nature*
575 *biotechnology* **36**, 1160-1163 (2018).
- 576 3. Ruiz, K.B. *et al.* Quinoa biodiversity and sustainability for food security under climate
577 change. A review. *Agronomy for sustainable development* **34**, 349-359 (2014).
- 578 4. González, J.A., Eisa, S., Hussin, S. & Prado, F.E. Quinoa: an Incan crop to face global
579 changes in agriculture. *Quinoa: Improvement and sustainable production*, 1-18 (2015).
- 580 5. James, L.E.A. Quinoa (*Chenopodium quinoa* Willd.): composition, chemistry, nutritional,
581 and functional properties. *Advances in food and nutrition research* **58**, 1-31 (2009).
- 582 6. Vega-Gálvez, A. *et al.* Nutrition facts and functional potential of quinoa (*Chenopodium*
583 *quinoa* willd.), an ancient Andean grain: a review. *Journal of the Science of Food and*
584 *Agriculture* **90**, 2541-2547 (2010).
- 585 7. Palomino, G., Hernández, L.T. & de la Cruz Torres, E. Nuclear genome size and
586 chromosome analysis in *Chenopodium quinoa* and *C. berlandieri* subsp. *nuttalliae*.
587 *Euphytica* **164**, 221 (2008).
- 588 8. Kolano, B., Siwinska, D., Pando, L.G., Szymanowska-Pulka, J. & Maluszynska, J. Genome
589 size variation in *Chenopodium quinoa* (Chenopodiaceae). *Plant systematics and evolution*
590 **298**, 251-255 (2012).
- 591 9. Jarvis, D.E. *et al.* The genome of *Chenopodium quinoa*. *Nature* **542**, 307 (2017).
- 592 10. Maughan, P.J. *et al.* Mitochondrial and chloroplast genomes provide insights into the
593 evolutionary origins of quinoa (*Chenopodium quinoa* Willd.). *Scientific Reports* **9**, 185
594 (2019).
- 595 11. Gandarillas, H., Alandia, S., Cardozo, A. & Mujica, A. Qinoa y Kaniwa cultivos Andinos.
596 *Instituto Interamericano de Ciencias Agrícolas, Bogotá, Colombia* (1979).
- 597 12. Silvestri, V. & Gil, F. Alogamia en quinoa. *Revista de la Facultad de Ciencias Agrarias*
598 **32**(2000).
- 599 13. Christensen, S.A. *et al.* Assessment of genetic diversity in the USDA and CIP-FAO
600 international nursery collections of quinoa (*Chenopodium quinoa* Willd.) using
601 microsatellite markers. *Plant Genetic Resources: Characterisation and Utilisation* **5**, 82-95
602 (2007).

- 603 14. Peterson, A., Jacobsen, S.E., Bonifacio, A. & Murphy, K. A crossing method for Quinoa.
604 *Sustainability (Switzerland)* **7**, 3230-3243 (2015).
- 605 15. Emrani, N. *et al.* An efficient method to produce segregating populations in quinoa
606 (*Chenopodium quinoa* Willd.). (2020).
- 607 16. Gandarillas, H. Botánica. in *Quinoa y kañiwa: cultivos Andinos* (ed. Tapia, M.E.) (CIID,
608 Bogotá, 1979).
- 609 17. Bonifacio, A., Gomez-Pando, L. & Rojas, W. Quinoa breeding and modern variety
610 development. *State of the Art Report on Quinoa Around the World* (2013).
- 611 18. Gomez-Pando, L. Quinoa breeding. *Quinoa: Improvement and Sustainable Production*, 87-
612 108 (2015).
- 613 19. Murphy, K.M. *et al.* Quinoa breeding and genomics. *Plant Breeding Reviews* **42**, 257-320
614 (2018).
- 615 20. Wilson, H.D. Allozyme variation and morphological relationships of *Chenopodium*
616 *hircinum* (s.l.). *Syst. Bot* **13**(1988).
- 617 21. Ruas, P.M., Bonifacio, A., Ruas, C.F., Fairbanks, D.J. & Andersen, W.R. Genetic
618 relationship among 19 accessions of six species of *Chenopodium* L., by Random Amplified
619 Polymorphic DNA fragments (RAPD). *Euphytica* **105**, 25-32 (1999).
- 620 22. Rodríguez, L.A. & Isla, M.T. Comparative analysis of genetic and morphologic diversity
621 among quinoa accessions (*Chenopodium quinoa* Willd.) of the South of Chile and highland
622 accessions. *Journal of Plant Breeding and Crop Science* **1**, 210-216 (2009).
- 623 23. Mason, S. *et al.* Development and use of microsatellite markers for germplasm
624 characterization in quinoa (*Chenopodium quinoa* Willd.). *Crop Science* **45**, 1618-1630
625 (2005).
- 626 24. Coles, N. *et al.* Development and use of an expressed sequenced tag library in quinoa
627 (*Chenopodium quinoa* Willd.) for the discovery of single nucleotide polymorphisms. *Plant*
628 *Science* **168**, 439-447 (2005).
- 629 25. Maughan, P.J. *et al.* Single Nucleotide Polymorphism Identification, Characterization, and
630 Linkage Mapping in Quinoa. *The Plant Genome Journal* **5**, 114 (2012).
- 631 26. Zhang, T. *et al.* Development of novel InDel markers and genetic diversity in *Chenopodium*
632 *quinoa* through whole-genome re-sequencing. *BMC Genomics* **18**, 685 (2017).
- 633 27. Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide
634 polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w¹¹¹⁸; iso-
635 2; iso-3. *Fly* **6**, 80-92 (2012).
- 636 28. Varshney, R.K. *et al.* Pearl millet genome sequence provides a resource to improve
637 agronomic traits in arid environments. *Nature biotechnology* **35**, 969-976 (2017).

