

1 **Title**

2 The genome of *Chenopodium quinoa*

3

4 **Authors**

5 David E. Jarvis*, Yung Shwen Ho*, Damien J. Lightfoot*, Sandra M. Schmöckel*, Bo Li*, Theo Borm,
6 Hajime Ohyanagi, Katsuhiko Mineta, Craig T. Mitchell, Noha Saber, Najeh M. Kharbatia, Ryan R. Rupper,
7 Aaron R. Sharp, Nadine Dally, Berin A Boughton, Yong H. Woo, Ge Gao, Elio Schijlen, Xiujie Guo, Afaque
8 A. Momin, Sónia Negrão, Salim Al-Babili, Christoph Gehring, Ute Roessner, Christian Jung, Kevin Murphy,
9 Stefan T. Arold, Takashi Gojobori, C. Gerard van der Linden, Eibertus N. van Loo, Eric N. Jellen, Peter J.
10 Maughan, Mark Tester

11

12 *These authors contributed equally to this work

13

14 **Summary**

15 *Chenopodium quinoa* (quinoa) is a highly nutritious grain crop with high abiotic stress tolerance that has
16 been identified as an important crop to improve world food security; unfortunately, few resources are
17 available to facilitate its genetic improvement. Here we report the assembly of a high-quality,
18 chromosome-scale reference genome sequence for quinoa, which was produced using single-molecule
19 real-time sequencing in combination with optical, chromosome contact, and genetic maps. We also
20 report the sequencing and assembly of two diploids from among the ancestral gene pools of quinoa –
21 which enabled the identification of sub-genomes in quinoa – and reduced-coverage genome sequences
22 for 22 other accessions of the allotetraploid goosefoot complex. The genome sequence facilitated the
23 identification of the transcription factor likely to control the production of anti-nutritional triterpenoid
24 saponins found in quinoa seeds, including a mutation that appears to cause alternative splicing and the
25 inclusion of a premature stop codon, thereby inactivating the protein and leading to the absence of
26 saponins in sweet quinoa accessions. These genomic resources are an important first step towards the
27 genetic improvement of quinoa to help increase global food security in the face of climate change and a
28 growing world population.

29 Almost 800 million people are undernourished, in part due to limited food production on marginal
30 lands¹. These poor agricultural environments are often disproportionately affected by climate change,
31 pressure from population growth, and intensification of agriculture, all of which exacerbate the
32 challenge of producing sufficient and nutritious food². Meeting the caloric and nutritional demands of
33 these growing populations will not only require increases in overall food production, but also the
34 development of new crops that can be grown sustainably in agricultural environments that are
35 increasingly susceptible to degradation.

36
37 Quinoa (*Chenopodium quinoa* Willd., $2n = 4x = 36$) is a highly nutritious crop that is adapted to thrive in
38 a wide range of agroecosystems. The centre of diversity of quinoa is near Lake Titicaca on the border of
39 Bolivia and Peru, where it has adapted to the arid and saline soils of the high plains of the Andean
40 Altiplano (>3,500 meters above sea level). Quinoa was presumably first domesticated more than 7,000
41 years ago by Pre-Columbian cultures and was known as the “Mother Grain” of Incan Empire³. Quinoa
42 was able to spread with pre-Incan cultures due to its adaptability to diverse environments, such as the
43 Mediterranean climates of coastal Chile as well as the intermediate highland valleys of Peru and
44 Ecuador. Abiotic stresses to which it is adapted include soil salinity^{4,5}, frost, high UV irradiance, and
45 drought⁶. Recently, quinoa has gained international attention because of the nutritional value of its
46 seeds, which are gluten-free, have a low glycaemic index⁷, and contain an excellent balance of essential
47 amino acids, fibre, lipids, carbohydrates, vitamins, and minerals⁸. Thus, quinoa has the potential to
48 provide a highly nutritious food source that can be grown on marginal lands not currently suitable for
49 other major crops. This potential was recognised when the United Nations declared 2013 as the
50 International Year of Quinoa (<http://www.fao.org/quinoa-2013/what-is-quinoa/nutritional-value/en/>),
51 this being one of only three times plants have received such a designation.

52
53 Despite its agronomic potential, quinoa is still an underutilised crop⁹, with relatively few active breeding
54 programs¹⁰. Breeding efforts to improve the crop for important agronomic traits, including tolerance to
55 heat and biotic stress, are needed to expand quinoa production worldwide and to relieve stress on
56 Andean production environments strained by intensive cultivation and the effects of climate change. To
57 accelerate the improvement of quinoa, we present here the allotetraploid quinoa genome, which was
58 sequenced and assembled using long-read sequencing technology in combination with optical,
59 chromosome contact, and genetic maps. We demonstrate the utility of the genome sequence by
60 identifying genes that underlie the production of betalain pigments and seed triterpenoid saponin
61 content. Moreover, we sequenced the genomes of additional diploid and tetraploid *Chenopodium*
62 species to characterise genetic diversity within the primary germplasm pool for quinoa and to
63 understand sub-genome evolution in quinoa. Together, these resources provide the foundation for
64 accelerating the genetic improvement of quinoa, with the objective of enhancing global food security for
65 a growing world population.

66 **Sequencing, assembly, and annotation**

67
68 We sequenced and assembled the genome of the coastal Chilean quinoa accession PI 614886 using
69 single-molecule real-time (SMRT) sequencing technology from Pacific Biosciences (PacBio). After
70 removal of contaminant and chloroplast sequences, the assembly consisted of 4,232 contigs, with a
71 contig N50 of 1.66 megabases (Mb) (Table 1). To further improve this assembly, we generated genome
72 maps using optical mapping technology¹¹ from BioNano Genomics and chromosome contact from
73 Dovetail Genomics¹². Scaffolding of the PacBio, BioNano, and Dovetail assemblies produced a final
74 assembly consisting of 3,483 scaffolds, with a scaffold N50 of 3.84 Mb and 90% of the assembled
75 genome contained in 439 scaffolds. The total assembly size of 1.39 gigabases (Gb) is similar to the
76 reported size estimates of the quinoa genome (1.45 – 1.50 Gb^{13,14}). To combine scaffolds into

77 pseudomolecules, an existing linkage map from quinoa¹⁵ was integrated with two new linkage maps
78 generated from independent populations. The resulting map (Extended Data Fig. 1) of 6,403 unique
79 markers spanned a total length of 2,034 centimorgans (cM) and consisted of 18 linkage groups (LGs)
80 (Supplementary Table 7), corresponding to the haploid chromosome number of quinoa; we refer to the
81 LGs hereafter as chromosomes, which are numbered based on numbering from a previously published
82 SNP linkage map¹⁵. The mapped markers came from 565 scaffolds, representing 1.18 Gb (85%) of the
83 total assembly length (Table 1, Supplementary File 1, Supplementary File 2).

84
85 Predicted protein-coding and miRNA genes (Supplementary Table 4) were annotated using a
86 combination of *ab initio* prediction and transcript evidence gathered from RNA sequenced from multiple
87 tissues using both RNA-Seq and Iso-Seq approaches. The inclusion of full-length reads from Iso-Seq
88 helped resolve several gene models (Extended Data Fig. 2a). The annotation contains 44,776 gene
89 models (Supplementary Table 2, Extended Data Fig. 2b), including 33,365 genes with annotation edit
90 distance (AED)^{16,17} values ≤ 0.3 (Extended Data Fig. 2c). Although the number of annotated genes is
91 greater than that of most sequenced diploid species, it is in line with sequenced tetraploid species¹⁸ and
92 is therefore likely a result of the presence of two sub-genomes in quinoa. Sixty-four percent of the
93 genome was found to be repetitive, including a large proportion of long terminal repeat (LTR)
94 transposable elements (Supplementary Table 1). To assess the completeness of the genome annotation,
95 we evaluated the coverage of a set of 956 highly conserved core eukaryotic genes found in a wide range
96 of plant taxa¹⁹. We identified 97.3% of the genes in the Plantae BUSCO dataset (Supplementary Table 3),
97 which is suggestive of a complete assembly and annotation. Eighty-seven percent of the single-copy
98 orthologs used in the BUSCO analysis were duplicated in the quinoa genome, which is likely a reflection
99 of the tetraploid nature of the quinoa genome.

100
101 The utility of the assembly, linkage maps, and annotation was demonstrated by mapping the betalain
102 locus, which controls stem pigmentation in one of the segregating mapping populations and is often
103 used as a morphological marker in breeding programs. The phenotype segregated as a single gene in the
104 F2 progeny (70 red, 22 green), and scoring stem colour in 51 F3 individuals enabled mapping of the trait
105 to chromosome 2 (CqA02), where it mapped to the same position as a SNP marker from Scaffold 1995
106 (3,473,993 bp) (Supplementary File 1). This SNP causes an amino acid change (Ala to Gly) in an
107 annotated peroxidase gene (AUR62012343) in the pigmented parent, and expression of the gene is
108 significantly lower in the pigmented compared to non-pigmented progeny (Supplementary File 3).
109 Peroxidase is known to regulate the stability of betalain pigments²⁰. An additional 65 candidate genes lie
110 within a 1-Mb window surrounding the mapped SNP marker (Supplementary File 4), including four other
111 peroxidase genes, and a gene (AUR62012346) annotated as being homologous to *CYP76AD1*, which
112 encodes for a cytochrome P450 which has also been shown to be required for production of the red
113 betalain pigment in *Beta vulgaris* (sugar beet)²¹, a member of the same family (Amaranthaceae) as
114 quinoa. RNA-Seq analysis showed that this gene is expressed at significantly higher levels in pigmented
115 plants than in non-pigmented plants (Supplementary File 3). Fine mapping is now required to enable
116 confident identification of the causative mutation and separation of primary and pleiotropic effects.

117 118 **Evolutionary history of quinoa**

119 Quinoa is an allotetraploid that resulted from the hybridisation of ancestral A- and B-genome diploid
120 species²². Single-gene sequencing studies previously identified pools of North American and Eurasian
121 diploids, respectively, as candidate sources of the A and B sub-genomes^{23–25}, with hybridisation
122 occurring somewhere in North America, in a scenario similar to that of cotton²⁶. To better understand
123 genome structure and evolution in quinoa, we sequenced, assembled, and annotated the A-genome
124 diploid *C. pallidicaule* and the B-genome diploid *C. suecicum*²⁴ (Fig. 1a, Table 1). We note that C.

125 *pallidicaule* (commonly called cañahua or kañiwa) is itself a regionally important and highly nutritious
126 alternative crop species native to the Andean Altiplano. A high proportion of orthologous gene pairs in
127 quinoa showed similar rates of synonymous substitutions per synonymous site (Ks), indicative of a
128 whole-genome duplication event (Fig. 1b). This most likely represents the hybridisation of ancestral
129 diploid species, because a similar peak was not observed in *C. pallidicaule* or *C. suecicum* (Fig. 1b). Using
130 mutation rates calculated for *Arabidopsis thaliana*²⁷ and for core eukaryotes²⁸, we estimate the
131 tetraploidisation to have occurred 3.3 – 6.3 million years ago.

