

Genomic and phenotypic evidence for an incomplete domestication of South American grain amaranth (*Amaranthus caudatus*)

Markus G. Stetter, Thomas Müller and Karl J. Schmid

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

Keywords

Amaranthus, genotyping-by-sequencing, genetic diversity, domestication, flow cytometry, orphan crop

Corresponding Author

Name: Karl Schmid, Prof. Dr.

Address: Institute of Plant Breeding, Seed Science and Population Genetics (350)

University of Hohenheim

Fruwirthstraße 21

D-70599 Stuttgart / Germany

Tel: +49 711 459 23487

Fax: +49 711 459 24458

E-Mail address: karl.schmid@uni-hohenheim.de

Running title: Incomplete domestication of *A. caudatus*

Incomplete Amaranth domestication

1 **Abstract**

2 The process of domestication leads to major morphological and genetic changes, which in
3 combination are known as domestication syndrome that differentiates crops from their wild an-
4 cestors. We characterized the genomic and phenotypic diversity of the South American grain
5 amaranth *Amaranthus caudatus*, which has been cultivated for thousands of years and is one
6 of the three native grain amaranths of South and Central America. Previously, several models
7 of domestication were proposed including a domestication from the close relatives and putative
8 ancestors *A. hybridus* or *A. quitensis*. To investigate the evolutionary relationship of *A. caudatus*
9 and its two close relatives, we genotyped 119 amaranth accessions of the three species from
10 the Andean region using genotyping-by-sequencing (GBS) and compared phenotypic variation
11 in two domestication-related traits, seed size and seed color. The analysis of 9,485 SNPs re-
12 vealed a strong genetic differentiation of cultivated *A. caudatus* from the relatives *A. hybridus*
13 and *A. quitensis*. The two relatives did not cluster according to the species assignment but
14 formed mixed groups according to their geographic origin in Ecuador and Peru, respectively.
15 *A. caudatus* had a higher genetic diversity than its close relatives and shared a high proportion
16 of polymorphisms with their wild relatives consistent with the absence of a strong bottleneck
17 or a high level of recent gene flow. Genome sizes and seed sizes were not significantly differ-
18 ent between *A. caudatus* and its relatives, although a genetically distinct group of *A. caudatus*
19 from Bolivia had significantly larger seeds. We conclude that despite a long history of human
20 cultivation and selection for white grain color, *A. caudatus* shows a weak genomic and pheno-
21 typic domestication syndrome and is an incompletely domesticated species because of weak
22 selection or high levels of gene flow from its sympatric close undomesticated relatives that
23 counteracted the fixation of key domestication traits.

Incomplete Amaranth domestication

24 **Introduction**

25 Research on the domestication of crop plants revealed that numerous traits can be affected by
26 domestication, which results in so-called domestication syndromes. The type and extent of do-
27 mestication syndromes depends on the life history and uses of crop plants (Meyer *et al*, 2012),
28 although crops from distantly related plant families frequently show similar domestication phe-
29 notypes. For example, the 'classical' domestication syndrome, which includes larger seeds,
30 loss of seed shattering, reduced branching, loss of seed dormancy and decreased photoperiod
31 sensitivity, is observed in legumes and grasses (Abbo *et al*, 2014; Hake & Ross-Ibarra, 2015).
32 Similar to phenotypic diversity, crops show variable genomic signatures of domestication be-
33 cause of differences in their biology and utilization by humans (Meyer *et al*, 2012). In particular,
34 domestication affects the level and structure of genetic diversity in crops because selection and
35 genetic drift contributed to strong genetic bottlenecks (Doebley *et al*, 2006; Olsen & Wendel,
36 2013; Sang & Li, 2013; Nabholz *et al*, 2014). The geographic expansion of domesticated crops
37 provided the opportunity for gene flow with new crop wild relatives, which further contributed to
38 genetic differentiation from wild ancestors. Such a diversity of phenotypic and genomic changes
39 associated with domestication suggest that the classical model of a single domestication event
40 in a short time span within a small geographic region may not apply to numerous crop plants
41 like barley, apple and olive trees (Besnard & Rubio de Casas, 2015; Cornille *et al*, 2012; Poets
42 *et al*, 2015). The motivation of the present study was to investigate both the phenotypic and
43 genomic consequences of amaranth cultivation in the light of these concepts.

44 The genus *Amaranthus* L. comprises between 50 and 75 species and is distributed worldwide
45 (Sauer, 1967; Costea & DeMason, 2001). Four species are cultivated as grain amaranths or
46 leaf vegetables (Sauer, 1967; Brenner *et al*, 2010). The grain amaranths *Amaranthus caudatus*,
47 *Amaranthus cruentus* and *Amaranthus hypochondriacus* originated from South and Central
48 America while *A. tricolor* is used as leafy vegetable in Africa. Amaranth is an ancient crop
49 because archaeological evidence in Northern Argentina suggested that wild amaranth seeds
50 were collected and used for human consumption during the initial mid-Holocene (8,000 - 7,000
51 BP; Arreguez *et al*, 2013). In the Aztec empire, amaranth was a highly valued crop and tributes
52 were collected from the farmers that were nearly as high as for maize (Sauer, 1967). Currently,

Incomplete Amaranth domestication

53 amaranth is promoted as a healthy food because of its favorable composition of essential amino
54 acids and high micronutrient content.

55 The three grain amaranth species differ in their geographical distribution. *A. cruentus* and
56 *A. hypochondriacus* are most common in Central America, whereas *A. caudatus* is cultivated
57 mainly in South America. In the Andean region, *A. caudatus* grows in close proximity to the
58 two related (i.e., wild) species *A. hybridus* and *A. quitensis*, which are considered as potential
59 ancestors (Sauer, 1967). *A. hybridus* has the widest distribution range from Central to South
60 America while *A. quitensis* is restricted to the central part of South America. *A. quitensis* was
61 tolerated and possibly cultivated in Andean home gardens and used for coloring in historical
62 times.

63 The history of amaranth cultivation and extent of its domestication are still under discussion
64 (Figure 1). Sauer (1967) proposed two domestication scenarios based on the morphology and
65 geographic distribution of the different species. One scenario postulates three independent
66 domestication events from three different wild ancestors, and another scenario proposes the
67 domestication of *A. cruentus* from *A. hybridus* followed by a migration and intercrossing of *A.*
68 *cruentus* with *A. powellii* in Central America and an intercrossing of *A. cruentus* with *A. quiten-*
69 *sis* resulting in *A. caudatus* in South America. A third scenario was based on genetic markers
70 and suggested that all three cultivated amaranths evolved from *Amaranthus hybridus*, but at
71 multiple locations (Maughan *et al*, 2011). Most recently, Kietlinski *et al* (2014) proposed a
72 single domestication of *A. hybridus* in the Andes or in Mesoamerica and a subsequent spatial
73 separation of two lineages leading to *A. caudatus* and *A. hypochondriacus*, or two indepen-
74 dent domestication events of *A. hypochondriacus* and *A. caudatus* from a single *A. hybridus*
75 lineage in Central and South America (Figure 1C and D). A more recent analysis based on the
76 phylogeny of the whole *Amaranthus* genus supports independent domestication of the South
77 American *A. caudatus* and the two Central American grain amaranths from two different, geo-
78 graphically separated lineages of *A. hybridus* as shown in Figure 1D (Stetter & Schmid, 2016).
79 Despite its long history of cultivation and the self-pollinating breeding system, the domestication
80 syndrome of cultivated amaranth is remarkably indistinct because it still shows strong photope-
81 riod sensitivity and has very small shattering seeds (Sauer, 1967; Brenner *et al*, 2010). Other
82 crops like maize that were cultivated at a similar time period in the same region exhibit the