- 638 29. Varshney, R.K. *et al.* Resequencing of 429 chickpea accessions from 45 countries provides
639 insights into genome diversity, domestication and agronomic traits. *Nature Genetics* **51**,
640 857-864 (2019).
- 641 30. Hatlestad, G.J. *et al.* The beet *R* locus encodes a new cytochrome P450 required for red
642 betalain production. *Nature Genetics* **44**, 816-820 (2012).
- 643 31. Bean, A. *et al.* Gain-of-function mutations in beet DODA 2 identify key residues for
644 betalain pigment evolution. *New Phytologist* **219**, 287-296 (2018).
- 645 32. Solovieff, N., Cotsapas, C., Lee, P.H., Purcell, S.M. & Smoller, J.W. Pleiotropy in complex
646 traits: challenges and strategies. *Nature Reviews Genetics* **14**, 483-495 (2013).
- 647 33. Devanathan, S., Erban, A., Perez-Torres Jr, R., Kopka, J. & Makaroff, C.A. *Arabidopsis*
648 *thaliana* glyoxalase 2-1 is required during abiotic stress but is not essential under normal
649 plant growth. *PLoS One* **9**, e95971 (2014).
- 650 34. Lu, X. *et al.* A *PP2C-1* allele underlying a quantitative trait locus enhances soybean 100-
651 seed weight. *Molecular Plant* **10**, 670-684 (2017).
- 652 35. Li, N., Xu, R. & Li, Y. Molecular networks of seed size control in plants. *Annual review of*
653 *plant biology* **70**, 435-463 (2019).
- 654 36. Zhang, R. *et al.* Evolution of disease defense genes and their regulators in plants.
655 *International Journal of Molecular Sciences* **20**, 335 (2019).
- 656 37. Rojas, W. *et al.* Quinoa genetic resources and ex situ conservation. (2015).
- 657 38. Milner, S.G. *et al.* Genebank genomics highlights the diversity of a global barley collection.
658 *Nature Genetics* **51**, 319-326 (2019).
- 659 39. Wu, D. *et al.* Whole-genome resequencing of a worldwide collection of rapeseed accessions
660 reveals the genetic basis of ecotype divergence. *Molecular plant* **12**, 30-43 (2019).
- 661 40. Jia, G. *et al.* A haplotype map of genomic variations and genome-wide association studies of
662 agronomic traits in foxtail millet (*Setaria italica*). *Nature genetics* **45**, 957-961 (2013).
- 663 41. Varshney, R.K. *et al.* Whole-genome resequencing of 292 pigeonpea accessions identifies
664 genomic regions associated with domestication and agronomic traits. *Nature Genetics* **49**,
665 1082 (2017).
- 666 42. Zhou, Z. *et al.* Resequencing 302 wild and cultivated accessions identifies genes related to
667 domestication and improvement in soybean. *Nature biotechnology* **33**, 408-414 (2015).
- 668 43. Mather, K.A. *et al.* The extent of linkage disequilibrium in rice (*Oryza sativa* L.). *Genetics*
669 **177**, 2223-2232 (2007).
- 670 44. Xu, X. *et al.* Resequencing 50 accessions of cultivated and wild rice yields markers for
671 identifying agronomically important genes. *Nature biotechnology* **30**, 105-111 (2012).
- 672 45. Huang, X. *et al.* Genome-wide association studies of 14 agronomic traits in rice landraces.
673 *Nature genetics* **42**, 961 (2010).