132
133 Multiple tetraploid species – or ecotypes of the biological species complex – have arisen from the
134 ancestral tetraploid following hybridisation, including *C. berlandieri* and *C. hircinum*, which, together
135 with quinoa, are commonly referred to as the allotetraploid goosefoot complex (ATGC). Long-range
136 dispersal of the ancestral weedy ATGC eventually resulted in the appearance and domestication of
137 quinoa in the Lake Titicaca Basin over 7,000 years ago^{3,29}. Five ecotypes of quinoa have been previously
138 described (Altiplano, Coastal, Salares, Inter-Andean Valley, and Yungas)³⁰, although quinoa is most
139 commonly grouped simply into highland and coastal types²⁹. In addition to quinoa, the ATGC was the
140 source of at least three other cultigens: Mesoamerican vegetable *huauzontle* and pseudocereal *chia roja*
141 (*C. berlandieri* ssp. *nuttalliae*³¹); and the extinct staple pseudocereal of the North American Hopewell
142 Culture, *C. berlandieri* ssp. *jonesianum*³². The evolutionary relationships among the ATGC taxa, including
143 quinoa, remain unclear²⁹. To begin to resolve these issues, we re-sequenced 15 additional quinoa
144 accessions. Together with the reference accession, PI 614886, these 16 sequenced quinoa accessions
145 represent four major cultivated ecotypes (Supplementary File 5). We also sequenced five accessions of
146 *C. berlandieri* representing four of its major ecotypes and one accession each of *C. hircinum* from the
147 Pacific and Atlantic Andean watersheds (Supplementary File 5). Phylogenetic analysis of these taxa
148 suggests that *C. hircinum* and quinoa have a parent-derivative relationship, and that North American *C.*
149 *berlandieri* is the basal member of the ATGC with Mesoamerican *huauzontle* as an independent
150 domesticate (Fig. 1c). Quinoa was thought to have been domesticated from *C. hircinum* in a single event
151 from which coastal quinoa was later derived (Fig. 1d, arrow 1); however, our sequencing data place a *C.*
152 *hircinum* accession basal to coastal ecotypes (Fig. 1c), suggesting the possibility that quinoa was
153 domesticated independently in highland and coastal environments (Fig. 1d, arrows 2a and 2b,
154 respectively). Future analyses with deeper sampling of quinoa and *C. hircinum* will help clarify the
155 relationship of *C. hircinum* with highland and coastal quinoa ecotypes, as well as provide germplasm for
156 breeding broadly adapted coastal quinoa cultivars for warm-season production. The single-nucleotide
157 polymorphisms (SNPs) identified between these accessions and the reference quinoa genome – a total
158 of 7,809,381 (Extended Data Fig. 3, Supplementary Table 5), including 2,668,694 that are specific to
159 quinoa – will be useful in assessing genetic diversity and identifying genomic regions associated with
160 desirable traits.

161 **Analysis of sub-genome structure**

162 The availability of the *C. pallidicaule* and *C. suecicum* genome sequences enabled the identification of
163 chromosomes belonging to the A and B sub-genomes of quinoa. By mapping sequencing reads from *C.*
164 *pallidicaule* and *C. suecicum* onto the quinoa scaffold assembly, and by performing BLASTN searches of
165 each diploid against the quinoa assembly, 156 and 410 scaffolds (totalling 202.6 and 646.3 Mb) were
166 assigned to the A and B sub-genomes, respectively (Fig. 2a, Supplementary File 6). A mini-satellite
167 repeat (18-24J) previously shown to be more abundant in the B sub-genome of quinoa and related
168 species³³ is over-represented in scaffolds assigned to the B sub-genome (Supplementary File 6). Nine
169 chromosomes were assigned to each sub-genome based on the sub-genome assignments of scaffolds
170 contained in each chromosome (chromosomes hereafter designated as CqA or CqB, followed by the
171 chromosome number). On a chromosome level, the B sub-genome accounted for a larger percentage of
172

173 both the genetic (1087 cM) and physical (696 Mb) sizes of the quinoa genome than the A sub-genome
174 (946 cM, 488 Mb). This result was not unexpected, given the differences in the estimated genome sizes
175 of *C. pallidicaule* (452 Mb) and *C. suecicum* (815 Mb) were predicted based on k-mer analyses.
176

177 A total of 5,807 orthologous gene pairs for which one gene was phylogenetically more closely related to
178 an ortholog from *C. pallidicaule* and the other was phylogenetically more closely related to an ortholog
179 from *C. suecicum* were used to elucidate homoeologous relationships among the quinoa chromosomes.
180 Visualization of the chromosomal locations of these homoeologous gene pairs revealed a high degree of
181 synteny between the A and B sub-genomes (Fig. 2b). Interestingly, a small number of gene pairs (5.7%)
182 mapped within the same sub-genome, suggesting that recombination and chromosomal rearrangement
183 have occurred between the A and B sub-genomes. For example, we identified homoeologous A and B
184 sub-genome regions located in the B sub-genome chromosomes CqB05 and CqB03. The genes in the
185 region of ~54-56 Mb of CqB03 are phylogenetically more similar to the A-genome diploid *C. suecicum*
186 and therefore likely originated from the A sub-genome chromosome CqA12 (Fig. 2c). In addition to
187 exchanges between sub-genomes, we also see evidence of large chromosomal rearrangements within
188 sub-genomes, complicating the assignment of homoeologous chromosome pairs. For example, ortholog
189 analysis clearly identifies CqB05 and CqA12 as homoeologous, although the same analysis is much more
190 complicated with other chromosomes, including between CqB01, CqA02, and CqA04, and between
191 CqA07, CqB11, and CqB17 (Fig. 2b). To clarify these relationships, we identified syntenic regions
192 between chromosomes of the diploid *B. vulgaris*³⁴ ($n = 9$) and the A and B sub-genome chromosomes of
193 quinoa (Fig. 2d). These results indicate that CqA02 and CqA04 are orthologous to *B. vulgaris*
194 chromosomes 8 and 2 (Bvchr8 and Bvchr2), respectively, whereas CqB01 appears to be the result of a
195 chromosome fusion. Likewise, CqA07 appears to be the result of a fusion between ancestral
196 chromosomes orthologous to Bvchr3 and Bvchr7.
197

198 **Analysis of sub-genome content**

199 Duplicated genes are often lost following genome duplication events³⁵. We used OrthoMCL³⁶ to identify
200 clusters of orthologous genes in related species in the Amaranthaceae (Extended Data Fig. 4), and
201 specifically investigated the retention and loss of orthologous genes in quinoa and the three diploid
202 species *C. pallidicaule*, *C. suecicum*, and *B. vulgaris* (Fig. 3a, Supplementary Table 6). To identify genes
203 lost from either the A or B sub-genome of quinoa, we identified orthologous sets of genes for which only
204 one copy could be found in quinoa and in each of the diploid species. Of these, we found a similar
205 number (1,031 and 849) of genes lost from the A and B sub-genomes, respectively (Fig. 3b). To identify
206 duplicated genes retained as single-copy in each of the quinoa sub-genomes, we identified sets of genes
207 for which one copy was found in each of the diploids and two copies were found in quinoa. Of these,
208 5,807 sets contained one quinoa gene assigned to each of the A and B sub-genomes (Fig. 3b, Fig. 2b),
209 thus representing a core set of single-copy genes retained in each genome and sub-genome. We also
210 investigated genes retained in multiple copies. Notably, we found that quinoa, like *B. vulgaris*³⁷, contains
211 two genes that are orthologous to the *Arabidopsis thaliana* *FT* flowering-time regulator gene, with
212 quinoa containing two homoeologous copies of each ortholog due to its tetraploid nature (Fig. 3c,
213 Extended Data Fig. 5). *FT* is known to promote flowering in *A. thaliana*, and functional orthologs have
214 been found in other species³⁸; however, *B. vulgaris*, which requires vernalisation to induce the transition
215 to flowering, was found to contain a second *FT* gene that acts antagonistically by repressing flowering
216 and whose expression is reduced following vernalisation³⁷. In quinoa, one pair of homoeologous *FT*
217 genes was found to be orthologous to the *FT* gene that promotes flowering in *B. vulgaris*, and a second
218 pair of homoeologous genes was found to be orthologous to the gene that represses flowering in *B.*
219 *vulgaris*. The identification of two pairs of *FT* orthologs in quinoa, as well as the previous identification of
220 two *FT* genes in *C. rubrum*³⁹, suggests that the common ancestor of quinoa and *B. vulgaris* also had two

221 *FT* genes. Given the lack of a vernalisation requirement to induce flowering in quinoa, it is possible that
222 both *FT* genes have retained the normal *FT* function in quinoa, while the second *FT* gene has acquired a
223 repressive function in *B. vulgaris* following its divergence from quinoa. Future functional studies will
224 help address these questions, which are important for many aspects determining yield and yield
225 maintenance under abiotic stress.

226

227 **Mechanisms underlying saponin production**

228 Quinoa seeds contain a mixture of triterpene glycosides called saponins⁴⁰. Although saponins may be
229 beneficial for plant growth – for example, by deterring herbivory^{41,42} – they are also haemolytic and
230 produce a bitter flavour in quinoa seeds and must therefore be removed before human consumption.
231 Because this process is costly, is often water-intensive, and can reduce the nutritional value of the
232 seeds⁴³, the development of saponin-free lines is a major breeding objective in quinoa¹⁰.

233

234 Using imaging mass spectrometry (MS)⁴⁴, we show that saponins accumulate in the seed pericarp (Fig.
235 4a, 4b, Extended Data Fig. 6), which is maternal tissue derived from the ovary wall. We found that
236 saponins accumulate in the seed early in seed development, between 20 and 24 days after anthesis (Fig.
237 4c), eventually accounting for 4% (w/w) of the mature seed mass (Supplementary Information 8.1.).

238

239 Triterpenoid saponins are synthesised as part of the mevalonic acid-dependent pathway, which also
240 produces phytosterols⁴⁵. They are comprised of an aglycone backbone to which various sugar side chains
241 can be attached, including glucose, galactose, and arabinose. Almost 100 different saponins have been
242 identified in different accessions of quinoa^{40,46}. Using liquid chromatography–MS (LC-MS), we identified
243 and annotated 47 different saponins in the reference accession (Supplementary Table 9).

244

245 Naturally-occurring sweet quinoa accessions that contain very low levels of saponins are present within
246 the quinoa germplasm⁴⁷, although the underlying genes regulating the absence of saponins in these
247 lines are unknown. To identify these genes, we performed linkage mapping and bulk segregant analysis
248 (BSA) using two populations segregating for the presence of saponins in the seeds: Kurmi (sweet) × 0654
249 (bitter), and Atlas (sweet) × Carina Red (bitter). Consistent with reports from other populations⁴⁸,
250 segregation ratios in these populations indicated that the presence and absence of seed saponins is
251 controlled by a single gene, with the presence of saponins being dominant (71 bitter and 21 sweet in
252 Kurmi × 0654; 567 bitter and 175 sweet in Atlas × Carina Red). We note that qualitative differences exist
253 in the types of saponins identified in bitter lines (Extended Data Fig. 7, Supplementary Table 8) and that
254 the presence and absence of saponins was correlated with differences in seed coat thickness, with bitter
255 lines having significantly thicker seed coats than sweet lines (Extended Data Fig. 8).