Incomplete Amaranth domestication

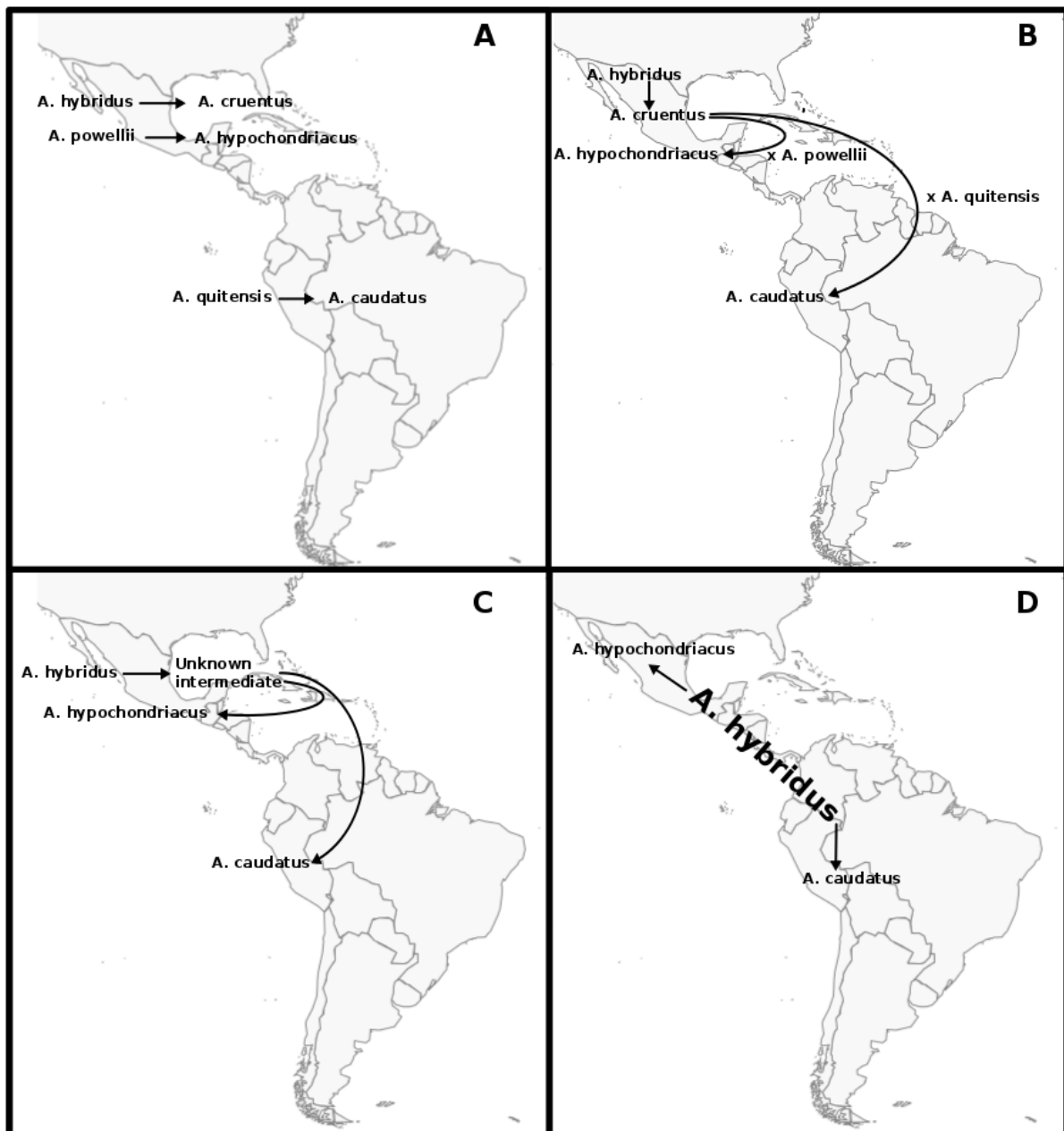


Figure 1: Models of amaranth domestication. (A) Independent domestication of three grain amaranths from different close relatives (Sauer, 1967). (B) Initial domestication from *A. hybridus* and subsequent migration and hybridization with additional close relatives (Sauer, 1967). (C) Single domestication in the Andes or in Mesoamerica and subsequent spatial separation of two lineages leading to *A. caudatus* and *A. hypochondriacus* (Kietlinski *et al*, 2014). (D) Two domestication events from a single *A. hybridus* lineage spanning Central and South America (Kietlinski *et al*, 2014).

83 classical domestication syndrome (Sang & Li, 2013; Lenser & Theißen, 2013). This raises the
84 question whether amaranth is domesticated at all or has a different domestication syndrome,
85 and if the latter is true whether genetic constraints, a lack of genetic variation or (agri-)cultural

Incomplete Amaranth domestication

86 reasons determined its domestication syndrome. The phenotypic analysis of amaranth do-
87 mestication is complicated by the taxonomic uncertainty of cultivated amaranth species and
88 their close relatives. Although *A. quitensis* was suggested to be the ancestor of *A. caudatus*,
89 the state of *A. quitensis* as a separate species is under debate. Sauer (1967) classified it as
90 species, but later it was argued that it is the same species as *A. hybridus* (Coons, 1978; Bren-
91 ner *et al*, 2010). However, until today *A. quitensis* is treated as separate species and since
92 genetic evidence for the status of *A. quitensis* as a separate species is based on few studies
93 with limited numbers of markers, this topic is still unresolved (Mallory *et al*, 2008; Kietlinski
94 *et al*, 2014).

95 The rapid development of sequencing technologies facilitates the large-scale investigation of
96 the genetic history of crops and their close relatives. Among available methods, reduced repre-
97 sentation sequencing approaches such as genotyping-by-sequencing (GBS) allow a genome-
98 wide and cost-efficient marker detection compared to whole genome sequencing (Elshire *et al*,
99 2011; Poland *et al*, 2012). Despite some biases associated with reduced representation se-
100 quencing, GBS and related methods are suitable and powerful approaches for studying inter-
101 specific phylogenetic relationships (Cruaud *et al*, 2014) and intraspecific patterns of genetic
102 variation in crop plants (Morris *et al*, 2013).

103 We used GBS and genome size measurements to characterize the genetic diversity and rela-
104 tionship of cultivated *A. caudatus* and its possible ancestors *A. quitensis* and *A. hybridus*, and
105 compared patterns of genetic structure with two domestication-related phenotypic traits (seed
106 color and hundred seed weight). For this study, we focussed on the South American amaranth
107 species, because *A. caudatus*, *A. quitensis* and South American accessions of *A. hybridus* form
108 a clade that is strongly separated from the two Central American grain amaranths in a phylo-
109 genetic analysis of the whole genus (Stetter & Schmid, 2016). For this reason, we reasoned
110 that the domestication of *A. caudatus*, which is native to South America, and its relationship
111 to the sympatric relatives, *A. hybridus* and *A. quitensis* can be conducted independently of the
112 Central American amaranth species. Our analysis includes a comparison of genetic diversity
113 and seed-related traits like size and color between cultivated and wild amaranths and analyses
114 the taxonomic relationship and gene flow among species. Our results indicate that *A. caudatus*
115 has a history of domestication that may be considered as incomplete.

116 **Material and Methods**

117 **Plant material**

118 A total of 119 amaranth accessions of three *Amaranthus* species originating from South Amer-
119 ica were obtained from the USDA genebank ([http://www.ars-grin.gov/npgs/searchgrin.](http://www.ars-grin.gov/npgs/searchgrin.html)
120 [html](http://www.ars-grin.gov/npgs/searchgrin.html)). Of these accessions, 89 were classified as *A. caudatus*, 17 as *A. hybridus*, seven as
121 *A. quitensis* and six as interspecific hybrids according to the passport information (Table S5).
122 We selected *A. caudatus* accessions based on the altitude of the collection site and focused
123 on high-altitude regions (2,200 to 3,700 m) where amaranth has been cultivated for thousands
124 of years and survived until today since it fell into disuse after the Spanish conquest (Kauffman
125 & Weber, 1990). Therefore, high-altitude accessions may represent a large proportion of the
126 species-wide genetic diversity. The smaller sample sizes of *A. hybridus* and *A. quitensis* ac-
127 cessions reflect that fewer accessions of these species than of *A. caudatus* are available from
128 the USDA genebank. However, the geographic origin of the two wild relatives covers the An-
129 dean highlands, which is the distribution range of *A. caudatus*, and we compared the population
130 structure of the sample derived from the genomic data with the passport information to test for
131 consistency between the population structure and geographic origin. Accessions were planted
132 in a field in Nürtingen (Germany), and one young leaf of one representative plant per accession
133 was sampled to avoid the sampling of potential seed cross-contamination. We sampled and
134 sequenced three plants each of 12 accessions independently for quality control.

135 **Genome size**

136 To compare genome sizes among the three diploid *Amaranthus* species, we measured the
137 genome size of 22 *A. caudatus*, 8 *A. hybridus* and 4 *A. quitensis* accessions. Genome size
138 differences of individuals within species are expected to be low, and we therefore estimated
139 species-specific genome sizes using 25% the total sample of *A. caudatus* and 50% of *A. hy-*
140 *bridus* and *A. quitensis* accessions, respectively. Plants were grown for four weeks in the
141 greenhouse before one young leaf was collected for cell extraction. A tomato cultivar (*Solanum*
142 *lycopersicum* cv Stupicke) was used as internal standard because it has a comparable genome

Incomplete Amaranth domestication

143 size that has been measured with high accuracy (DNA content = 1.96 pg; Dolezel *et al*, 1992).
144 Fresh leaves were cut up with a razor blade and cells were extracted with CyStain PI Abso-
145 lute P (Partec, Muenster/Germany). Approximately 0.5 cm² of the sample leaf was extracted
146 together with similar area of tomato leaf in 0.5 ml of extraction buffer. The DNA content was
147 determined with CyFlow Space (Partec, Muenster/Germany) flow cytometer and analyzed with
148 FlowMax software (Partec, Muenster/Germany). For each sample, 10,000 particles were mea-
149 sured each time. Two different plants were measured for each accession. The DNA content
150 was calculated as:

$$151 \quad \text{DNA content 2C [pg]} = \text{genome size tomato} \times \frac{\text{fluorescence amaranth}}{\text{fluorescence tomato}}$$

152 and the genome size (in Mbp) was calculated as:

$$153 \quad \text{genome size 1C [Mbp]} = (0.978 * 10^3) \times \frac{\text{DNA content 2C [pg]}}{2}$$

154 The conversion from pg to bp was calculated with 1pg DNA = 0.978 × 10⁹ bp (Dolezel *et al*,
155 2003). Means were calculated using R software (Team, 2014) and an ANOVA was performed
156 to infer differences in genome size for the species.