- 674 46. Zhao, G. *et al.* A comprehensive genome variation map of melon identifies multiple
675 domestication events and loci influencing agronomic traits. *Nature Genetics* **51**, 1607-1615
676 (2019).
- 677 47. Choi, Y.-J. *et al.* Morphological and molecular characterization of the causal agent of
678 downy mildew on quinoa (*Chenopodium quinoa*). *Mycopathologia* **169**, 403-412 (2010).
- 679 48. Colque-Little, C.X. *et al.* Genetic variation for tolerance to the downy mildew pathogen
680 *Peronospora variabilis* in genetic resources of quinoa (*Chenopodium quinoa*). *bioRxiv*
681 (2020).
- 682 49. Chakravarty, T. & Sopory, S. Blue light stimulation of cell proliferation and glyoxalase I
683 activity in callus cultures of *Amaranthus paniculatus*. *Plant science* **132**, 63-69 (1998).
- 684 50. Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting Linear Mixed-Effects Models Using
685 lme4. *J Stat Soft* **67**, 48 (2015).
- 686 51. Lê, S., Josse, J. & Husson, F. FactoMineR: an R package for multivariate analysis. *Journal*
687 *of statistical software* **25**, 1-18 (2008).
- 688 52. Bolger, A.M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
689 sequence data. *Bioinformatics* **30**, 2114-2120 (2014).
- 690 53. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows–Wheeler
691 transform. *Bioinformatics* **26**, 589-595 (2010).
- 692 54. Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078-
693 2079 (2009).
- 694 55. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing
695 next-generation DNA sequencing data. *Genome research* **20**, 1297-1303 (2010).
- 696 56. Van der Auwera, G.A. *et al.* From FastQ data to high-confidence variant calls: the genome
697 analysis toolkit best practices pipeline. *Current protocols in bioinformatics* **43**, 11.10. 1-
698 11.10. 33 (2013).
- 699 57. Li, H. Tabix: fast retrieval of sequence features from generic TAB-delimited files.
700 *Bioinformatics* **27**, 718-719 (2011).
- 701 58. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156-2158
702 (2011).
- 703 59. Kalyanamoorthy, S., Minh, B.Q., Wong, T.K., von Haeseler, A. & Jermini, L.S.
704 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nature methods* **14**,
705 587 (2017).
- 706 60. Hoang, D.T., Chernomor, O., Von Haeseler, A., Minh, B.Q. & Vinh, L.S. UFBoot2:
707 improving the ultrafast bootstrap approximation. *Molecular Biology and Evolution* **35**, 518-
708 522 (2017).
- 709 61. Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool for the display and
710 annotation of phylogenetic and other trees. *Nucleic Acids Research* **44**, W242-W245 (2016).

- 711 62. Ruden, D.M. *et al.* Using *Drosophila melanogaster* as a model for genotoxic chemical
712 mutational studies with a new program, SnpSift. *Frontiers in genetics* **3**, 35 (2012).
- 713 63. Wang, W. *et al.* Genomic variation in 3,010 diverse accessions of Asian cultivated rice.
714 *Nature* **557**, 43-49 (2018).
- 715 64. Alexander, D.H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in
716 unrelated individuals. *Genome research* **19**, 1655-1664 (2009).
- 717 65. Jakobsson, M. & Rosenberg, N.A. CLUMPP: a cluster matching and permutation program
718 for dealing with label switching and multimodality in analysis of population structure.
719 *Bioinformatics* **23**, 1801-1806 (2007).
- 720 66. Zheng, X. *et al.* A high-performance computing toolset for relatedness and principal
721 component analysis of SNP data. *Bioinformatics* **28**, 3326-3328 (2012).
- 722 67. Zhang, C., Dong, S.-S., Xu, J.-Y., He, W.-M. & Yang, T.-L. PopLDdecay: a fast and
723 effective tool for linkage disequilibrium decay analysis based on variant call format files.
724 *Bioinformatics* **35**, 1786-1788 (2018).
- 725 68. Kang, H.M. *et al.* Variance component model to account for sample structure in genome-
726 wide association studies. *Nature Genetics* **42**, 348-354 (2010).
- 727 69. Li, M.-X., Yeung, J.M., Cherny, S.S. & Sham, P.C. Evaluating the effective numbers of
728 independent tests and significant p-value thresholds in commercial genotyping arrays and
729 public imputation reference datasets. *Human genetics* **131**, 747-756 (2012).
- 730 70. Turner, S.D. qqman: an R package for visualizing GWAS results using QQ and manhattan
731 plots. *Biorxiv*, 005165 (2014).