256

257 To identify the gene controlling the absence of saponins in these populations, complementary
258 genotyping analyses were performed, including RNA-Seq using RNA extracted from clusters of flowers
259 and seeds of the parents and of 14 bitter and 15 sweet homozygous F3 lines (Kurmi × 0654 population;
260 Extended Data Fig. 7), and DNA-Seq performed using DNA extracted from leaves of the parents and 94
261 sweet F2 lines (Atlas × Carina Red population). Linkage mapping and BSA in each population identified
262 the same region on CqB16 that distinguishes the bitter and sweet lines (Fig. 5a). Frequencies of the
263 sweet allele in both populations reached 100% for markers located in CqB16 on Scaffold 3489. We
264 investigated the genes in a 700-kb window surrounding this region of 100% sweet allele frequency. Of
265 the 54 annotated genes in this region (Supplementary File 7), two are similar to genes previously shown
266 to play a role in saponin biosynthesis. Specifically, AUR62017204 and AUR62017206 are neighbouring
267 genes annotated as basic helix-loop-helix (bHLH) transcription factors sharing homology (Extended Data
268 Fig. 9a) with the class IVa *bHLH* genes that are known to regulate triterpene saponin biosynthesis in

269 *Medicago truncatula*⁴⁹. In *M. truncatula*, overexpression of the TRITERPENE SAPONIN BIOSYNTHESIS
270 ACTIVATING REGULATOR 1 (TSAR1) and TSAR 2 bHLH transcription factors was recently shown to
271 increase expression of genes in the saponin biosynthetic pathway, resulting in increased accumulation of
272 triterpene saponins⁴⁹. The overexpression of *TSAR1* and *TSAR2* elicited different patterns of
273 transactivation of downstream saponin biosynthetic genes in *M. truncatula*, suggesting distinct
274 functionalities within the triterpenoid biosynthetic pathway, including the accumulation of different
275 classes of saponins⁴⁹. Similarly, our RNA-Seq data also shows distinct expression patterns for these two
276 genes. Specifically, AUR62017206 (hereafter *C. quinoa TSAR-LIKE2*, *CqTSARL2*) was expressed in root
277 tissue but not in flowers or immature seeds, whereas AUR62017204 (hereafter *C. quinoa TSAR-LIKE1*,
278 *CqTSARL1*) was almost exclusively expressed in seeds, with significantly lower expression levels in sweet
279 lines (Supplementary File 8). Together, these results suggest that *CqTSARL1* might be a functional TSAR
280 ortholog, although whether this is due to shared ancestry or convergent evolution is unclear.

281
282 In *M. truncatula*, TSAR1 and TSAR2 were found to bind to the DNA motif 5'-CACGHG-3' (where H can be
283 A, C, or T)⁴⁹. We identified this motif within 2 kb upstream of the start codon in several saponin
284 biosynthetic pathway genes in quinoa, including genes encoding hydroxymethylglutaryl (HMG) CoA
285 synthase, HMG CoA reductase, mevalonate kinase, mevalonate-5-pyrophosphate decarboxylase,
286 isopentenyl diphosphate isomerase, farnesyl pyrophosphate synthase, squalene synthase, and squalene
287 epoxidase (Fig. 5b). Expression levels of these genes – as well as several other genes in the saponin
288 biosynthetic pathway – were significantly downregulated in sweet lines (Fig. 5b, Supplementary File 8).

289
290 Analysis of RNA-Seq reads in the sweet progeny of Kurmi and 0654 revealed that the *CqTSARL1*
291 transcript was alternatively spliced in exon 3 due to a SNP (G2078C) in the last position of the exon,
292 thereby altering the canonical intron/exon splice boundary (Fig. 5c). The SNP co-segregates with the
293 presence of saponins in our Kurmi × 0654 F3 population, and is likely the causative mechanism for the
294 alternative splicing at an upstream cryptic splice site in the sweet lines (Fig. 5c). This alternative splicing
295 of *CqTSARL1* results in a premature stop codon (Extended Data Fig. 9b) and a truncated protein that
296 modelling predicts to be compromised in its ability to form homodimers and/or to bind DNA (Extended
297 Data Fig. 9b-9d, Supplementary Information 8.8.), which are both necessary for regulation of
298 transcription. All bitter accessions in our re-sequencing pool share the same allele (G) found in the bitter
299 progeny of Kurmi and 0654, whereas all sweet accessions (Chucapaca, G205-95DK, Salcedo INIA, as well
300 as the mapping population parents Kurmi and Atlas) but one (Pasankalla) contain the same allele
301 (G2078C) as the sweet progeny. Interestingly, however, although the G2078C allele is present in the
302 sequenced Atlas line, none of the sweet progeny in the Atlas × Carina Red population were found to
303 have the G2078C allele. Additional sequencing of individual plants of the Atlas variety revealed a low
304 level of heterogeneity within the variety for the *CqTSARL1* gene, with some plants containing the
305 G20178C allele and others containing sequence insertions (Supplementary Information 7.2.4). Thus, it is
306 likely that the Atlas plant used in the cross with Carina Red – which, importantly, was not the same plant
307 used for sequencing – possessed a different mutation other than the G2078C allele. Indeed, we found
308 strong evidence of insertions in and around the *CqTSARL1* gene in all the sweet progeny of the Atlas ×
309 Carina Red population (Fig. 5c, Extended Data Fig. 10). In particular, two exonic insertions in *CqTSARL1* in
310 the sweet progeny are likely to inactivate the gene and result in a sweet phenotype. The identification of
311 multiple, independent mutations in *CqTSARL1* that co-segregate with the sweet phenotype strongly
312 suggests that this gene regulates the presence and absence of saponins in quinoa seeds. PCR markers
313 designed from sequences surrounding *CqTSARL1* on CqB16 are perfectly linked to the presence and
314 absence of saponins, and can now be used in marker-assisted selection (MAS) to accelerate the
315 development of sweet commercial quinoa varieties.

316

317 **Conclusions**

318 As an emerging international crop, quinoa has great potential to enhance global food security. The high-
319 quality reference genome assembly presented here will accelerate improvements of quinoa. Major
320 breeding objectives for quinoa improvement include the development of shorter plants with fewer
321 branches and more compact seed heads, increased heat and biotic stress tolerance, and the
322 introgression of the sweet phenotype into commercial varieties. The identification of the likely causative
323 mutation underlying the sweet phenotype not only provides insights into triterpenoid saponin
324 biosynthesis, but also enables accelerated breeding of sweet commercial varieties using marker-assisted
325 selection. The diversity present in the primary gene pool of quinoa, which we have begun to
326 characterise, will also help direct future breeding strategies. The resources presented here also help to
327 make quinoa a useful model for studying polyploid genome evolution and mechanisms of abiotic stress
328 tolerance, in particular salinity tolerance.

329 **Methods**

330 **Quinoa sequencing and assembly**

331 We sequenced *Chenopodium quinoa* Willd. (quinoa) accession PI 614886 (also known as NSL 106399 and
332 QQ74). DNA was extracted from leaf and flower tissue of a single plant, as described in the “Preparing
333 *Arabidopsis* Genomic DNA for Size-Selected ~20 kb SMRTbell™ Libraries” protocol
334 ([http://www.pacb.com/wp-content/uploads/2015/09/Shared-Protocol-Preparing-Arabidopsis-DNA-for-](http://www.pacb.com/wp-content/uploads/2015/09/Shared-Protocol-Preparing-Arabidopsis-DNA-for-20-kb-SMRTbell-Libraries.pdf)
335 [20-kb-SMRTbell-Libraries.pdf](http://www.pacb.com/wp-content/uploads/2015/09/Shared-Protocol-Preparing-Arabidopsis-DNA-for-20-kb-SMRTbell-Libraries.pdf)). DNA was purified twice with Beckman Coulter Genomics AMPure XP
336 magnetic beads and assessed by standard agarose gel electrophoresis and Thermo Fisher Scientific
337 Qubit Fluorometry. 100 Single-Molecule Real-Time (SMRT) cells were run on the PacBio RS II system
338 with the P6-C4 chemistry by DNALink (Seoul, Republic of Korea). *De novo* assembly was conducted using
339 the smrtmake assembly pipeline (<https://github.com/PacificBiosciences/smrtmake>) and the Celera
340 Assembler, and the draft assembly was polished using the quiver algorithm.

341
342 DNA was also sequenced using an Illumina HiSeq 2000 machine. For this, DNA was extracted from leaf
343 tissue of a single soil-grown plant using the Qiagen DNeasy Plant Mini Kit. 500-bp paired-end (PE)
344 libraries were prepared using the NEBNext Ultra DNA Library Prep Kit for Illumina. Sequencing reads
345 were processed with Trimmomatic (v0.33)⁵⁰, and reads < 75 nucleotides in length after trimming were
346 removed from further analysis. The remaining high-quality reads were assembled with Velvet (v1.2.10)⁵¹
347 using a k-mer of 75.

348
349 **Integrating BioNano optical maps with the PacBio assembly**

350 High molecular weight DNA was isolated and labelled from leaf tissue of three-week old quinoa plants
351 according to standard BioNano protocols, using the single-stranded nicking endonuclease Nt.BspQI.
352 Labelled DNA was imaged automatically using the BioNano Irys system and *de novo* assembled into
353 consensus physical maps using the BioNano IrysView analysis software. The final *de novo* assembly used
354 only single molecules with a minimum length of 150 kb and eight labels per molecule. PacBio-BioNano
355 hybrid scaffolds were identified using IrysView’s hybrid scaffold alignment subprogram.

356
357 **Chicago library preparation and sequencing**

358 Using the same DNA prepared for PacBio sequencing, a Chicago library was prepared as described
359 previously¹². The library was sequenced on an Illumina HiSeq 2500.

360
361 **Scaffolding the PacBio and BioNano assemblies with HiRise**

362 Chicago sequence data (in FASTQ format) was used to scaffold the PacBio-BioNano hybrid assembly
363 using HiRise, a software pipeline designed specifically for using Chicago data to assemble genomes¹².
364 Chicago library sequences were aligned to the draft input assembly using a modified SNAP read mapper
365 (<http://snap.cs.berkeley.edu>). The separations of Chicago read pairs mapped within draft scaffolds were
366 analysed by HiRise to produce a likelihood model, and the resulting likelihood model was used to
367 identify putative mis-joins and score prospective joins.

368
369 **Linkage mapping and genetic marker analyses**

370 **Kurmi × 0654 population.** A population was developed by crossing Kurmi (green, sweet) and 0654 (red,
371 bitter). Homozygous high- and low-saponin F2 lines were identified by planting 12 F3 seeds derived from
372 each F2 line, harvesting F4 seed from these F3 plants, and then performing foam tests on the F4 seed.
373 Phenotyping was validated using gas chromatography/mass spectrometry (GC/MS). RNA was extracted
374 from inflorescences containing a mixture of flowers and seeds at various stages of development from
375 the parents and 45 individual F3 progeny. RNA extraction and Illumina sequencing were performed as
376 described above. Sequencing reads from all lines were trimmed using Trimmomatic and mapped to the
377 reference assembly using TopHat⁵², and SNPs were called using SAMtools mpileup (v1.1)⁵³.

378
379 For linkage mapping, markers were assigned to linkage groups on the basis of the grouping by JoinMap
380 v4.1. Using the maximum likelihood algorithm of JoinMap, the order of the markers was determined;
381 using this as start order and fixed order, regression mapping in JoinMap was used to determine the cM
382 distances.