157 **DNA extraction and library preparation**

158 Genomic DNA was extracted using a modified CTAB protocol (Saghai-Marooft *et al*, 1984). The
159 DNA was dried and dissolved in 50-100 μl TE and diluted to 100 ng/μl for further usage. Two-
160 enzyme GBS libraries were constructed with a modified protocol from the previously described
161 two-enzyme GBS protocol (Poland *et al*, 2012). DNA was digested with a mix of 2 μl DNA,
162 2 μl NEB Buffer 2 (NEB, Frankfurt/Germany), 1 μl ApeKI (4U/μl, NEB), 1 μl HindIII (20U/μl,
163 NEB) and 14 μl ddH₂O for 2 hours at 37°C before incubating for 2 hours at 75°C. Adapters
164 were ligated with 20 μl of digested DNA 5 μl ligase buffer (NEB), T₄- DNA ligase (NEB), 4 μl
165 ddH₂O and 20 μl of adapter mix containing 10μl barcode adapter (0.3 ng/μl) and 10 μl common
166 adapter (0.3ng/μl). Samples were incubated at 22°C for 60 minutes before deactivating ligase
167 at 65°C for 30 minutes. Subsequently, samples were cooled down to 4°C. For each sequencing
168 lane, 5μl of 48 samples with different barcodes were pooled after adapter ligation. Samples of
169 the different species were randomized over the 3 pools and different barcode lengths. The 12
170 replicated samples were added to each pool. The pooled samples were purified with QIAquick

Incomplete *Amaranth* domestication

171 PCR purification kit (Qiagen, Hilden/Germany) and eluted in 50 μ l elution buffer before PCR
172 amplification of the pools. The PCR was performed with 10 μ l of pooled DNA, 25 μ l 2x Taq
173 Master Mix (NEB), 2 μ l PCR primer mix (25pmol/ μ l of each primer) and 13 μ l ddH₂O for 5 min
174 at 72°C and 30 sec at 98°C before 18 cycles of 10 sec at 98°C, 30 sec at 65°C and 30 sec at
175 72°C after the 18 cycles 5 min of 72°C were applied and samples were cooled down to 4°C.
176 Samples were purified again with QIAquick PCR purification kit (Qiagen) and eluted in 30 μ l
177 elution buffer. Three lanes with 48 samples per lane were sequenced on an Illumina HighScan
178 SQ with single end and 105 cycles on the same flow cell (see supporting data).

179 **Data preparation**

180 Raw sequence data were filtered with the following steps. First, reads were divided into sepa-
181 rate files according to the different barcodes using Python scripts. Read quality was assessed
182 with fastQC (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). Due to lower
183 read quality towards the end of the reads, they were trimmed to 90 bp. Low quality reads were
184 excluded if they contained at least one N (undefined base) or if the quality score after trimming
185 was below 20 in more than 10% of the bases. Data from technical replicates were combined
186 and individuals with less than 10,000 reads were excluded from further analysis (Table S5).
187 The 12 replicated samples were used to detect a lane effect with an analysis of variance.

188 **SNP calling and filtering**

189 Since no high quality reference genome for *Amaranthus* sp. was available for read mapping,
190 we used *Stacks* 1.19, for the *de novo* identification of SNPs in GBS data (Catchen *et al*,
191 2011, 2013). The SNP calling pipeline provided by *Stacks* `denovo_map.pl` was used to call
192 SNPs from the processed data. Highly repetitive GBS reads were removed in the `ustacks`
193 program with option `-t`. Additionally, the minimum number of identical raw reads required to
194 create a stack was set to three ($m=3$) and the number of mismatches allowed between loci
195 when processing a single individual was two ($M=2$). Four mismatches were allowed between
196 loci when building the catalog ($n=4$). The catalog is a set of non redundant loci representing all
197 loci in the accessions and used as reference for SNP calling. SNPs were called with the *Stacks*

Incomplete Amaranth domestication

198 tool `populations 1.19` without filtering for missing data using option `-r 0`. One individual, PI
199 511754, was excluded from further analysis because it appeared to be misclassified. According
200 to its passport information it belonged to *A. hybridus*, but with all clustering methods it was
201 placed into a separate cluster consisting only of this individual, which suggested it belongs to
202 a different species. Therefore, we repeated the SNP calling without this individual. The SNPs
203 were filtered over the whole sample for missing data with `vcftools` (Danecek *et al*, 2011) by
204 allowing a maximum of 60% missing values per SNP position. Given the stringent filtering
205 criteria for SNP calling and the restricted number of *A. quitensis* individuals, we did not filter
206 SNPs by their minor allele frequency for further analysis.

207 Inference of genetic diversity and population structure

208 Nucleotide diversity (π) weighted by coverage was calculated with a Python script that imple-
209 ments the formula of Begun *et al* (2007) which corrects for different sampling depths of SNPs in
210 sequencing data. The confidence interval of π was calculated by bootstrapping the calculation
211 10,000 times. To account for the difference in sampling between wild and cultivated amaranths,
212 we sub-sampled *A. caudatus* 100 times with the the same number of individuals (23) as used
213 for wild amaranth. The pairwise difference in π between *A. caudatus* and the close relatives was
214 calculated for each site. Mean expected (H_{exp}) and observed (H_{obs}) heterozygosities based on
215 SNPs were calculated with the R package `adegenet 1.4-2` (Jombart & Ahmed, 2011). The
216 inbreeding coefficient (F) was calculated as:

$$217 \quad \frac{H_{exp} - H_{obs}}{H_{exp}}$$

218 Weir and Cockerham weighted F_{ST} estimates were calculated with `vcftools` (Weir & Cocker-
219 ham, 1984; Danecek *et al*, 2011). To infer the population structure, we used ADMIXTURE for a
220 model-based clustering (Alexander *et al*, 2009) and conducted the analysis with different num-
221 bers of predefined populations ranging from $K = 1$ to $K = 9$ to find the value of K that was most
222 consistent with the data using a cross-validation procedure described in the ADMIXTURE man-
223 ual. To avoid convergence effects we ran ADMIXTURE 10 times with different random seeds
224 for each value of K . As a multivariate clustering method, we applied discriminant analysis of
225 principal components (DAPC) implemented in the R-package `adegenet` (Jombart *et al*, 2010;

Incomplete Amaranth domestication

226 Jombart & Ahmed, 2011) and determined the number of principal components (PCs) used
227 in DAPC with the `optim.a.score` method. To investigate the phylogenetic relationship of the
228 species, we calculated an uncorrected neighbor joining network using the algorithm Neighbor-
229 Net (Bryant & Moulton, 2004) as implemented in the `SplitsTree4` program (Huson & Bryant,
230 2006). The Euclidean distance was calculated from the genetic data to construct a neighbor
231 joining tree, which was bootstrapped 1,000 times with the `pegas` R-package (Paradis *et al*,
232 2004). The migration between genetic groups was modeled with TreeMix (Pickrell & Pritchard,
233 2012). For the TreeMix analysis we used the groups that were identified by ADMIXTURE (K
234 = 5) without an outgroup, and allowed four migration events, as preliminary runs indicates four
235 migration events to be the highest number. The tree was bootstrapped 1,000 times.

236 **Seed color and hundred seed weight**

237 For each accession we calculated the hundred seed weight (HSW) by weighting three samples
238 of 200 seeds. Seed color was determined from digital images taken with a binocular (at 6.5x
239 magnification) and by visual comparison to the GRIN descriptors for amaranth (<http://www.ars-grin.gov/cgi-bin/npgs/html/desclist.pl?159>). There were three colors present in
240 the set of accessions, white, pink, which also indicates a white seed coat and dark brown.
241 To infer how the species, assigned genetic groups or seed color influenced seed size, we
242 conducted an ANOVA. Differences were tested with a LSD test implemented in the R package
243 `agricolae` (<http://tarwi.lamolina.edu.pe/~fmendiburu/>).

245 **Results**

246 **Genome size measurements**

247 Although the genomic history of amaranth species still is largely unknown, genome sizes and
248 chromosome numbers are highly variable within the genus *Amaranthus* (<http://data.kew.org/cvalues/>). We therefore tested whether a change in genome size by polyploidization or
249 large-scale insertions or deletions played a role in the speciation history of *A. caudatus* and the
250 two relatives *A. quitensis* and *A. hybridus* by measuring the genome size of multiple individuals
251

Incomplete Amaranth domestication

252 from all three species with flow cytometry. The mean genome size of *A. caudatus* was 501.93
253 Mbp, and the two relatives did not differ significantly from this value (Table 1) indicating that
254 all measured individuals are diploid and that polyploidization did not play a role in the recent
255 evolution of cultivated amaranth.

Table 1: Genome size of representative group of individuals for each species. There are no significant differences between genome sizes ($p \leq 0.05$). The number of individuals per population is N and SD is the standard deviation for each parameter.