732

733 **Supplementary data**

734 **Supplementary tables**

735 **Supplementary Table 1:** Accessions from the quinoa diversity panel and results from re-
736 sequencing

737 **Supplementary Table 2:** Summary of high-quality SNPs identified in quinoa accessions

738 **Supplementary Table 3:** Variance components analysis of 12 quantitative traits

739 **Supplementary Table 4:** Summary of marker trait associations (MTA)

740 **Supplementary Table 5:** Candidate genes linked to SNP with significant trait associations

741 **Supplementary Table 6:** Summary of MTA associated with DTF, DTM, PD and PH identified on
742 chromosome Cq2A

743 **Supplementary Table 7:** Candidate genes located within the 50kb flanking regions of significantly
744 associated SNPs from the multivariate GWAS analysis

745 **Supplementary figures**

746 **Supplementary Fig. 1:** Geographical origin of the accessions forming the quinoa diversity panel.

747 **Supplementary Fig. 2:** Overview of the field experiment and exemplary images demonstrating
748 phenotypic traits; (A) and (B): Overview of the field and phenotypic variation among accession;
749 (C): Bolting (BBCH51) and (D) flowering (BBCH60) stage; Glomerulate (E) and amarantiform (F)
750 panicle shapes; red (G) and green (H) stem color ; red (I) and green (J) flower/inflorescence;
751 Growth type 1 (K) and type 5 (L); (M): Plant height and maturity variation between two accessions.

752 **Supplementary Fig. 3:** SNP density heat map across the 18 quinoa chromosomes. Different colors
753 depict SNP density.

754 **Supplementary Fig. 4:** Chromosome wide LD decay in genome A (A) and genome B (B). Colors
755 are depicting different chromosomes. (C) Genome-wide average LD decay of the A sub-genome
756 (blue) and B sub-genome (red).

757 **Supplementary Fig. 5:** SNP based PCA across all 18 quinoa chromosomes. Red circles are
758 depicting the two clusters of Lowland accessions.

759 **Supplementary Fig. 6** (A) ADMIXTURE ancestry coefficients for K ranging from 3 to 7 and 9.
760 Each vertical bar represents an accession, and color proportions on the bar correspond to the genetic
761 ancestry. (B) Cross-validation error in ADMIXTURE run.

762 **Supplementary Fig. 7:** Diversity of populations along chromosomes measured based on 10 kb non-
763 overlapping windows. Nucleotide diversity (π) distribution of 10 kb windows in population
764 Highland (A) and Lowland (B). (C) Nucleotide diversity ratios (π Lowland/ π Highland). (D)
765 Pairwise genome-wide fixation index (F_{ST}) between Highland and Lowland. The broken horizontal
766 line represents the top 1% threshold.

767 **Supplementary Fig. 8:** Distribution of Tajima's D along chromosomes in Highland (B) and
768 Lowland (D) populations. Density distribution of Tajima's D between populations. Different colors
769 represent the quartiles.

770 **Supplementary Fig. 9:** Graphical presentation of correlations between years among 12 traits.
771 Pearson correlation value (R) with P -values are shown. DTB: days to bolting (inflorescence
772 emergence), DTF: days to flowering, DTB to DTF: days between bolting and flowering, DTM;
773 days to maturity, PH: plant height (cm), PL: panicle length (cm), PD: panicle density (cm), NoB:
774 Number of branches, STL: stem lying, Saponin: saponin content as foam height (mm), Seed yield:
775 seed yield per plant (g), TSW: thousand seed weight (g),

776 **Supplementary Fig. 10:** Pearson correlations among 12 quinoa traits. Best linear unbiased
777 estimates across two years were used. Below the diagonal, scatter plots are shown with the fitted
778 line in red. Above the diagonal, the Pearson correlation coefficients are shown with significance
779 levels, *** = $P < 0.001$, ** = $P < 0.01$.

780 **Supplementary Fig. 11:** PCA of 12 quantitative phenotypes. A: Individual factor map colored
781 according to populations identified from SNP analysis. B: Variables factor map of the PCA.

782 **Supplementary Fig. 12:** Manhattan plots from GWAS with data from 2018 (left), 2019 (center),
783 and the mean of both years (right): The blue horizontal line indicates the suggestive threshold -

784 $\log_{10}(8.98E-7)$. The red horizontal line indicates the significant threshold (Bonferroni correction) -
785 $\log_{10}(1.67e-8)$.

786 **Supplementary Fig. 13:** Quantile-quantile plots of GWAS in two years, 2018 (left) and 2019
787 (center), and BLUE (right).