383
384 Genes differentially expressed between bitter and sweet lines and between green and red lines were
385 identified using default parameters of the Cuffdiff function of the Cufflinks program⁵⁴.

386
387 **Atlas x Carina Red population.** A second mapping population was developed by crossing Atlas (sweet)
388 and Carina Red (bitter). Bitter and sweet F2 lines were identified by performing foam and taste tests on
389 the F3 seed. DNA sequencing was performed with DNA from the parents and 94 sweet F2 lines, as
390 described above, and sequencing reads were mapped to the reference assembly using BWA. SNPs were
391 called in the parents and in a merged file containing all combined F2 lines

392
393 Genotype calls were generated for the 94 F2 genotypes by summing up read counts over a sliding
394 window of 500 variants, at all variant positions for which the parents were homozygous and
395 polymorphic. Over each 500-variant stretch, all reads with Atlas alleles were summed, and all reads with
396 the Carina Red allele were summed. Markers were assigned to linkage groups using JoinMap, with
397 regression mapping used to obtain the genetic maps per linkage group.

398
399 **Integrated linkage map.** The Kurmi × 0654 and Atlas × Carina Red maps were integrated with the
400 previously published quinoa linkage map¹⁵, with the Kurmi × 0654 map being used as the reference for
401 the positions of anchor markers and scaling. We selected markers from the same scaffold that were in
402 the same 10,000-bp bin in the assembly. The anchor markers on the alternative map received the
403 position of the Kurmi × 0654 map anchor marker in the integrated map. This process was repeated with
404 anchor markers at the 100,000-bp bin level. The assumption is that at the 100,000-bp bin level
405 recombination should essentially be zero. On this level, a regression of cM position on both maps
406 yielded $R^2 > 0.85$ and often > 0.9 , so the regression line can easily be used for interpolating the positions
407 of the alternative map towards the corresponding position on the Kurmi × 0654 map. All Kurmi × 0654
408 markers went into the integrated map on their original position.

409
410 **Chromosome pseudomolecules.** Pseudomolecules were assembled by concatenating scaffolds based on
411 their order and orientation as determined from the integrated linkage map, with 2,000 N's inserted

412 between each scaffold. Scaffolds that mapped to multiple linkage groups were assigned to the linkage
413 group to which the greatest number of markers from that scaffold were mapped. A custom Perl script
414 was used to convert the coordinates from the scaffold-level annotation.

415

416 **Sequencing and assembly of *C. pallidicaule* and *C. suecicum***

417 DNA was extracted from *C. pallidicaule* (PI 478407) and *C. suecicum* (BYU 1480) and was sent to the
418 Beijing Genomic Institute (BGI, Hong Kong, China) where one 180-bp PE library and two mate-pair
419 libraries with insert sizes of 3 and 6 kb were prepared and sequenced on the Illumina HiSeq platform to
420 obtain 2 X 100-bp reads for each library. The generated reads were trimmed using the quality-based
421 trimming tool Sickle (<https://github.com/najoshi/sickle>). The trimmed reads were then assembled using
422 the ALLPATHS-LG assembler⁵⁵, and GapCloser v1.12⁵⁶ was used to resolve N spacers and gap lengths
423 produced by the ALLPATHS-LG assembler.

424

425 **Genome annotation**

426 Repeat families found in the genome assemblies of quinoa, *C. pallidicaule*, and *C. suecicum* (see
427 Supplementary Information 3.) were first independently identified *de novo* and classified using the
428 software package RepeatModeler⁵⁷. RepeatMasker⁵⁸ was used to discover and identify repeats within
429 the respective genomes.

430

431 AUGUSTUS⁵⁹ was used for *ab initio* gene prediction, using model training based on coding sequences
432 from *Amaranthus hypochondriacus*, *Beta vulgaris*, *Spinacia oleracea* and *Arabidopsis thaliana*. RNA-Seq
433 and Iso-Seq reads generated from RNA of different tissues were mapped onto the reference genome
434 using Bowtie 2⁶⁰ and GMAP⁶¹, respectively. Hints with locations of potential intron-exon boundaries
435 were generated from the alignment files with the software package BAM2hints in the MAKER package⁶².
436 MAKER with AUGUSTUS (intron-exon boundary hints provided from RNA-Seq and Iso-Seq) was then
437 used to predict genes in the repeat-masked reference genome. To help guide the prediction process,
438 peptide sequences from *B. vulgaris* and the original quinoa full-length transcript (provided as EST
439 evidence) were used by MAKER during the prediction. Genes were characterised for their putative
440 function by performing a BLAST search of the peptide sequences against the UniProt database. PFAM
441 domains and InterProScan ID were added to the gene models using the scripts provided in the MAKER
442 package.

443

444 **Re-sequencing**

445 The following quinoa accessions were chosen for DNA re-sequencing: 0654, Ollague, Real, Pasankalla
446 (BYU 1202), Kurmi, CICA-17, Regalona (BYU 947), Salcedo INIA, G-205-95DK, Cherry Vanilla (BYU 1439),
447 Chucapaca, Ku-2, PI 634921 (Ames 22157), Atlas, and Carina Red. The following accessions of *C.*
448 *berlandieri* were sequenced: var. *boscianum* (BYU 937), var. *macrocalycium* (BYU 803), var. *zschackei*
449 (BYU 1314), var. *sinuatum* (BYU 14108), and subspecies *nuttaliae* ("Huazontle"). Two accessions of *C.*
450 *hircinum* (BYU 566 and BYU 1101) were also sequenced. All sequencing was performed with an Illumina
451 HiSeq 2000 machine, using either 125-bp (Atlas and Carina Red) or 100-bp (all other accessions) paired-
452 end libraries. Reads were trimmed using Trimmomatic and mapped to the reference assembly using

453 BWA (v0.7.10)⁶³. Read alignments were manipulated with SAMtools, and the mpileup function of
454 SAMtools was used to call SNPs.

455

456 **Identification of orthologous genes**

457 Orthologous and paralogous gene clusters were identified using OrthoMCL³⁶. Recommended settings
458 were used for all-against-all BLASTP comparisons (Blast+ v2.3.0⁶⁴) and OrthoMCL analyses. Custom Perl
459 scripts were utilised to process OrthoMCL outputs for visualisation with InteractiVenn⁶⁵.

460

461 **Phylogenetic inference**

462 Using OrthoMCL, orthologous gene sets containing two copies in quinoa and one copy each in *C.*
463 *pallidicaule*, *C. suecicum*, and *B. vulgaris* were identified. In total, 7,433 gene sets were chosen, and their
464 amino acid sequences were aligned individually for each set using MAFFT⁶⁶. The 7,433 alignments were
465 converted into PHYLIP format files by the seqret command in the EMBOSS package⁶⁷. Individual gene
466 trees were then constructed using the maximum likelihood method using proml in PHYLIP⁶⁸.

467

468 In addition, the genomic variants of all 25 sequenced taxa (Supplementary File 5) relative to the
469 reference sequence were called based on the mapped Illumina reads in 25 bam files using SAMtools. To
470 call variants in the reference accession (PI 614886), Illumina sequencing reads were mapped to the
471 reference assembly. Variants were then filtered using VCFtools⁶⁹ and SAMtools, and the qualified SNPs
472 were combined into a single VCF file which was used as an input into SNPhylo⁷⁰ to construct the
473 phylogenetic relationship using maximum likelihood and 1,000 bootstrap iterations.

474

475 To identify *FT* homologs, the protein sequence from *A. thaliana* flowering time gene *FT* was used as a
476 BLAST query. Filtering for hits with an E-value < 1e⁻³ and with RNA-Seq evidence resulted in the
477 identification of four quinoa proteins. For the construction of the phylogenetic tree, protein sequences
478 from these four quinoa *FT* homologs were aligned using MAFFT⁶⁶ (progressive method G-INS-1) along
479 with two *B. vulgaris* (gene models: BvFT1-miuf.t1, BvFT2-ewwx.t1), one *A. thaliana* (AT1G65480.1), and
480 one *O. sativa* (AFK31087.1) homolog. Phylogenetic analysis was performed with MEGA⁷¹ (v6.06). The
481 JTT+G model was selected as the best fitting model. The phylogenetic tree was estimated using the
482 maximum likelihood method with 1,000 bootstrap replicates. The syntenic relationships between the
483 coding sequences of the chromosomal regions surrounding these *FT* genes were visualised using the
484 GEvo tool of CoGE⁷².

485

486 The alignment of bHLH domains was performed with Clustal Omega⁷³, using sequences from Mertens *et*
487 *al*⁴⁹. The phylogeny was inferred using the maximum likelihood method based on the JTT matrix-based
488 model⁷⁴. Initial tree(s) for the heuristic search were obtained automatically by applying Neighbor-Join
489 and BioNJ algorithms to a matrix of pairwise distances estimated using a JTT model, and then selecting
490 the topology with superior log likelihood value. All positions containing gaps and missing data were
491 eliminated.

492

493 **Distinguishing and analysing the quinoa sub-genomes**

494 Trimmed PE Illumina sequencing reads that were used for the *de novo* assembly of *C. suecicum* and *C.*
495 *pallidicaule* were mapped onto the reference quinoa genome using the default settings of BWA. For
496 every base in the quinoa genome, the depth coverage of properly paired reads from the *C. suecicum* and
497 *C. pallidicaule* mapping was calculated using the program GenomeCoverage in the BEDtools package⁷⁵. A
498 custom Perl script was used to calculate the percentage of each scaffold with more than 5X coverage
499 from both diploids. Scaffolds were assigned to the A or B sub-genome if > 65% of the bases were
500 covered by reads from one diploid and < 25% of the bases were covered by reads from the other diploid.
501 The relationship between the quinoa sub-genomes and the diploid species *C. pallidicaule* and *C.*
502 *suecicum* was presented in a circle proportional to their sizes using Circos⁷⁶. Orthologous regions in the
503 three species were identified using BLASTN searches of the quinoa genome against each diploid genome
504 individually. Single top BLASTN hits longer than 8kb were selected and presented as links between the
505 quinoa genome assembly (arranged in chromosomes, see Supplementary Information 7.3.) and the two
506 diploid genome assemblies on the Circos plot (Fig. 2a).

507
508 Sub-genome synteny was analysed by plotting the positions of homoeologous pairs of A- and B-sub-
509 genome pairs within the context of the 18 chromosomes using Circos. Synteny between the sub-
510 genomes and *B. vulgaris* was assessed by first creating pseudomolecules by concatenating scaffolds
511 which were known to be ordered and oriented within each of the nine chromosomes. Syntenic regions
512 between these *B. vulgaris* chromosomes and those of quinoa were then identified using the
513 recommended settings of the CoGe SynMap tool⁷⁷ and visualized using MCScanX⁷⁸ and VGSC⁷⁹. For the
514 purposes of visualization, quinoa chromosomes CqB05, CqA08, CqB11, CqA15, and CqB16 were inverted.