	N	DNA content (pg)	SD	genome size (Mbp)	SD
<i>A. caudatus</i>	22	1.026	0.026	501.93	12.74
<i>A. hybridus</i>	8	1.029	0.025	502.96	12.20
<i>A. quitensis</i>	4	1.021	0.016	499.07	7.91

256 SNP identification by GBS

257 To investigate genome-wide patterns of genetic diversity in *A. caudatus* and its two closest
258 relatives, we genotyped a diverse panel of 119 amaranth accessions from the three species
259 that were initially collected in the Andean region and then obtained from the USDA genebank.
260 The sequencing data generated with a two-enzyme GBS protocol consisted of 210 Mio. raw
261 reads with an average of 1.5 Mio. reads per accession (Supporting information S2). We tested
262 for a lane effect of the Illumina flow cell by sequencing the same 12 individuals on each of the
263 three lanes used for sequencing of the whole sample. An ANOVA of the read number did not
264 show a lane effect (Table S1). Since a high-quality reference genome of an amaranth species
265 was not available, we aligned reads *de novo* within the dataset to unique tags using Stacks
266 (Catchen *et al*, 2011). The total length of the aligned reads was 16.6 Mb, which corresponds
267 to approximately 3.3 % of the *A. caudatus* genome. For SNP calling, reads of each individual
268 were mapped to the aligned tags. SNPs were called with parameters described in Materials
269 and Methods, which resulted in 63,956 SNPs and a mean read depth of 40.28 per site. Since
270 GBS data are characterized by a high proportion of missing values, we removed SNPs with
271 more than 60% of missing values. After this filtering step, we obtained 9,485 biallelic SNPs
272 with an average of 35.3 % missing data for subsequent analyses (Figure S1). The folded site
273 frequency spectrum showed an expected distribution but *A. quitensis* had more sites with low

Incomplete Amaranth domestication

274 frequency due to the restricted number of individuals (Figure S2)

275 **Inference of population structure**

276 To infer the genetic relationship and population structure of the three amaranth species, we
277 used three different methods that included ADMIXTURE, Discriminant Analysis of Principal
278 Components (DAPC) and phylogenetic reconstruction with an uncorrected neighbor-joining
279 network. The ADMIXTURE analysis with three predefined groups ($K = 3$) that corresponds
280 to the number of *Amaranthus* species included in the study did not cluster accessions by their
281 species, but combined the two relatives *A. hybridus* and *A. quitensis* into a single cluster and
282 grouped the *A. caudatus* accessions into two distinct clusters. Higher values of K did not lead
283 to subdivision of the two close relatives into separate groups that correspond to the species
284 assignment (Figure 2), however, they were split according to their geographic origin. Cross-
285 validation showed that $K = 5$ was most consistent with the data (Figure S3), which produced
286 three different groups of *A. caudatus* accessions that included a few accessions from the close
287 relatives, and two clusters that both consist of *A. hybridus* and *A. quitensis* accessions. These
288 two clusters are not separated by the species assignment but by the geographic origin of ac-
289 cessions because the clusters consist of *A. hybridus* and *A. quitensis* accessions from Peru
290 and Ecuador, respectively, which indicates a strong geographic differentiation among the close
291 relatives.

292 The groups of *A. caudatus* accessions also showed a clear geographic differentiation. The first
293 cluster consisted of individuals from Bolivia (Figures 2 and 3; $K = 5$, red color). *A. caudatus*
294 accessions from Peru were split into two clusters of which one predominantly represents a
295 region from North Peru (Huari Province; Figures 2 and 3; $K = 5$, yellow color), whereas the
296 second cluster contains individuals distributed over a wide geographic range that extends from
297 North to South Peru ($K = 5$, green color). Ten *A. caudatus* accessions from the Cuzco region
298 clustered with the three accessions of the close relatives from Peru ($K = 5$, blue color). These
299 ten accessions showed admixture with the other cluster of *A. hybridus/A. quitensis* and with
300 a Peruvian cluster of *A. caudatus*. Accessions that were labeled as 'hybrids' in their passport
301 data, because they express a set of phenotypic traits of different species, clustered with different
302 groups. 'Hybrids' from Bolivia were highly admixed, whereas 'hybrids' from Peru clustered with

Incomplete Amaranth domestication

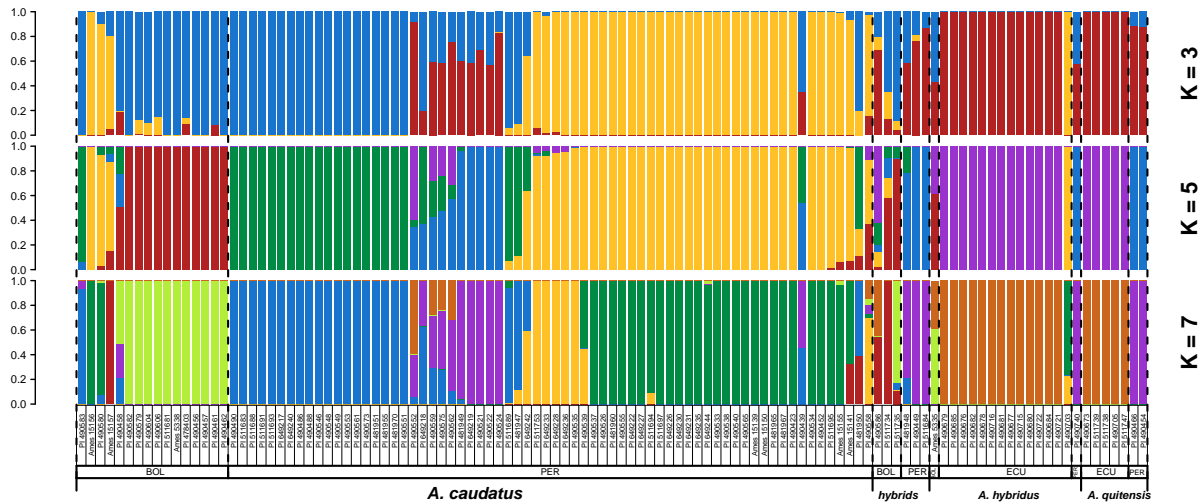


Figure 2: Model based clustering analysis with different numbers of clusters ($K=3, 5, 7$) with ADMIXTURE. The clusters reflect the number of species in the study ($K=3$), the number of single populations (species per country of origin, $K=7$) and the optimal number as determined by cross validation ($K=5$). Individuals are sorted by species and country of origin (BOL=Bolivia, PER = Peru and ECU = Ecuador) as given by their passport data.

303 the close relatives from Peru (Figure 2). Taken together, the population structure inference
304 with ADMIXTURE identified a clear separation between the cultivated *A. caudatus* and its to
305 close relatives, and a high level of genetic differentiation among cultivated amaranths with some
306 evidence for gene flow between groups.

307 The inference of population structure with a discriminant analysis of principal components
308 (DAPC) and Neighbor-Joining network produced very similar results as ADMIXTURE. The first
309 principal component of the DAPC analysis which we used to cluster accessions based on their
310 species explained 96% of the variation and separated the cultivated *A. caudatus* from its two
311 relatives (Figure S4A). In a second DAPC analysis that was based on information on species
312 and country of origin (Figure S4B) the first principal component explained 55% of the variation
313 and separated most cultivated amaranth accessions from the close relatives. The second prin-
314 cipal component explained 35% of the variation and separated the Peruvian from the Bolivian
315 *A. caudatus* accessions.

316 The phylogenetic network outlines the relationships between the different clusters (Figure 4).
317 It shows two distinct groups of mainly Peruvian *A. caudatus* accessions and a group of acces-
318 sions with a wide geographic distribution (Figure 3; green color). The latter is more closely

Incomplete Amaranth domestication

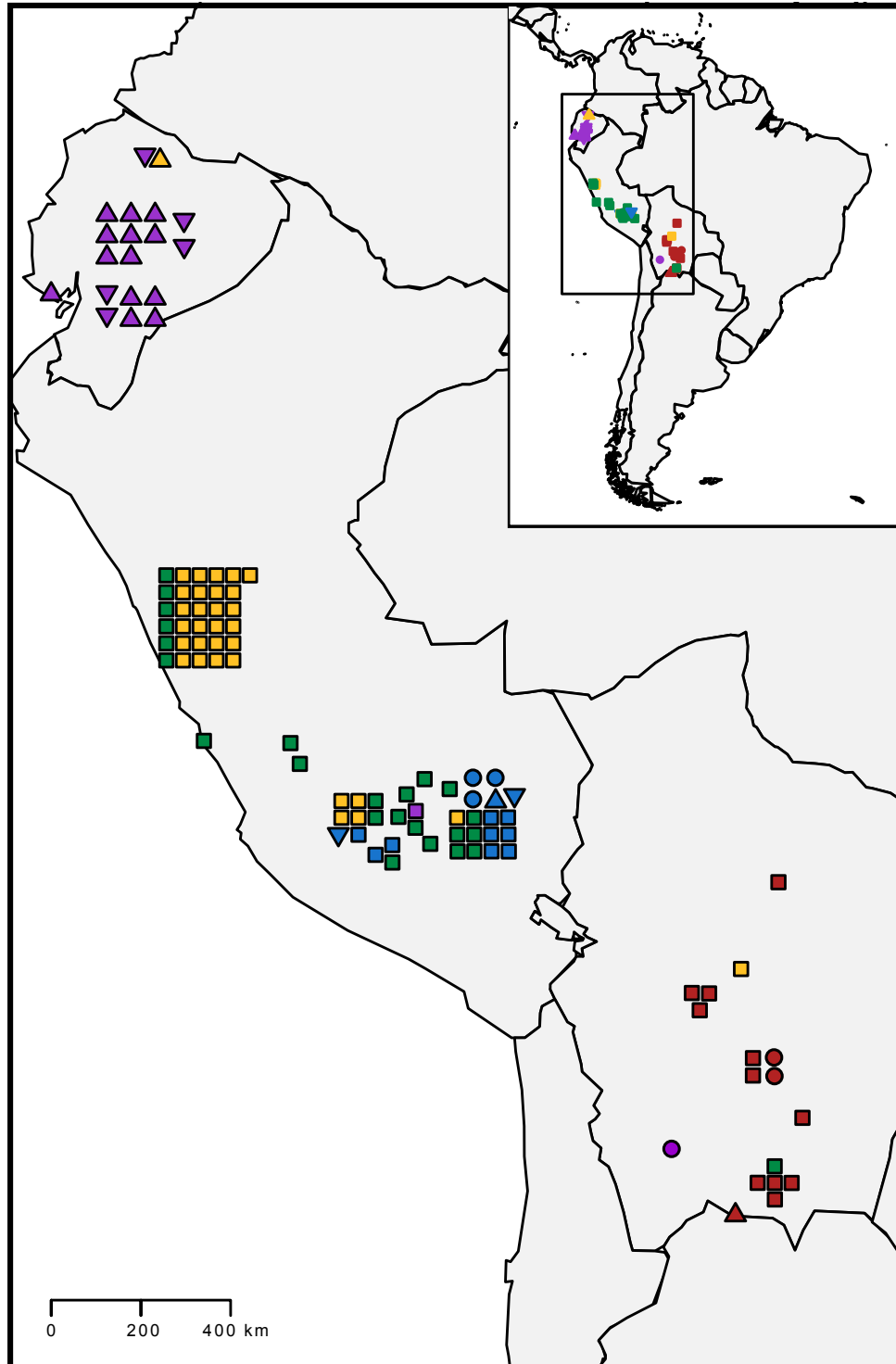


Figure 3: Geographic distribution of accessions for which data was available from passport information. Locations are not exact geographic locations because location data was given as country province. Colors are given by ADMIXTURE with $K=5$ (Figure 2). Species are indicated by shapes. *A. caudatus* (\square), *A. hybridus* (\triangle), *A. quitensis* (∇) and hybrids between species (\circ)