788 **Supplementary Fig. 14:** Local Manhattan plots for (A) flower color, (B) saponin content, and (C)
789 mildew infection. Candidate genes are shown in the color legend. LD heat maps are placed at the
790 Bottom. The colors of the heat map represent the pairwise correlation between individual SNPs.

791 **Supplementary Fig. 15:** PCA of 4 quantitative traits (DTF, DTM, PH, and PL). A: Individual
792 factor map, B: variables factor map of the PCA, C: distribution of the first three principal
793 components which were used for GWAS analysis.

794 **Supplementary Fig. 16:** GWAS analysis of principal components, PC1 (A), PC2 (B), PC3 (C):
795 Manhattan plots (left), and quantile-quantile plots (right): The blue horizontal line in the Manhattan
796 plots indicates the suggestive threshold $-\log_{10}(8.98E-7)$. The red horizontal line indicates the
797 significance threshold (Bonferroni correction) $-\log_{10}(1.67e-8)$.

798 **Supplementary Fig. 17:** Haplotypes of two genes, *CqPP2C* and *CqRING* controlling seed size in
799 quinoa. Geographic origin of the accessions and haplotype networks are displayed below the gene
800 structure.

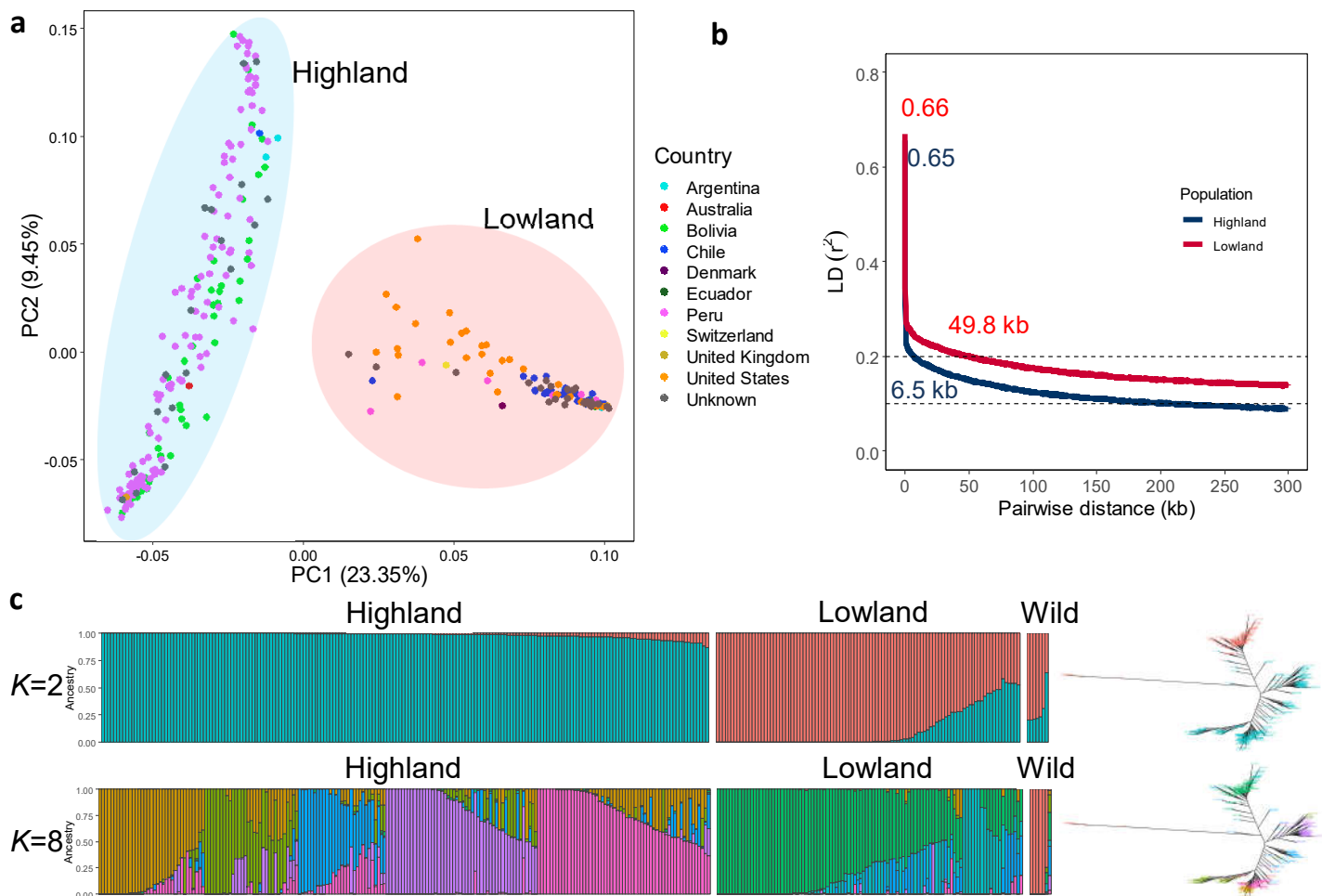


Fig. 1: Genetic diversity and population structure of the quinoa diversity panel. (a) PCA of 303 quinoa accessions. PC1 and PC2 represent the first two components of analysis, accounting for 23.35% and 9.45% of the total variation, respectively. The colors of dots represent the origin of accessions. Two populations are highlighted by different colors: Highland (light blue) and Lowland (pink). (b) Subpopulation wise LD decay in Highland (blue) and Lowland population (red). (c) Population structure is based on ten subsets of SNPs, each containing 50,000 SNPs from the whole-genome SNP data. Model-based clustering was done in ADMIXTURE with different numbers of ancestral kinships ($K=2$ and $K=8$). $K=8$ was identified as the optimum number of populations. Left: Each vertical bar represents an accession, and color proportions on the bar correspond to the genetic ancestry. Right: Unrooted phylogenetic tree of the diversity panel. Colors correspond to the subpopulation.