516 **Saponin analyses**

517 Quinoa seeds were embedded in a 2% carboxymethylcellulose solution and frozen above liquid
518 nitrogen. Sections of 50 µm thickness were obtained using a Reichert-Jung Frigocut 2800N, modified to
519 use a Feather C35 blade holder and blades at -20°C using a modified Kawamoto method⁸⁰. A 2,5-
520 dihydroxybenzoic acid (Sigma-Aldrich) matrix (40 mg ml⁻¹ in 70% methanol) was applied using a HTX TM-
521 Sprayer (HTX Technologies LLC, Carrboro, NC, USA) with attached LC20-AD HPLC pump (Shimadzu
522 Scientific Instruments, Ermington, NSW, Australia). Sections were vacuum dried in a desiccator prior to
523 analysis. The optical image was generated using an Epson 4400 Flatbed Scanner at 4800 dpi. For mass
524 spectrometric analyses, a Bruker Solarix XR with 7T magnet was used. Images were generated using
525 Bruker Compass FlexImaging 4.1. Data were normalised to the TIC, and brightness optimisation was
526 employed to enhance visualisation of the distribution of selected compounds. Individual spectra were
527 recalibrated using Bruker Compass DataAnalysis 4.4 to internally lock masses of known DHB clusters:
528 C₁₄H₉O₆ = 273.039364 and C₂₁H₁₃O₉ = 409.055408 *m/z*. Accurate mass measurements for individual
529 saponins and identified compounds were run using Continuous Accumulation of Selected Ions (CASI)
530 using mass windows of 50-100 *m/z* and a transient of 4 Megaword generating a transient of 2.93 s
531 providing a mass resolving power of approximately 390,000 @ 400 *m/z*. Lipids were putatively assigned
532 by searching the LipidMaps database⁸¹ (www.lipidmaps.org) and lipid class confirmed by Collision
533 Induced Dissociation using a 10 *m/z* window centred around the monoisotopic peak with collision
534 energy of between 15-20 V.

535

536 Quinoa flowers were marked at anthesis, and seeds were sampled at 12, 16, 20, and 24 days after
537 anthesis. A pool of 5 seeds from each time point was analysed using GC/MS.

538
539 Quantification of saponins was performed indirectly by quantifying oleanolic acid (OA) derived from the
540 hydrolysis of saponins extracted from quinoa seeds. Derivatized solution was analysed using single
541 quadrupole GC-MS system (Agilent 7890 GC/5975C MSD) equipped with EI source at ionisation energy
542 of 70 eV. Chromatography separation was performed using DB-5MS fused silica capillary column (30m x
543 0.25 mm I.D., 0.25 µm film thickness; Agilent J&W Scientific), chemically bonded with 5% phenyl 95%
544 methylpolysiloxane cross-linked stationary phase. Helium was used as the carrier gas with constant flow
545 rate of 1.0 ml min⁻¹. The quantification of OA in each sample was performed using a standard curve
546 based on standards of OA.

547
548 Specific, individual saponins were identified in quinoa using a preparation of 20 mg of seeds performed
549 according a modified protocol from Giavalisco *et al.*⁸². Samples were measured with a Waters ACQUITY
550 Reversed Phase Ultra Performance Liquid Chromatography (RP-UPLC) coupled to a Thermo-Fisher
551 Exactive mass spectrometer which consists of an electrospray ionisation source and an Orbitrap mass
552 analyser. A C18 column was used for the hydrophilic measurements. Chromatograms were recorded in
553 Full Scan MS mode (Mass Range [100–1,500]). Extraction of the LC-MS data was accomplished with the
554 software REFINER MS 7.5 (GeneData).

555
556 SwissModel⁸³ was used to produce homology models for the bHLH region of AUR62017204,
557 AUR62017206 and AUR62010677. RaptorX⁸⁴ was used for prediction of secondary structure and
558 disorder. QUARK⁸⁵ was used for *ab initio* modelling of the C-terminal domain, and the DALI server⁸⁶ was
559 used for 3D homology searches of this region. Models were manually inspected and evaluated using the
560 Pymol program (pymol.org).

561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607

References

1. FAO, IFAD & WFP. *The State of Food Insecurity in the World 2015. Meeting the 2015 international hunger targets: taking stock of uneven progress*. Rome, FAO (2015).
2. Pretty, J. Agricultural sustainability: concepts, principles and evidence. *Phis. Trans. R. Soc. B.* **363**, 447–465 (2008).
3. Risi, C. & Galwey, N. W. The Chenopodium grains of the Andes: Inca crops for modern agriculture. *Adv. Appl. Biol.* **10**, 145–216 (1984).
4. Adolf, V. I., Shabala, S., Andersen, M. N., Razzaghi, F. & Jacobsen, S.-E. Varietal differences of quinoa's tolerance to saline conditions. *Plant Soil* **357**, 117–129 (2012).
5. Hariadi, Y., Marandon, K., Tian, Y., Jacobsen, S.-E. & Shabala, S. Ionic and osmotic relations in quinoa (*Chenopodium quinoa* Willd.) plants grown at various salinity levels. *J. Exp. Bot.* **62**, 185–193 (2011).
6. Jacobsen, S.-E., Mujica, A. & Jensen, C. R. The resistance of quinoa (*Chenopodium quinoa* Willd.) to adverse abiotic factors. *Food Rev. Int.* **19**, 99–109 (2003).
7. Gordillo-Bastidas, E., Díaz-Rizzolo, D. A., Roura, E., Massanés, T. & Gomis, R. Quinoa (*Chenopodium quinoa* Willd), from nutritional value to potential health benefits: an integrative review. *J. Nutr. Food Sci.* **6**, 497 (2016).
8. Vega-Gálvez, A. *et al.* Nutrition facts and functional potential of quinoa (*Chenopodium quinoa* Willd.), an ancient Andean grain: a review. *J. Sci. Food Agr.* **90**, 2541–2547 (2010).
9. Massawe, F., Mayes, S. & Cheng, A. Crop diversity: an unexploited treasure trove for food security. *Trends Plant Sci.* **21**, 365–368 (2016).
10. Zurita-Silva, A., Fuentes, F., Zamora, P., Jacobsen, S.-E. & Schwember, A. R. Breeding quinoa (*Chenopodium quinoa* Willd.): potential and perspectives. *Mol. Breeding* **34**, 13–30 (2014).
11. Lam, E. *et al.* Genome mapping on nanochannel arrays for structural variation analysis and sequence assembly. *Nat. Biotechnol.* **30**, 771–776 (2012).
12. Putnam, N. H. *et al.* Chromosome-scale shotgun assembly using an in vitro method for long-range linkage. *Genome Res.* **26**, 342–350 (2016).
13. Palomino, G., Hernández, L. T. & Torres, E. d. I. C. Nuclear genome size and chromosome analysis in *Chenopodium quinoa* and *C. berlandieri* subsp. *nuttalliae*. *Euphytica* **164**, 221–230 (2008).
14. Kolano, B., Siwinska, D., Pando, L. G., Szymanowska-Pulka, J. & Maluszynska, J. Genome size variation in *Chenopodium quinoa* (Chenopodiaceae). *Plant Syst. Evol.* **298**, 251–255 (2012).
15. Maughan, P. J. *et al.* Single nucleotide polymorphism identification, characterization, and linkage mapping in quinoa. *Plant Genome* **5**, 114–125 (2012).
16. Eilbeck, K., Moore, B., Holt, C. & Yandell, M. Quantitative measures for the management and comparison of annotated genomes. *BMC Bioinformatics* **10**, 67 (2009).
17. Yandell, M & Ence, D. A beginner's guide to eukaryotic genome annotation. *Nat. Rev. Genet.* **13**, 329–342 (2012).
18. Li, F. *et al.* Genome sequence of cultivated Upland cotton (*Gossypium hirsutum* TM-1) provides insights into genome evolution. *Nat. Biotechnol.* **33**, 524–530 (2015).
19. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* **31**, 3210–3212 (2015).
20. Gandía-Herrero, F & García-Carmona, F. Biosynthesis of betalains: yellow and violet plant pigments. *Trends Plant Sci.* **18**, 334–343 (2013).
21. Hatlestad, G. J. *et al.* The beet *R* locus encodes a new cytochrome P450 required for red betalain production. *Nat. Genet.* **44**, 816–820 (2012).

- 608 22. Kolano, B. *et al.* Molecular and cytogenetic evidence for an allotetraploid origin of *Chenopodium*
609 *quinoa* and *C. berlandieri* (Amaranthaceae). *Mol. Phylogenet. Evol.* **100**, 109–123 (2016).
- 610 23. Brown, D. C., Cepeda-Cornejo, V., Maughan, P. J. & Jellen, E. N. Characterization of the *Granule-*
611 *Bound Starch Synthase I* gene in *Chenopodium*. *Plant Genome* **8**, 1 (2014).
- 612 24. Štorchová, H., Drabešová, J., Cháb, D., Kolář, J. & Jellen, E. N. The introns in *FLOWERING LOCUS*
613 *T-LIKE (FTL)* genes are useful markers for tracking paternity in tetraploid *Chenopodium quinoa*
614 Willd. *Genet. Resour. Crop Evol.* **62**, 913–925 (2015).
- 615 25. Walsh, B. M., Adhikary, D., Maughan, P. J., Emshwiller, E. & Jellen, E. N. *Chenopodium* polyploidy
616 inferences from *Salt Overly Sensitive 1 (SOS1)* data. *Am. J. Bot.* **102**, 533–543 (2015).
- 617 26. Wendel, J. F., Brubaker, C., Alvarez, I., Cronn, R. & Stewart, J. M. Evolution and natural history of
618 the cotton genus. In *Genetics and Genomics of Cotton* (ed. Paterson, A. H.) 3–22 (Springer, New
619 York, 2009).
- 620 27. Koch, M. A., Haubold, B. & Mitchell-Olds, T. Comparative evolutionary analysis of chalcone
621 synthase and alcohol dehydrogenase loci in *Arabidopsis*, *Arabis*, and related genera
622 (Brassicaceae). *Mol. Biol. Evol.* **17**, 1483–1498 (2000).
- 623 28. Lynch, M. & Conery, J.S. The evolutionary fate and consequences of duplicate genes. *Science*
624 **10**, 1151–1155 (2000).
- 625 29. Wilson, H. D. Quinoa biosystematics II: free-living populations. *Econ. Bot.* **42**, 478–494 (1988).
- 626 30. Tapia, M.E., Mujica, A. & Canahua, A. Origen y distribución geográfica y sistemas de producción
627 de la quinua (*Chenopodium quinoa* Willd.). *Publicación Universidad Nacional Técnica del Altiplano*
628 (Perú) **1**, 1–5 (1980).
- 629 31. Wilson, H. D. & Heiser Jr., C. B. The origin and evolutionary relationships of ‘Huauzontle’
630 (*Chenopodium nuttalliae* Safford), domesticated chenopod of Mexico. *Amer. J. Bot.* **66**, 198–206
631 (1979).
- 632 32. Smith, B.D. *Chenopodium berlandieri* ssp. *jonesianum*: evidence for a Hopewellian domesticate
633 from Ash Cave, Ohio. *Southeastern Archaeology* **4**, 107–133 (1985).
- 634 33. Kolano, B *et al.* Chromosomal localization of two novel repetitive sequences isolated from the
635 *Chenopodium quinoa* Willd. genome. *Genome* **54**, 710–717 (2011).
- 636 34. Dohm, J. C. *et al.* The genome of the recently domesticated crop plant sugar beet (*Beta*
637 *vulgaris*). *Nature* **505**, 546–549 (2014).
- 638 35. Sankoff, D., Zheng, C. & Zhu, Q. The collapse of gene complement following whole genome
639 duplication. *BMC Genomics* **11**, 313 (2010).
- 640 36. Li, L., Stoeckert Jr., C. J. & Roos, D. S. OrthoMCL: identification of ortholog groups for eukaryotic
641 genomes. *Genome Res.* **13**, 2178–2189 (2003).
- 642 37. Pin, P. A. *et al.* An antagonistic pair of *FT* homologs mediates the control of flowering time in
643 sugar beet. *Science* **330**, 1397–1400 (2010).
- 644 38. Pin, P.A. & Nilsson, O. The multifaceted roles of *FLOWERING LOCUS T* in plant development.
645 *Plant Cell Environ.* **35**, 1742–1755 (2012).
- 646 39. Cháb, D., Kolář, J., Olson, M. S. & Štorchová, H. Two *FLOWERING LOCUS T (FT)* homologs in
647 *Chenopodium rubrum* differ in expression patterns. *Planta* **228**, 929–940 (2008).
- 648 40. Kuljanabagavad, T., Thongphasuk, P., Chamulitrat, W. & Wink, M. Triterpene saponins from
649 *Chenopodium quinoa* Willd. *Phytochemistry* **69**, 1919–1926 (2008).
- 650 41. de Geyter, E., Lambert, E., Geelen, D. & Smagghe, G. Novel advances with plant saponins as
651 natural insecticides to control pest insects. *Pest Technol.* **1**, 96–105 (2007).
- 652 42. Kuljanabagavad, T. & Wink, M. Biological activities and chemistry of saponins from
653 *Chenopodium quinoa* Willd. *Phytochem. Rev.* **8**, 473–490 (2009).
- 654 43. Konishi, Y., Hirano, S., Tsuboi, H. & Wada, M. Distribution of minerals in quinoa (*Chenopodium*
655 *quinoa* Willd.) seeds. *Biosci. Biotechnol. Biochem.* **68**, 231–234 (2004).