319 related to the Bolivian *A. caudatus* and the close relatives. The strong network structure be-

320 tween these three groups suggests a high proportion of shared ancestral polymorphisms or

Incomplete Amaranth domestication

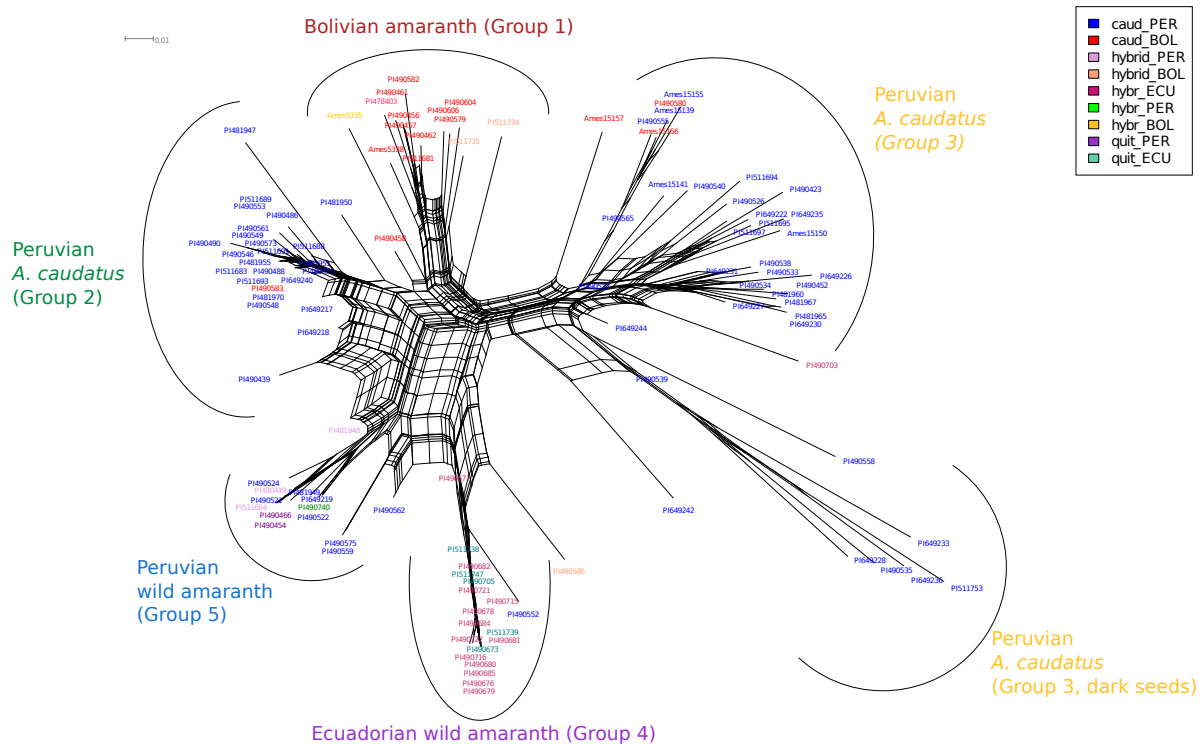


Figure 4: Neighbor-joining network of 113 amaranth accessions from six potential populations. Different colors indicate the species and origin according to gene bank information. *A. caudatus* from Peru (blue) and from Bolivia (red), *A. hybridus* from Ecuador (magenta), from Peru (green) and Bolivia (yellow), *A. quitensis* from Ecuador (turquoise) and Peru (purple) and hybrids between species from Peru (salmon) and Bolivia (light orange). Arches show genetic clusters as inferred with ADMIXTURE ($K = 5$).

321 a high level of recent gene flow. In contrast, *A. caudatus* accessions from Northern Peru are
 322 more strongly separated from the other groups (Figure 3; yellow color) and are split into two
 323 subgroups, of which the smaller one includes only accessions with dark seeds. In a bifurcat-
 324 ing phylogenetic tree, ten cultivated amaranth accessions clustered within the same clade as
 325 the close relatives *A. quitensis* and *A. caudatus* (Figure S5). The same clustering was also
 326 obtained with ADMIXTURE and $K = 7$ (Figure 2).

327 To quantify the level of genetic differentiation between the species and groups within *A. cau-*
 328 *datu*s, we estimated weighted F_{ST} values using the method of Weir and Cockerham (Weir &
 329 Cockerham, 1984). F_{ST} values between *A. caudatus* and *A. hybridus* and *A. quitensis* species
 330 were 0.31 and 0.32, respectively (Table 2), and 0.041 between *A. hybridus* and *A. quitensis*
 331 based on the taxonomic assignment. The latter reflects the high genetic similarity of the acces-
 332 sions from both species observed above. Within *A. caudatus* subpopulations, the F_{ST} between

Incomplete Amaranth domestication

Table 2: Weir and Cockerham weighted F_{ST} estimates between populations based on the taxonomic assignment of their passport data. The group of close relatives are *A. hybridus* and *A. quitensis* taken together.

	F_{ST}
<i>A. caudatus</i> x <i>A. hybridus</i>	0.319
<i>A. caudatus</i> x <i>A. quitensis</i>	0.274
<i>A. caudatus</i> x close relatives	0.322
<i>A. hybridus</i> x <i>A. quitensis</i>	0.041
<i>A. caudatus</i> (PER) x <i>A. caudatus</i> (BOL)	0.132

333 *A. caudatus* populations from Peru and Bolivia was 0.132, three times higher than between *A.*
334 *hybridus* and *A. quitensis*. The above analyses suggested that some individuals may be mis-
335 classified in the passport information, and we therefore calculated F_{ST} values of population sets
336 defined by ADMIXTURE. Although such F_{ST} values are upward biased, they allow to evaluate
337 the relative level of differentiation between groups defined by their genotypes. The comparison
338 of F_{ST} values showed that the three *A. caudatus* groups (groups 1-3) are less distant to the
339 group of *A. quitensis/A. hybridus* accessions from Peru (group 5) than from Ecuador (group 4;
340 Table S2). A TreeMix analysis, which is based on allele frequencies within groups (Figure 5),
341 suggests gene flow from the Peruvian *A. caudatus* (group 2) to *A. quitensis* and *A. hybridus*
342 amaranths from Peru (group 5) and, with a lower confidence level, from *A. quitensis* and *A.*
343 *hybridus* from Ecuador (group 4) into Bolivian *A. caudatus* (group 1), as well as from Bolivian
344 *A. caudatus* to Peruvian *A. caudatus* (Group 2).

345 Analysis of genetic diversity

346 We further investigated whether domestication reduced genetic diversity in cultivated *A. cau-*
347 *datatus* (Table 3). All measures of diversity were higher for *A. caudatus* than its relatives. For
348 example, nucleotide diversity (π) was about two times higher in *A. caudatus* than in the two
349 relatives combined. The diversity values of the accessions classified as 'hybrids' showed in-
350 termediate values between cultivated amaranth and its relatives supporting their hybrid nature.
351 The inbreeding coefficient, F , was highest in the cultivated amaranth but did not differ from the
352 two close relatives if they are combined. In contrast, accessions classified as 'hybrids' and *A.*