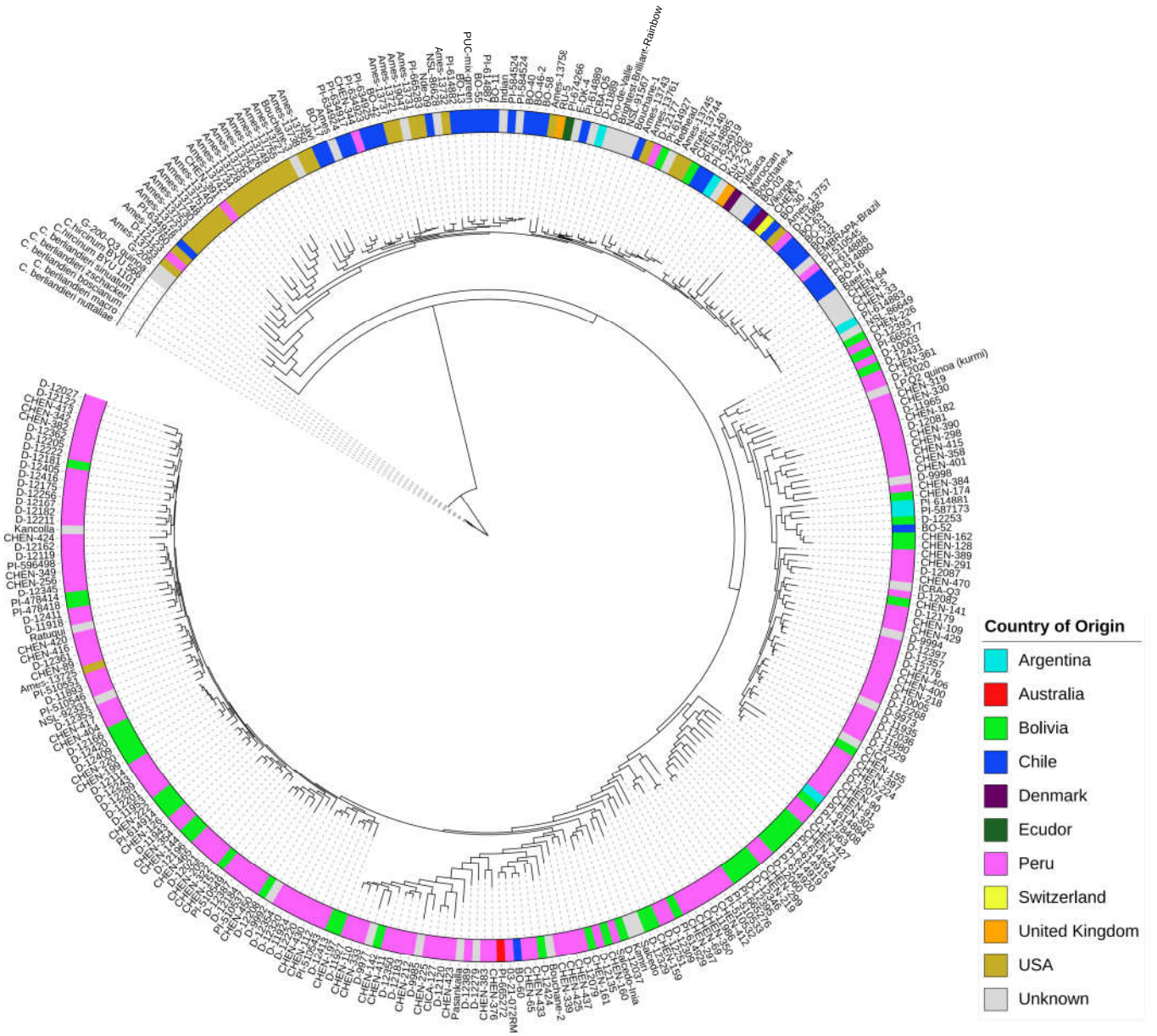
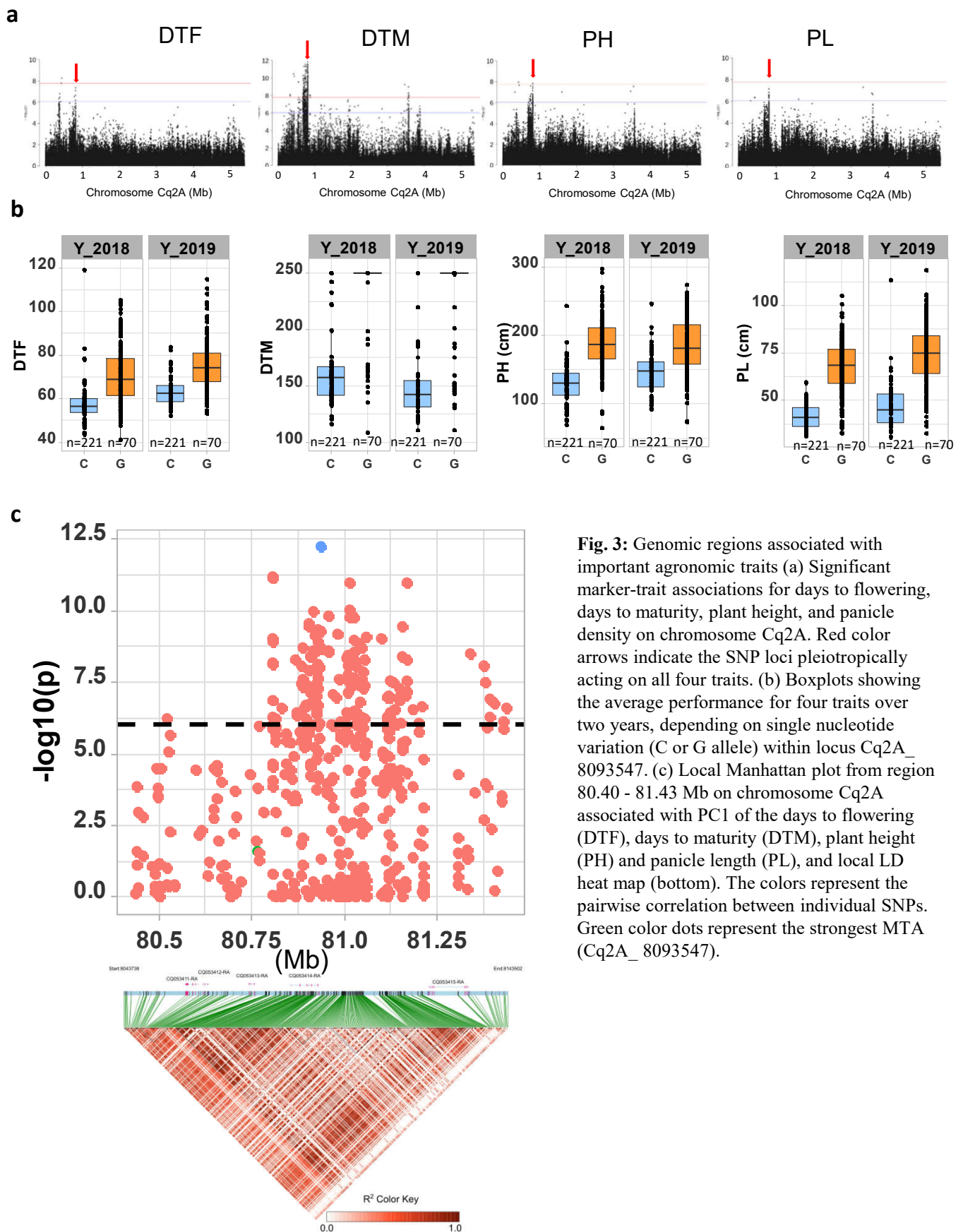


Fig. 2: Maximum likelihood tree of 303 quinoa and seven wild *Chenopodium* accessions from the diversity panel. Colors are depicting the geographical origin of accessions.