- 656 44. Boughton, B. A., Thinagaran, D., Sarabia, D., Bacic, A. & Roessner, U. Mass spectrometry imaging
657 for plant biology: a review. *Phytochem. Rev.* **15**, 445–488 (2015).
- 658 45. Augustin, J. M., Kuzina, V., Andersen, S. B. & Bak, S. Molecular activities, biosynthesis and
659 evolution of triterpenoid saponins. *Phytochemistry* **72**, 435–457 (2011).
- 660 46. Madl, T., Sterk, H., Mittelbach, M. & Rechberger, G. N. Tandem mass spectrometric analysis of a
661 complex triterpene saponin mixture of *Chenopodium quinoa*. *J. Am. Soc. Mass Spectrom.* **17**,
662 765–806 (2006).
- 663 47. Mastebroek, D. H., Limburg, H., Gilles, T. & Marvin, H. J. P. Occurrence of saponin in leaves
664 and seeds of quinoa (*Chenopodium quinoa* Willd.). *J. Sci. Food Agr.* **80**, 152–156 (2000).
- 665 48. Ward, S. M. A recessive allele inhibiting saponin synthesis in two lines of Bolivian quinoa
666 (*Chenopodium quinoa* Willd.). *J. Hered.* **92**, 83–86 (2001).
- 667 49. Mertens, J. *et al.* The bHLH transcription factors TSAR1 and TSAR2 regulate triterpene saponin
668 biosynthesis in *Medicago truncatula*. *Plant Physiol.* **170**, 194–210 (2016).
- 669 50. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence
670 data. *Bioinformatics* **30**, 2114–2120 (2014).
- 671 51. Zerbino, D. R. & Birney, E. Velvet: algorithms for de novo short read assembly using de Bruijn
672 graphs. *Genome Res.* **18**, 821–829 (2008).
- 673 52. Trapnell, C., Pachter, L. & Salzberg, S. L. TopHat: discovering splice junctions with RNA-Seq.
674 *Bioinformatics* **25**, 1105–1111 (2009).
- 675 53. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079
676 (2009).
- 677 54. Trapnell, C. *et al.* Transcript assembly and quantification by RNA-Seq reveals unannotated
678 transcripts and isoform switching during cell differentiation. *Nat. Biotechnol.* **28**, 511–515
679 (2010).
- 680 55. Gnerre, S. *et al.* High-quality draft assemblies of mammalian genomes from massively parallel
681 sequence data. *Proc. Natl. Acad. Sci. U.S.A.* **108**, 1513–1518 (2011).
- 682 56. Luo, R. *et al.* SOAPdenovo2: an empirically improved memory-efficient short-read *de novo*
683 assembler. *Gigascience* **1**, 18 (2012).
- 684 57. Smit, A. F. A. & Hubley, R. *RepeatModeler Open-1.0*. <<http://www.repeatmasker.org>> (2008-
685 2015).
- 686 58. Smit, A. F. A., Hubley, R. & Green, P. *RepeatMasker Open-4.0*. <<http://www.repeatmasker.org>>
687 (2013-2015).
- 688 59. Stanke, M., Diekhans, M., Baertsch, R. & Haussler, D. Using native and syntenically mapped
689 cDNA alignments to improve *de novo* gene finding. *Bioinformatics* **24**, 637–644 (2008).
- 690 60. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–
691 359 (2012).
- 692 61. Wu, T. D. & Watanabe, C. K. GMAP: a genomic mapping and alignment program for mRNA and
693 EST sequences. *Bioinformatics* **21**, 1859–1875 (2005).
- 694 62. Cantarel, B. L. *et al.* MAKER: an easy-to-use annotation pipeline designed for emerging model
695 organism genomes. *Genome Res.* **18**, 188–196 (2008).
- 696 63. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows–Wheeler transform.
697 *Bioinformatics* **26**, 589–595 (2010).
- 698 64. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009).
- 699 65. Heberle, H., Meirelles, G. V., da Silva, F. R., Telles, G. P. & Minghim, R. InteractiVenn: a web-
700 based tool for the analysis of sets through Venn diagrams. *BMC Bioinformatics* **16**, 169 (2015).
- 701 66. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:
702 improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).

- 703 67. Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European Molecular Biology Open Software
704 Suite. *Trends Genet.* **16**, 276–277 (2000).
- 705 68. Felsenstein, J. *PHYLP (Phylogeny Inference Package) Version 3.6a3*. Distributed by the author:
706 <http://evolution.genetics.washington.edu/phylip.html>
- 707 69. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).
- 708 70. Lee, T.-H., Guo, H., Wang, X., Kim, C. & Paterson, A. H. SNPhylo: a pipeline to construct a
709 phylogenetic tree from huge SNP data. *BMC Genomics* **15**, 162 (2014).
- 710 71. Tamura, K., Stecher, G., Peterson, D., Filipowski, A. & Kumar, S. MEGA6: Molecular Evolutionary
711 Genetics Analysis version 6.0. *Mol. Biol. Evol.* **30**, 2725–2729 (2013).
- 712 72. Lyons, E. & Freeling, M. How to usefully compare homologous plant genes and chromosomes as
713 DNA sequences. *Plant J.* **53**, 661–673 (2008).
- 714 73. Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple sequence alignments
715 using Clustal Omega. *Mol. Syst. Biol.* **7**, 539 (2011).
- 716 74. Jones, D. T., Taylor, W. R. & Thornton, J. M. The rapid generation of mutation data matrices
717 from protein sequences. *Method. Biochem. Anal.* **8**, 275–282 (1992).
- 718 75. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features.
719 *Bioinformatics* **26**, 841–842 (2010).
- 720 76. Krzywinski, M. *et al.* Circos: an information aesthetic for comparative genomics. *Genome Res.*
721 **19**, 1639–1645 (2009).
- 722 77. Lyons, E., Pedersen, B., Kane, J. & Freeling, M. The value of nonmodel genomes and an example
723 using SynMap within CoGe to dissect the hexaploidy that predates the rosids. *Trop. Plant Biol.* **1**,
724 181–190 (2008).
- 725 78. Wang, Y. *et al.* MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and
726 collinearity. *Nucl. Acids Res.* **40**, e49 (2012).
- 727 79. Xu, Y. *et al.* VGSC: a web-based vector graph toolkit of genome synteny and collinearity. *Biomed.*
728 *Res. Int.* **2016**, 7823429 (2016).
- 729 80. Kawamoto, T. Use of a new adhesive film for the preparation of multi-purpose fresh-frozen
730 sections from hard tissues, whole-animals, insects and plants. *Arch. Histol. Cytol* **66**, 123–143
731 (2003).
- 732 81. Sud, M. *et al.* LMSD: LIPID MAPS structure database. *Nucleic Acids Res.* **35**, D527–D532 (2007).
- 733 82. Giavalisco, P., Köhl, K., Hummel, J., Seiwert, B. & Willmitzer, L. ¹³C isotope-labeled metabolomes
734 allowing for improved compound annotation and relative quantification in liquid
735 chromatography-mass spectrometry-based metabolomic research. *Anal. Chem.* **81**, 6546–6551
736 (2009).
- 737 83. Arnold, K., Bordoli, L., Kopp, J. & Schwede, T. The SWISS-MODEL workspace: a web-based
738 environment for protein structure homology modelling. *Bioinformatics* **22**, 195–201 (2006).
- 739 84. Källberg, M., Margaryan, G., Wang, S., Ma, J. & Xu, J. RaptorX server: a resource for template-
740 based protein structure modeling. In *Protein Structure Prediction* (ed. Kihara, K.) 17–27
741 (Springer, New York, 2014).
- 742 85. Xu, D. & Zhang, Y. *Ab initio* protein structure assembly using continuous structure fragments and
743 optimized knowledge-based force field. *Proteins* **80**, 1715–1735 (2012).
- 744 86. Holm, L & Rosenström, P. Dali server: conservation mapping in 3D. *Nucleic Acids Res.* **38**, W545–
745 W549 (2010).
- 746
747
748

749 **Supplementary Information** is linked to the online version of the paper at www.nature.com/nature.
750

751 **Acknowledgements**

752 Research reported in this publication was supported by the King Abdullah University of Science and
753 Technology (KAUST), by USDA/NIFA-REEIS grant #2012-51300-20100 (WSU/BYU), and by a grant from
754 the German Research Foundation, DFG (grant No. JU205/24-1). We thank Thiruvarangan Ramaraj from
755 the National Center for Genomic Resources (Santa Fe, NM) for his bioinformatics assistance, Ohoud
756 Mohammed Eid Alharbi at the KAUST Imaging and Characterization Lab for generating SEM images,
757 Khadija Zemmouri at the KAUST Greenhouse for plant care, and Hyun Ho (Heno) Hwang and Ivan D.
758 Gromicho at the KAUST Academic Writing Department for assistance with quinoa illustrations and
759 photographs. We also thank Helena Štorchová, Frank Morton, Daniel Bertero, Francisco Fuentes and
760 David Brenner of USDA-ARS-NPGS for their assistance in providing seed.