Incomplete Amaranth domestication

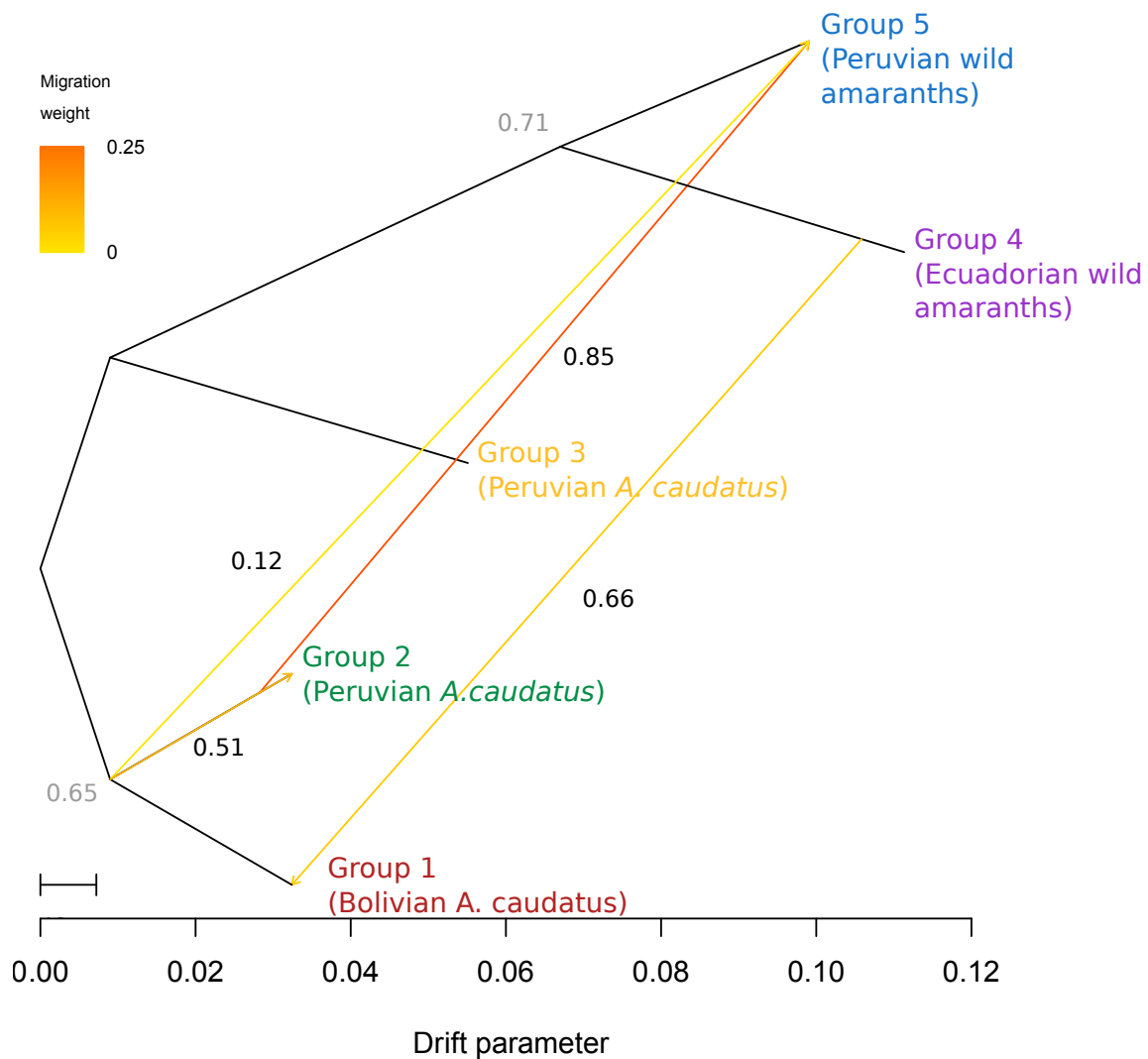


Figure 5: Tree of five genetic clusters of South American amaranths inferred with TreeMix. The genetic clusters which were used to calculate the tree were inferred with ADMIXTURE. Groups 1 to 3 represent *A. caudatus* clusters from Peru and Bolivia, group 4 represents accessions of *A. quitensis* and *A. hybridus* from Ecuador and group 5 wild amaranth from Peru, respectively. The migration events are colored according to their weight. Numbers at branching points and on the migration arrow represent bootstrapping results based on 1,000 runs.

353 *quitensis* had lower inbreeding coefficients. Within the groups of accessions defined by AD-
354 MIXTURE, genetic diversity differed substantially. The close relatives from Ecuador had the
355 lowest ($\pi = 0.00031$) while the group from northern Peru showed the highest level of nucleotide
356 diversity ($\pi = 0.00111$; Table S3). Figure 6 shows that even though the overall diversity of *A.*
357 *caudatus* was higher, a substantial proportion of sites were more diverse in the close relatives
358 ($\pi_{caud} - \pi_{hyb/quit} < 0$; Figure 6).

Incomplete Amaranth domestication

Table 3: Genetic diversity parameters for the cultivated *A. caudatus* and the close relatives *A. hybridus* and *A. quitensis*. π is the nucleotide diversity over all sites, CI_{π} is the 95% confidence interval of π , H_{exp} the mean expected heterozygosity for the variant sites and SD_{H_e} its standard deviation, H_{obs} the mean observed heterozygosity and SD_{H_o} its standard deviation. F is the inbreeding coefficient and SD_F its standard deviation.

Population	N	π	CI_{π}	H_{exp}	SD_{H_e}	H_{obs}	SD_{H_o}	F	SD_F	θ_w
<i>A. caudatus</i>	84	0.00117	± 0.00002	0.175	0.167	0.049	0.140	0.688	0.462	0.00123
<i>A. hybridus</i>	16	0.00061	± 0.00001	0.085	0.135	0.041	0.170	0.679	0.608	0.00073
<i>A. quitensis</i>	7	0.00059	± 0.00001	0.076	0.169	0.040	0.170	0.451	0.763	0.00048
Close relatives combined	23	0.00062	± 0.00002	0.090	0.140	0.041	0.166	0.681	0.591	0.00070
Hybrids	6	0.00091	± 0.00001	0.112	0.179	0.060	0.173	0.436	0.645	0.00107

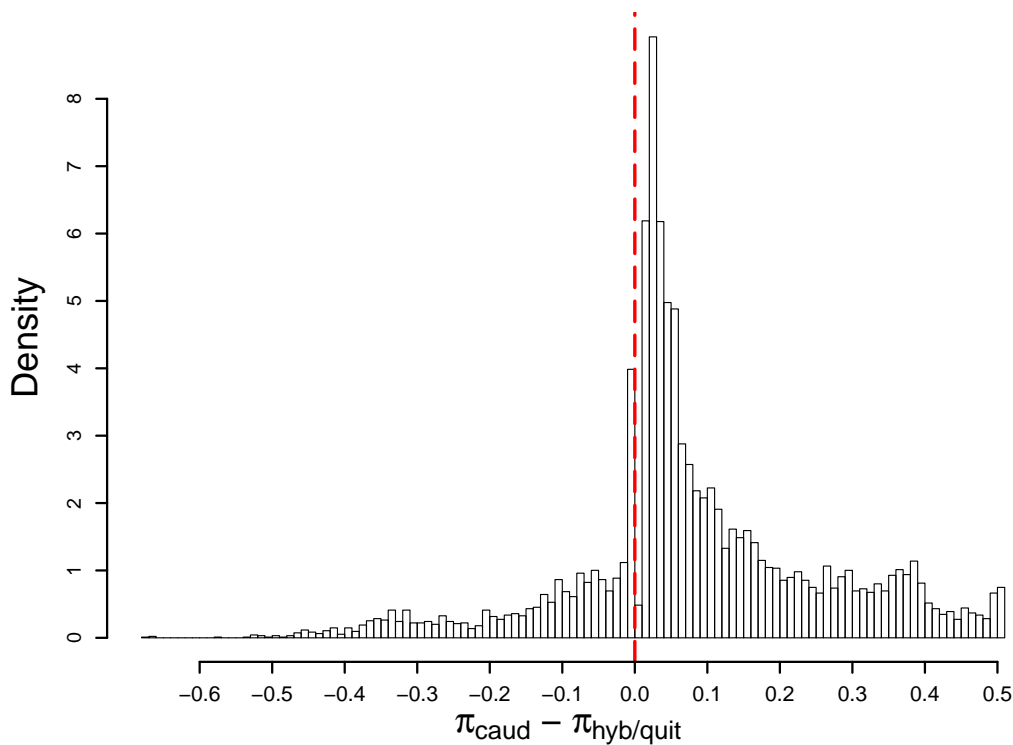


Figure 6: Per site difference in nucleotide diversity (π) between cultivated amaranth (*A. caudatus*) and the close relatives *A. hybridus* and *A. quitensis*.

359 Seed color and seed size as potential domestication traits.

360 In grain crops, grain size and seed color are important traits for selection and likely played
 361 a central role in domestication of numerous plants (Abbo *et al*, 2014; Hake & Ross-Ibarra,

Incomplete Amaranth domestication

2015). To investigate whether these two traits are part of the domestication syndrome in grain amaranth, we compared the predominant seed color of the different groups of accessions and measured their seed size. The seeds could be classified into three colors, white, pink and brown. The white and pink types have both a white seed coat, but the latter has red cotyledons that are visible through the translucent seed coat. A substantial number of seed samples (19) from the genebank contained seeds of other color up to a proportion of 20%. One *A. caudatus* accession from Peru (PI 649244) consisted of 65% dark seeds and 35% white seeds in the sample. No accession from the two close relatives *A. hybridus* and *A. quitensis*, or from the hybrid accessions had white seeds, whereas the majority (74%) of *A. caudatus* accessions had white (70%) or pink (4%) seeds, and the remaining (26%) brown seeds (Figure 7 A). We also compared the seed color of groups defined by ADMIXTURE ($K = 5$; Figure 2), which reflect their genetic relationship and may correct for mislabeling of accessions (Figure 7 B). No group had only white seeds, but clusters consisting mainly of *A. hybridus* and *A. quitensis* had no white seeds at all. In contrast to seed color, the hundred seed weight (HSW) of the different *Amaranthus* species did not significantly differ between cultivated *A. caudatus* and the two relatives. The mean HSW of *A. caudatus* was 0.056 g and slightly higher than the HWS of *A. hybridus* (0.051 g) and *A. quitensis* (0.050 g; Figure 7 C and Table S4). Among the groups identified by ADMIXTURE ($K = 5$), one group showed a significantly higher HSW than the other groups, while the other four groups did not differ in their seed size. The group with the higher HSW consisted mainly of Bolivian *A. caudatus* accessions and had a 21 % and 35 % larger HSW than the two groups consisting mainly of Peruvian *A. caudatus* accessions, respectively (Figure 7 D). An ANOVA also revealed that seed color has an effect on seed size because white seeds are larger than dark seeds (Table 4).