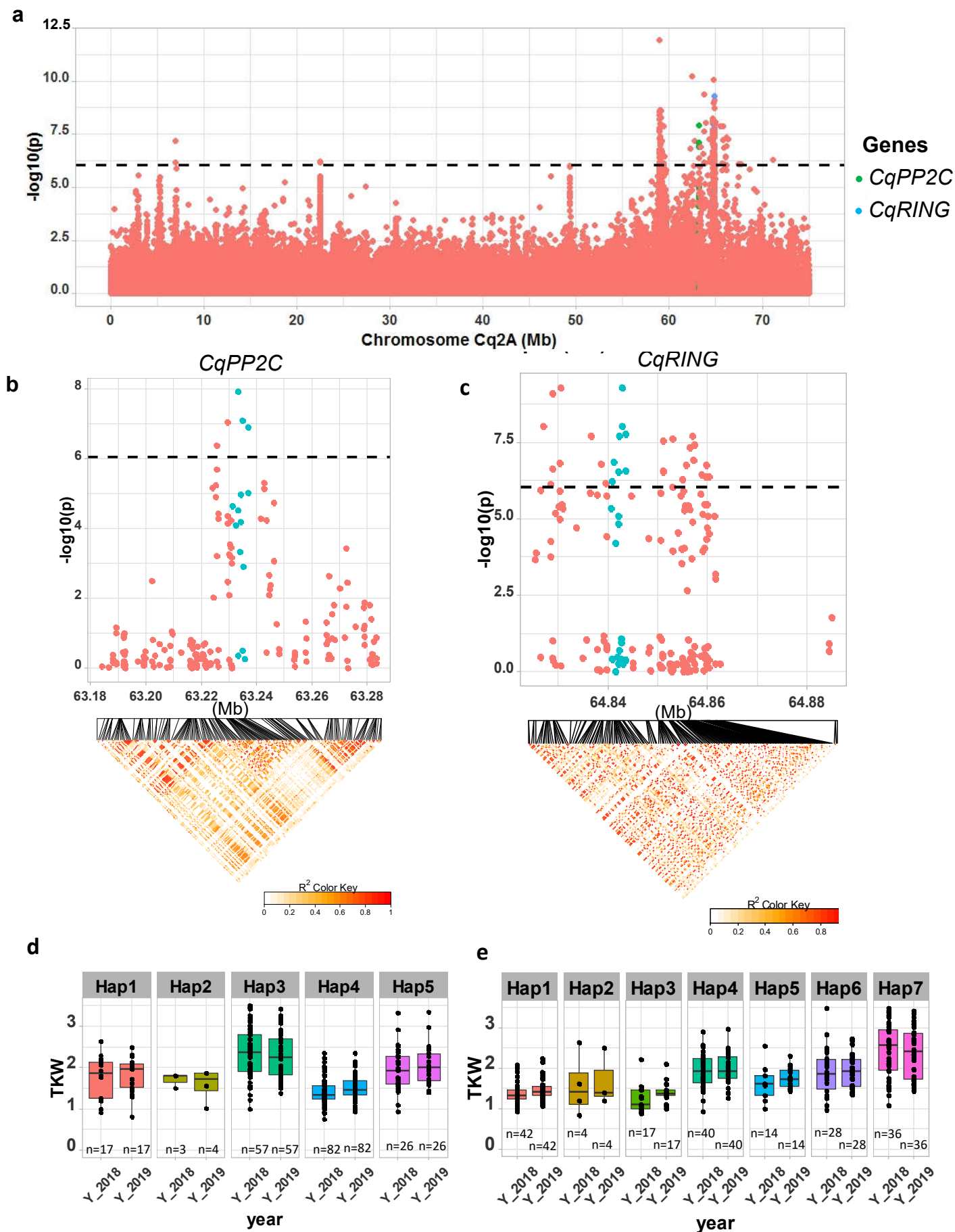


Fig. 4: Identification of candidate genes for thousand seed weight. (a) Manhattan plot from chromosome Cq8B. Green and blue dots are depicting the *CqPP2C5* and the *CqRING* gene, respectively. (b) Top: Local Manhattan plot in the neighborhood of the *CqPP2C* gene. Bottom: LD heat map. (c) Top: Local Manhattan plot in the neighborhood of the *CqRING* gene. Bottom: LD heat map. Differences in thousand seed weight between five *CqPP2C* (d) and seven *CqRING* haplotypes (e).