761

762 **Author contributions**

763 MT and DEJ conceived the project. MT supervised the research. MT, DEJ, YSH, DJL, SMS, PJM, ENJ, ENVL,
764 CGvdL and TG conceived and designed the experiments and managed particular components of the
765 project. YSH led the bioinformatics analyses. YSH, DJL, DEJ and BL did most compilation of the genome
766 scaffolds and the genomic analyses. YSH annotated the genome and undertook the analysis of repetitive
767 elements. ENVL did the final genetic mapping. ENJ and PJM provided germplasm, oversaw sequencing of
768 the diploid genomes, oversaw the BioNano mapping, led the comparative genomics work and
769 generously shared their deep knowledge of quinoa. SMS oversaw all saponin-related analyses. CTM,
770 ARS, TB, ES, KM, ENJ and PJM prepared materials and undertook sequencing activities. KM, HO and SN
771 oversaw all phylogenetic analyses. YHW and GG analysed the microRNAs. NS, NMK, RRR, XG and SA-B
772 did saponin analyses. ND and CJ analysed the genes related to flowering time. STA and MAM did the
773 computational structure-function analysis. BAB and UR did the metabolomics imaging. All authors
774 contributed to the writing of the paper. DEJ, YSH, DJL, SMS, BL and MT organized the manuscript. DEJ
775 and MT coordinated the project. DEJ, YSH, DJL, SMS and BL contributed equally.

776

777 **Author information**

778 The genome assemblies and sequence data for *C. quinoa*, *C. pallidicaule* and *C. suecicum* were deposited
779 at NCBI under BioProject codes PRJNA306026, PRJNA326220 and PRJNA326219, respectively. Additional
780 accessions numbers for deposited data can be found in Supplementary File 9. The quinoa genome can
781 also be accessed at <http://www.cbrc.kaust.edu.sa/chenopodiumdb/> and on the Phytozome database
782 (www.phytozome.net/). Reprints and permission information is available at www.nature.com/reprints.
783 The authors declare no competing financial interests. Readers are welcome to comment on the online
784 version of the paper. Correspondence and requests for materials should be addressed to MT
785 (mark.tester@kaust.edu.sa).

786

787

788
789
790

Tables

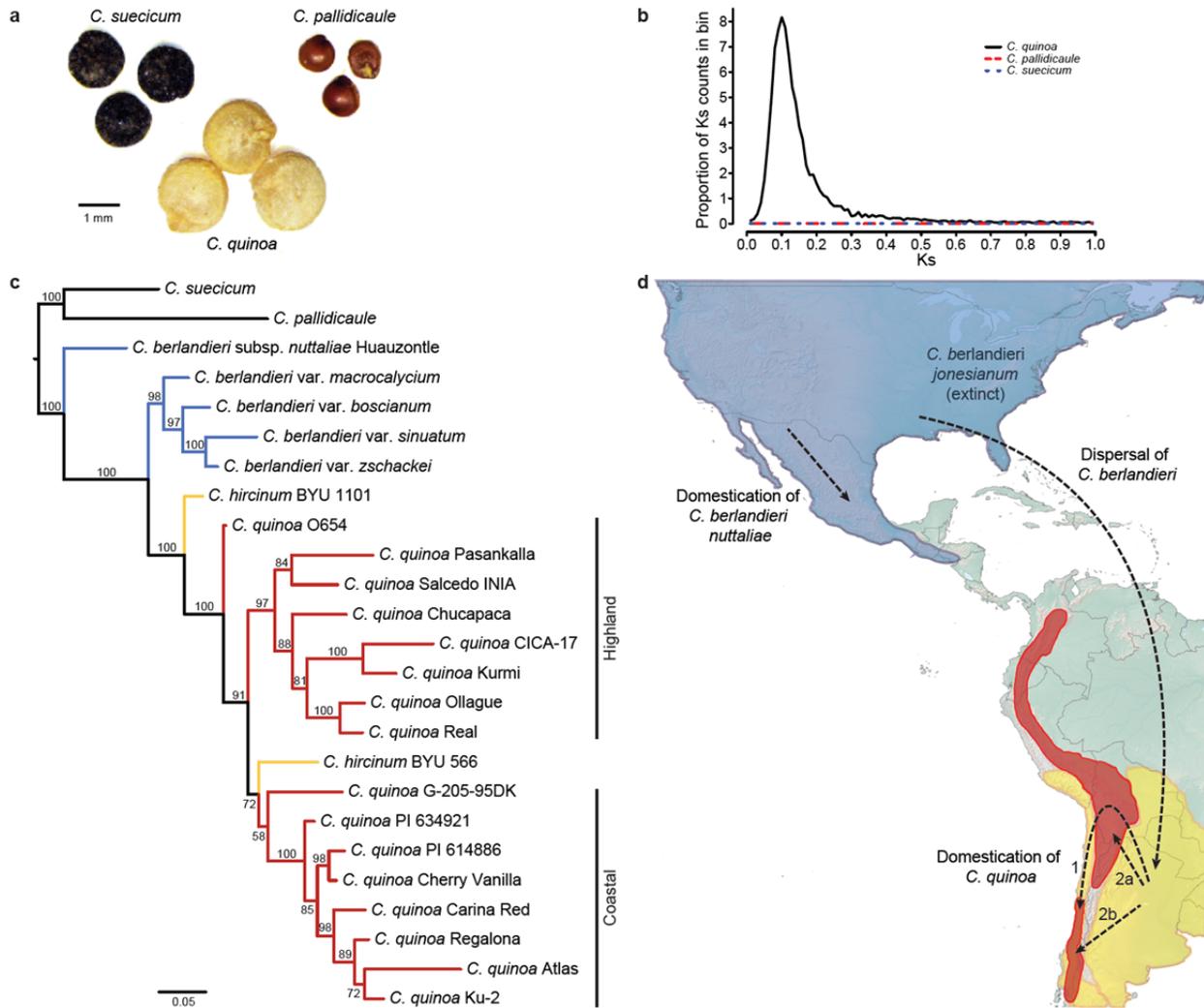
Table 1 | Assembly statistics for quinoa, *C. pallidicaule*, and *C. suecicum*

	<i>C. pallidicaule</i> Diploid (2n=2x=18)	<i>C. suecicum</i> Diploid (2n=2x=18)	<i>C. quinoa</i> Allotetraploid (2n=4x=36)			
Illumina	✓	✓				
PacBio			✓	✓	✓	✓
BioNano				✓	✓	✓
Dovetail					✓	✓
Linkage map						✓
Total assembly size (bp)	337,010,935	536,949,265	1,325,007,020	1,395,179,653	1,385,456,844	1,385,456,844
Longest scaffold (bp)	2,949,784	1,614,553	11,561,360	11,561,360	23,816,425	133,230,289*
Number of contigs	-	-	4,232	-	-	-
N50 contig length (bp)	-	-	1,663,340	-	-	-
L50 contig count	-	-	216	-	-	-
Number of scaffolds	3,013	11,198	-	4,014	3,486	18*
N50 scaffold length (bp)	356,818	105,389	-	2,450,933	3,846,917	64,593,919*
L50 scaffold count	243	1,285	-	177	105	10*
N90 scaffold length (bp)	55,204	27,807	-	157,165	249,904	26,142,578*
L90 scaffold count	1,215	5,075	-	800	439	17*
Missing bases (%)	2.52	7.49	0.00	4.53	4.56	4.56