Table 4: Analysis of variance for the hundred seed weight in dependence of the seed color and population as determined by ADMIXTURE

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Seed color	2	0.000657	0.0003285	4.657	0.0116 *
Group	4	0.003151	0.0007877	11.165	1.46e-07 ***
Seed color:Group	2	0.000042	0.0000209	0.297	0.7440
Residuals	103	0.007266	0.0000705		

Incomplete Amaranth domestication

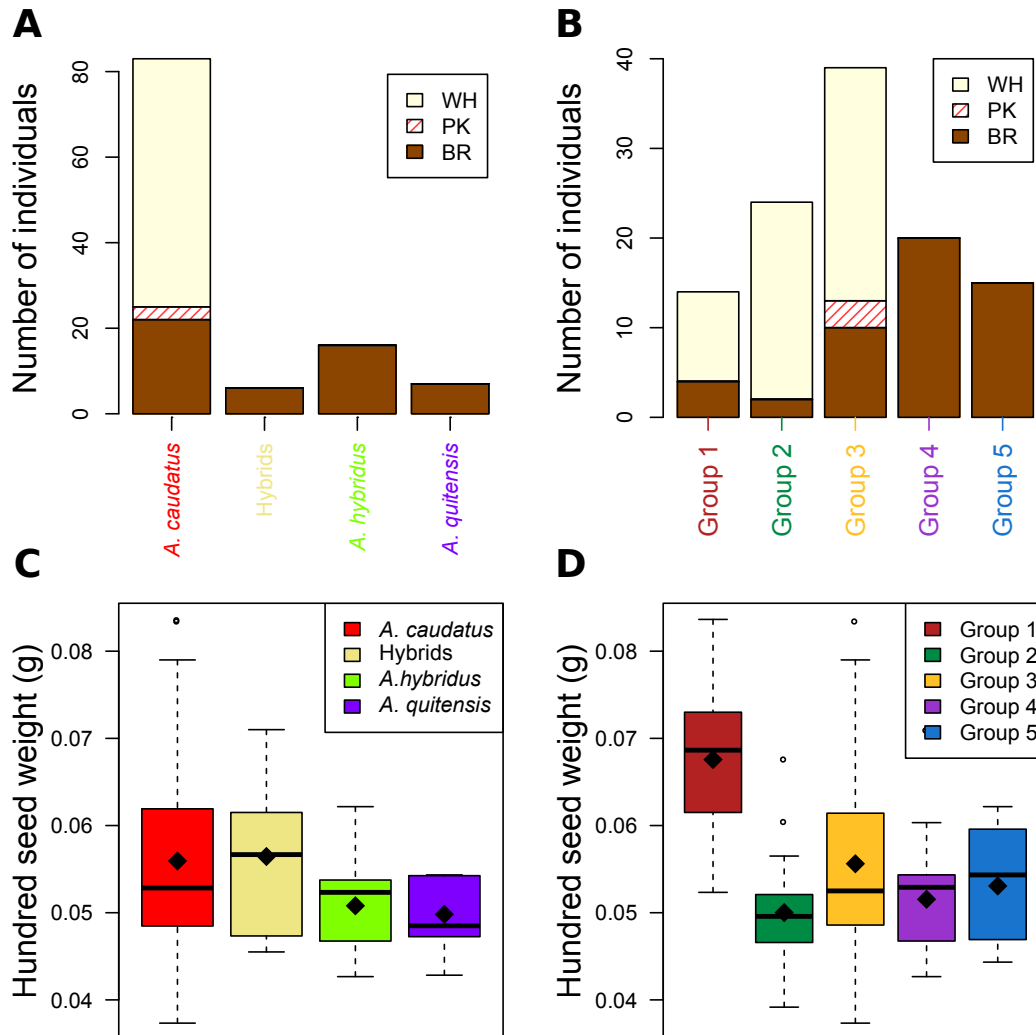


Figure 7: Predominant seed color (**A,B**) and hundred seed weight (**C,D**) by *Amaranthus* species (**A,C**) and groups identified with ADMIXTURE for $K = 5$ (**B,D**) where group 1 (red) resembles *A. caudatus* from Bolivia, group 2 (green) and 3 (yellow) *A. caudatus* from Peru, group 4 (purple) represents the close relatives *A. quitensis* and *A. hybridus* from Ecuador and group 5 (blue) from Peru, respectively. Seed colors were white (WH), pink (PK) and dark brown (BR). While there were no significant differences in seed size between the species, group 1 had significantly higher hundred seed weight ($p \leq 0.05$) than the other groups.

385 Discussion

386 Genotyping-by-sequencing of amaranth species

387 The genotyping of cultivated amaranth *A. caudatus* and two close relatives *A. quitensis* and
388 *A. hybridus* revealed a strong genetic differentiation between both groups and a high level of
389 genetic differentiation within cultivated *A. caudatus*. We based our sequence assembly and
390 SNP calling on a *de novo* assembly of GBS data with Stacks because currently no high quality
391 reference sequence of these species is available. Stacks allows SNP calling without a ref-
392 erence genome by constructing a reference catalog from the data and includes all reads in
393 the analysis (Catchen *et al*, 2011). *De novo* assembled fragments without a reference are
394 unsorted and can not be used to investigate genetic differentiation along along the genomic
395 regions, but they are suitable for analysing genetic diversity and population structure (Catchen
396 *et al*, 2013). GBS produces a large number of SNPs (Poland *et al*, 2012; Huang *et al*, 2014),
397 albeit with a substantial proportion of missing values. Missing data lead to biased estimators of
398 population parameters such as π and θ_w (Arnold *et al*, 2013) and need to be accounted for if
399 different studies are compared. Additionally, variable error rates in different GBS data sets can
400 inflate differentiation estimates (Mastretta-Yanes *et al*, 2015). The comparison of accessions
401 and groups within a study is possible, however, if all individuals were treated with the same
402 experimental protocol. We filtered out sites with high levels of missing values to obtain a robust
403 dataset for subsequent population genomic analysis. The SNPs were called based on the total
404 sample without accounting for the species which should not bias diversity estimates. Since a
405 smaller sample from the close relatives may underestimate their diversity compared to culti-
406 vated *A. caudatus*, we compared diversity estimates by repeated random sampling of 23 out
407 of 84 *A. caudatus* accessions and calculating π from the smaller sample. Diversity estimates
408 of the smaller *A. caudatus* did not differ from the full sample and estimates were in all cases
409 higher than in the close relatives (Figure S6). We conclude that the different sample sizes of
410 the two groups do not introduce a bias on diversity estimates.

Incomplete Amaranth domestication

411 **Genetic relationship of *A. quitensis* and *A. hybridus***

412 Coons (1978) suggested that *A. quitensis* is the same species as *A. hybridus*, but in the
413 genebank passport data *A. quitensis* is still considered as a separate species. The taxo-
414 nomic differentiation between the two species rests on two minor morphological trait, namely
415 the shape of the tepals and the short utricles, which are very small and prone to misidentifi-
416 cation (Sauer, 1967; Adhikary & Pratt, 2015). The high phenotypic similarity of *A. quitensis*
417 and *A. hybridus* is supported by the GBS data because accessions from the two species are
418 closely related. They are not separated by their species assignment but cluster into two groups
419 that both consist of accessions from the two species and reflect their geographic origin from
420 Peru and Ecuador, respectively. The F_{ST} value between *A. quitensis* and *A. hybridus* was lower
421 than between the two *A. caudatus* groups from Peru and Bolivia (Tables 2 and S2). The highly
422 similar genome sizes of all three diploid species is consistent with genetic relationship inferred
423 from the GBS data and indicates that large-scale genomic changes like polyploidization events
424 did not occur in the recent history of these species. For comparison, other species in the genus
425 *Amaranthus* have very different genome sizes due to variation in chromosome numbers and
426 ploidy levels (Baohua & Xuejie, 2002; Rayburn *et al*, 2005).

427 In contrast to our analysis, Kietlinski *et al* (2014) found stronger evidence for a genetic differen-
428 tiation between *A. hybridus* and *A. quitensis* based on the 11 SSR markers. However, their data
429 also show that both species are distinct groups that do not cluster their species assignment but
430 by geographic origin. These differences may result from the different marker types (SNPs vs.
431 SSRs) and a different sample composition because our sample consists of accessions from
432 the Andean region, whereas Kietlinski *et al.* included putative wild amaranth accessions with
433 little geographic overlap between the two species. The groups of *A. hybridus* and *A. quitensis*
434 accessions from Peru and Ecuador show a high level of differentiation ($F_{ST} = 0.579$; Table S2),
435 which is similar to the differentiation between one of two Peruvian *A. caudatus* groups and the
436 *A. hybridus/A. quitensis* accessions from Peru ($F_{ST} = 0.553$). Although the sample size of *A.*
437 *quitensis* and *A. hybridus* is small, genetic differentiation between species should be stronger
438 than between individuals within species in the ADMIXTURE and phylogenetic analyses. In
439 summary, our analysis and the work of Kietlinski *et al* (2014) show that *A. quitensis* and *A. hy-*
440 *bridus* do not have a simple genetic relationship that follows species assignment. The high level

Incomplete Amaranth domestication

441 of intraspecific differentiation in both cultivated amaranth and their relatives is relevant for in-
442 vestigating domestication because the genetic distance between groups of cultivated amaranth
443 is related to the geographic distance of the putative wild ancestors. Therefore, future studies
444 of these two close relatives of the grain amaranths should include large number of accessions
445 from the whole species range to model genetic differentiation within the two species as well as
446 the relationship between species.