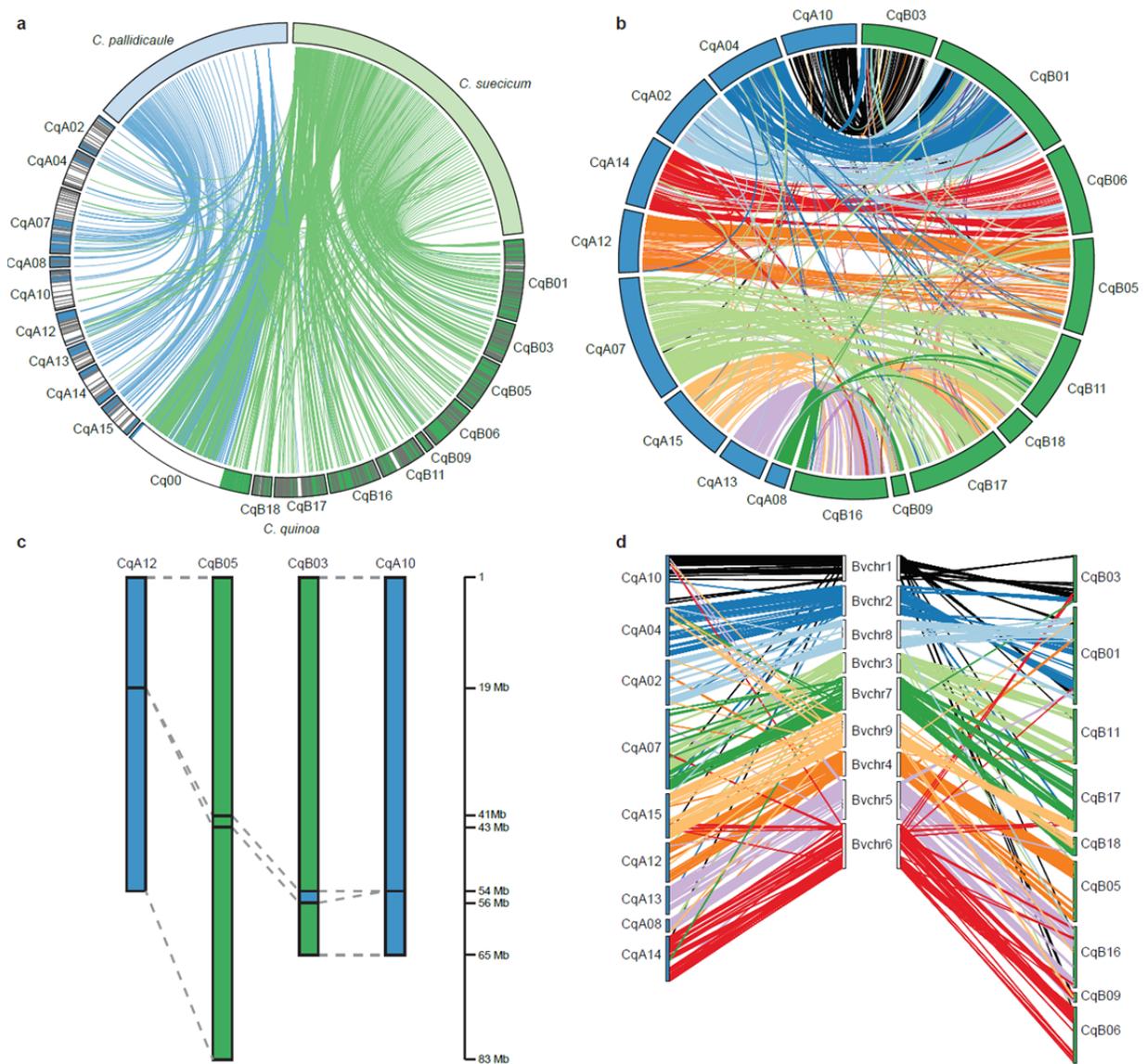
791

*Based on scaffolds assigned to LGs

792 **Figure legends**
 793

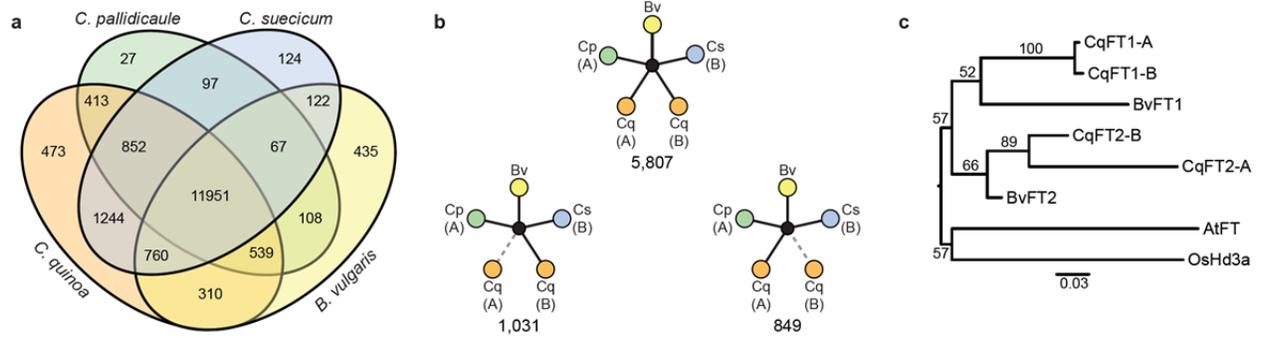


794 **Figure 1 | Evolutionary history of quinoa.** **a**, Seeds of *C. suecicum*, *C. pallidicaule*, and quinoa. **b**, The
 795 proportion of gene pairs in each species binned according to Ks values. **c**, Maximum likelihood tree
 796 generated from 3,132 SNPs. Black branches, diploid species. Coloured branches, tetraploid species: red,
 797 quinoa accessions; blue, *C. berlandieri* accessions; yellow, *C. hircinum* accessions. Branch values
 798 represent the percentage of 1,000 bootstrap replicates that support the topology. Scale bar represents
 799 substitutions per site. **d**, Evolutionary relationships of *Chenopodium* species, showing the hypothesised
 800 long-range dispersal of an ancestral *C. berlandieri* to South America, and the eventual domestication of
 801 quinoa from *C. hircinum*, either from a single event (1) that gave rise to highland and subsequently
 802 coastal quinoa, or in two events that gave rise to highland (2a) and coastal quinoa (2b) independently.
 803 Blue, red, and yellow shading represents the geographic distribution of *C. berlandieri*, quinoa, and
 804 *C. hircinum*, respectively.
 805
 806
 807



808
809

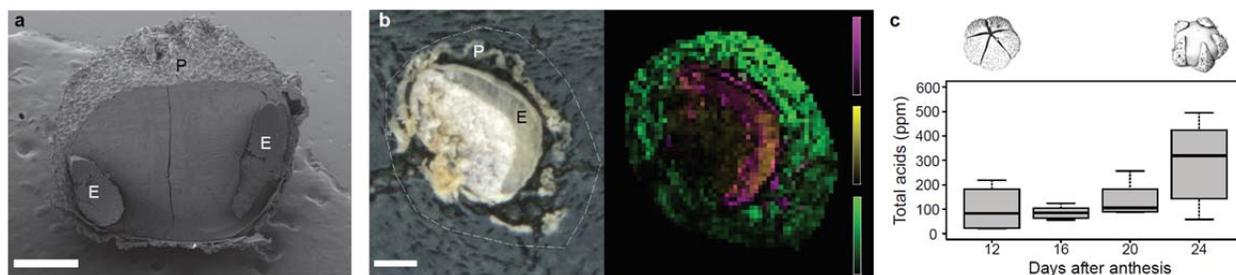
810 **Figure 2 | Identification and characterisation of quinoa sub-genomes.** **a**, Blue lines and green lines
 811 connect regions of the *C. pallidicaule* and *C. suecicum* genomes, respectively, with their orthologous
 812 regions in the quinoa genome based on BLASTN. Quinoa scaffolds are arranged into chromosomes, with
 813 blue- and green-coloured bars indicating sub-genome assignment based on mapped reads from the two
 814 diploid species. Scaffolds that could not be unambiguously assigned to a sub-genome based on read
 815 mapping are shown in white. Grey bars separate neighbouring scaffolds. **b**, Homoeologous gene pairs in
 816 the A (blue chromosomes) and B (green chromosomes) sub-genomes. **c**, Simplified representation of
 817 synteny between CqA12, CqB05, CqB03, and CqA10, highlighting the homeologous regions between the
 818 B sub-genome chromosomes. Dotted lines connect large-scale syntenic regions between the A (blue)
 819 and B (green) sub-genomes. The scale bar indicates approximate positions defining the indicated
 820 syntenic blocks. For purposes of visualization, CqB05 and Cq10A were inverted. **d**, Syntenic relationships
 821 between *B. vulgaris* (Bv) and the A and B sub-genomes of quinoa. Colours distinguish quinoa regions
 822 syntenic to each *B. vulgaris* chromosome. Blue and green quinoa chromosomes indicate the A and B
 823 sub-genomes, respectively.



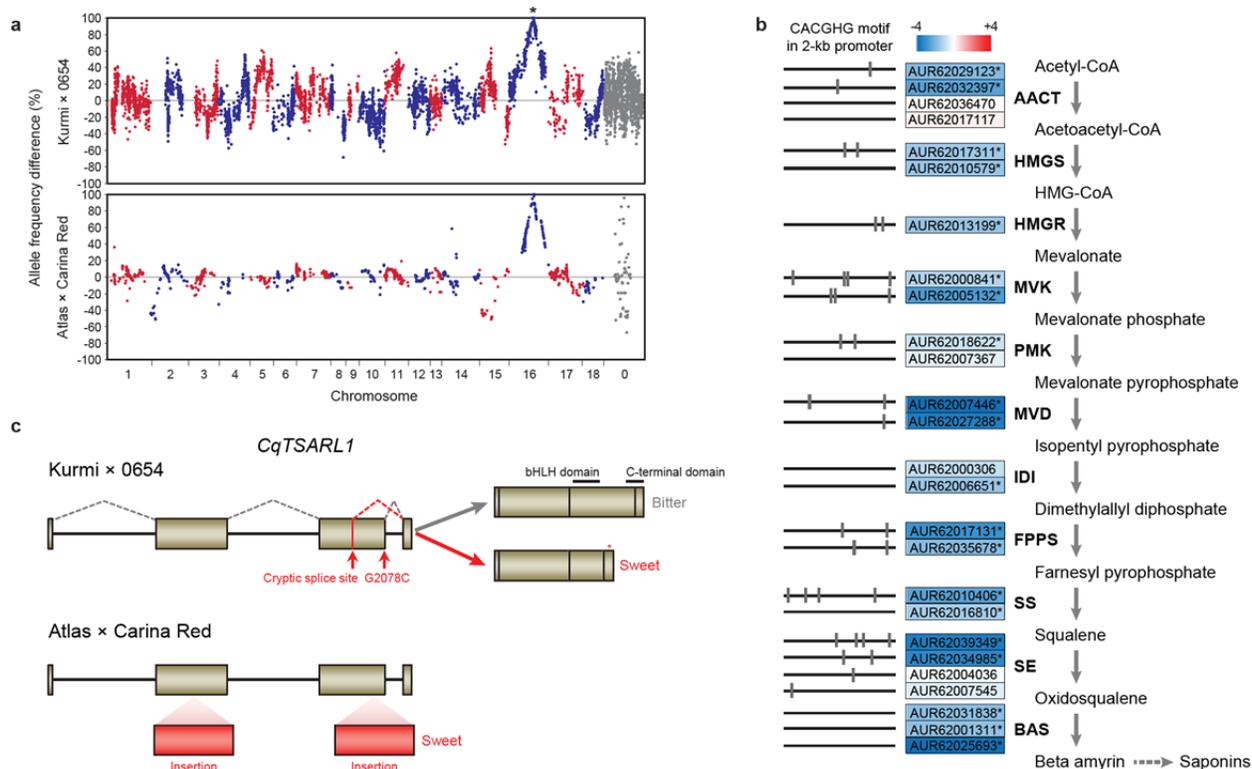
824
825

826 **Figure 3 | Sub-genome gene loss and retention.** **a**, The number of orthologous protein-coding gene
827 clusters shared between or unique to quinoa, *C. pallidicaule*, *C. suecicum*, and *B. vulgaris*. **b**, The number
828 of gene sets for which each gene has been retained as a single copy in each genome/sub-genome
829 (middle), or lost from the quinoa A (left) or B (right) sub-genome. **c**, Maximum likelihood tree of
830 FLOWERING LOCUS T (FT) sequences, indicating the presence of two sets of orthologs in quinoa (Cq) and
831 *B. vulgaris* (Bv). The tree is rooted on the branch containing orthologous FT sequences from *Arabidopsis*
832 *thaliana* (AtFT) and *Oryza sativa* (OsHd3a). Branch values represent the percentage of 1,000 bootstrap
833 replicates that support the topology. Scale bar represents substitutions per site.

834
835



836 **Figure 4 | Saponins in the seeds of the reference quinoa accession PI 614886.** **a**, SEM image of a quinoa
837 seed cross-section, showing the pericarp (P) and embryo (E). **b**, Imaging MS visualization of selected
838 masses, including saponins in the pericarp of a quinoa seed. Purple gradient bar, tentative
839 phosphatidylcholine-(34:1), $[M+Na]^+$ $m/z = 782.5610$, calc. 782.5670, 7.7 ppm error); yellow gradient
840 bar, tentative triacylglycerol-(54:6), $[M+K]^+$ $m/z = 917.6971$, calc. 917.6995, 2.6 ppm error); green
841 gradient bar indicates a representative saponin phytolaccagenic acid with sugar chains hexose-pentose-
842 hexose $[M+K]^+$ $m/z = 1173.5114$, calc. 1173.509, -2.0 ppm error). Coloured bars represent the ion signal
843 intensity scaled from 0% (bottom) to 50% (top) of maximum signal. Scale bars in **a** and **b**, 500 μm . **c**,
844 Accumulation of saponins as measured by total acids during seed development. Illustrations represent
845 fruit development at 12 and 24 days after anthesis.



846
847 **Figure 5 | Candidate gene underlying saponin production.** **a**, Mapping of the saponin production locus.
848 The percentage difference in allele frequency of sweet progeny compared to bitter progeny in the Kurmi
849 × 0654 (top) and Atlas × Carina Red (bottom) populations. Alternating red and blue dots indicate
850 positions of markers along alternating chromosomes, with unmapped markers in chromosome 0 shown
851 in grey. Asterisk above the top panel indicates the approximate position of *CqTSARL1*. **b**, The saponin
852 biosynthetic pathway, showing enzymes that catalyse each step of the pathway and the quinoa gene ID
853 for genes encoding each enzyme. Boxes surrounding each gene ID are coloured according to their fold
854 change in expression (\log_2) in sweet lines compared to bitter lines of Kurmi × 0654. Horizontal lines to
855 the left of each gene ID represent the 2-kb region upstream of the start codon of each gene, with tick
856 marks indicating the positions of motifs putatively recognized by *CqTSARL1*. **c**, Gene models of *CqTSARL1*
857 in bitter and sweet lines. In bitter lines of both populations, normal splicing of the four exons produces a
858 full-length transcript that contains the bHLH and conserved C-terminal domains (indicated). In sweet
859 lines of Kurmi × 0654, alternative splicing occurs at a cryptic splice site in the third exon, likely because
860 of the G20178C mutation at the last position of the exon. This results in an alternative transcript that
861 lacks the C-terminal domain and contains a premature stop codon (red asterisk). In Atlas × Carina Red,
862 sweet lines contain insertions in the second and third exons, which consequently likely disrupt the
863 protein function.