447 **Diversity of South American amaranth**

448 In numerous crops, domestication was associated with a decrease in genome-wide levels of
449 diversity due to bottleneck effects and strong artificial selection of domestication traits (Gepts,
450 2014). Under the assumption that the cultivated grain amaranth *A. caudatus* was domesti-
451 cated, genetic diversity should be reduced compared to the two close relatives. In contrast,
452 the overall genetic diversity in our sample of cultivated amaranths was higher than in the two
453 close relatives. The distribution of diversity between the GBS fragments includes genomic
454 regions with reduced diversity in *A. caudatus*, which may reflect selection in some genomic
455 regions (Figure 6). Without a reference genome it is not possible to position reads on a map to
456 identify genomic regions that harbor putative targets of selection based on an inference of the
457 demographic history. Despite the indirect phenotypic evidence for selection, the higher genetic
458 diversity of cultivated grain amaranth may result from a strong gene flow between cultivated
459 amaranths and its relatives. Gene flow between different amaranth species has been observed
460 before (Trucco *et al*, 2005) and is also consistent with the observation of six highly admixed
461 accessions classified as 'hybrids' in the passport data and which appear to be interspecific hy-
462 brids (Figure 2 and Table 3). Gene flow between *A. caudatus* and the relatives *A. quitensis* and
463 *A. hybridus* in different areas of the distribution range, not only from populations included in this
464 study, could explain a higher genetic diversity in cultivated amaranth. This is also consistent
465 with the strong network structure (Figure 4) and the TreeMix analysis (Figure 5). In summary,
466 cultivated *A. caudatus* is unusual in its higher overall genetic diversity compared to populations
467 of its putative wild ancestors originating from the same geographic region. The high genetic
468 diversity of *A. caudatus* is in contrast to other domesticated plants and suggests that a domes-
469 tication bottleneck in its cultivation history absent (i.e., no domestication), very weak or masked

Incomplete Amaranth domestication

470 by recurrent gene flow. We consider these results to be robust, because in comparison to
471 previous work (Maughan *et al*, 2009, 2011; Khaing *et al*, 2013; Jimenez *et al*, 2013; Kietlinski
472 *et al*, 2014), our study includes a larger number of accessions and more genetic markers. This
473 allowed us to assess the genetic diversity and population structure of South American ama-
474 ranth on a genome-wide basis, but we note that a more complete geographic sampling of all
475 cultivated amaranths and their relatives is required for a complete understanding of amaranth
476 history.

477 **Evidence for a weak domestication syndrome**

478 Despite their long history of cultivation, diverse uses for food and feed and a high importance
479 for the agriculture of ancient cultures (i.e., *A. hypochondriacus* during the Aztec period), grain
480 amaranths do not display the classical domestication syndrome as strongly as other crops
481 (Sauer, 1967). Cultivated amaranth shows morphological differentiation from putative wild an-
482 cestors like larger and more compact inflorescences (Sauer, 1967) and a level of genetic dif-
483 ferentiation (Table 2) which is comparable to other domesticated crops and their wild relatives
484 (Sunflower: $F_{ST}=0.22$ (Mandel *et al*, 2011); common bean: 0.1-0.4 (Papa *et al*, 2005), pi-
485 geonpea: 0.57-0.82 (Kassa *et al*, 2012)). However, the numerous amaranth flowers mature
486 asynchronously and produce very small seeds that are shattered (Brenner *et al*, 2010). All
487 putative wild amaranths have dark brown seeds, whereas the predominant seed color of cul-
488 tivated grain amaranth is white, which suggests that selection for seed color played a role in
489 the history of the latter. Dark-seeded accessions are present in all three groups of *A. caudatus*
490 defined by the genotypic data, which indicates that white seed color is not a fixed trait. Seed
491 sizes between cultivated amaranth and its relatives are not significantly different with the excep-
492 tion of white-seeded *A. caudatus* accessions from Bolivia (Figure 7), which have larger seeds.
493 The larger seeds in this group and of white seeds in general (Table 4) indicates past selection
494 for domestication-related traits, but only in specific geographic regions or in certain types of
495 amaranth, and not in the whole cultivated crop species.

496 These findings suggest that some selection occurred in the history of amaranth cultivation that
497 may reflect domestication. Possible explanations for the incomplete fixation of domestication
498 traits in South American grain amaranth include weak selection, genetic constraints or ongoing

Incomplete Amaranth domestication

499 gene flow. First, weak selection of putative domestication traits indicate that they were not es-
500 sential for cultivation. Although white seeds are predominant in cultivated amaranth and most
501 likely a selected trait, other seed colors may have been preferred for different uses. Since ama-
502 ranths are also an important leaf vegetable in Mexico, the grain use of *A. caudatus* may not
503 have been a main target of selection during domestication, thereby allowing a diversity of seed
504 traits due to weak or incomplete selection. Second, genetic constraints may limit phenotypic
505 variation in domestication traits. In contrast to genes with strong pleiotropic effects or epistatic
506 interactions, domestication genes that are part of simple molecular pathways, have minimal
507 pleiotropic effects, and show standing functional genetic variation have a higher chance of fix-
508 ation by selection (Doebley *et al*, 2006; Lenser & Theißen, 2013). Numerous genes with these
509 characteristics were cloned and characterized in major crops like rice, barley and maize. They
510 contribute to the distinct domestication syndrome such as a loss of seed shattering, larger seed
511 size and compact plant architecture. The molecular genetics of amaranth domestication traits
512 remains unknown, but the absence a strong domestication syndrome may reflect genetic con-
513 straints despite a long period of cultivation. A third explanation is ongoing gene flow between
514 cultivated amaranth and its relatives that may prevent or delay the formation of a distinct do-
515 mestication syndrome and contributes to the high genetic diversity (Table 3), similar seed size
516 (Figure 7 C), and the presence of dark seeds (Figure 7) in cultivated amaranth. Both historical
517 and ongoing gene flow are likely because amaranth has an outcrossing rate between 5% and
518 30% (Jain *et al*, 1982; Stetter *et al*, 2016). In South America, cultivated amaranth and its rela-
519 tives are sympatric over wide areas and the latter were tolerated in the fields and home gardens
520 with *A. caudatus* (Sauer, 1967), where they may have intercrossed. Gene flow between wild
521 and domesticated plants has also been observed in maize and teosinte in the Mexican high-
522 lands, but did not have a major influence on the maize domestication syndrome (Hufford *et al*,
523 2013). Further support for ongoing gene flow in amaranth is given by the presence of hybrids
524 and admixed accessions in our sample with evidence for genetic admixture and dark seeds
525 that demonstrate the phenotypic effects of introgression. Since the dark seed color is dominant
526 over white color (Kulakow *et al*, 1985) and *A. caudatus* is predominantly self-pollinating, dark
527 seeds could have efficiently removed by selection despite gene flow. Additionally, amaranth
528 was grown in small plots in the Andean highlands, which favors fixation of traits (Kauffman &
529 Weber, 1990). Thus, gene flow is a plausible explanation for the absence of a distinct domesti-

Incomplete Amaranth domestication

530 cation syndrome.

531 Although our sample does not include *A. hypochondriacus* or *A. cruentus* accessions, our data
532 are consistent with the model by Kietlinski *et al* (2014) who proposed the domestication of
533 *A. caudatus* and *A. hypochondriacus* from different *A. hybridus* lineages in Central and South
534 America (Figure 1D). Gene flow between *A. caudatus* and its close relative *A. quitensis* in the
535 Southern distribution range (Peru and Bolivia) may explain the higher genetic diversity of the
536 latter despite a strong genetic differentiation.

537 Conclusions

538 The genotypic and phenotypic analysis of cultivated South American grain amaranth and its
539 close relatives suggests that *A. caudatus* is an incompletely domesticated crop species. Key
540 domestication traits such as the shape of inflorescences, seed shattering and seed size are
541 rather similar between cultivated amaranths and their close relatives and there is strong evi-
542 dence of ongoing gene flow between these species despite selection for domestication-related
543 traits like white seeds. Grain amaranth is an ancient crop of the Americas but genomic and
544 phenotypic signatures of domestication differ from other, highly domesticated crops that orig-
545 inated from single domestication events like maize (Hake & Ross-Ibarra, 2015). In contrast,
546 the history of cultivated amaranth may include multiregional, multiple and incomplete domes-
547 tication events with frequent and ongoing gene flow from sympatric relatives, which is more
548 similar to the history of species like rice, apple or barley (Londo *et al*, 2006; Cornille *et al*,
549 2012; Poets *et al*, 2015). The classical model of a single domestication in a well-defined center
550 of domestication may not sufficiently reflect the history of numerous ancient crops. Our study
551 further highlights the importance of a comprehensive sampling to study the domestication of
552 amaranth. The three cultivated amaranths and all close relatives should be included in further
553 studies for a full understanding of amaranth domestication and its broader implications for crop
554 plant domestication.

Incomplete Amaranth domestication

555 **Acknowledgments**

556 We thank David Brenner (USDA-ARS) and Julie Jacquemin for discussions and Elisabeth
557 Kokai-Kota for support with the GBS library preparation and sequencing. The work was funded
558 by an endowment of the Stifterverband für die Deutsche Wissenschaft to K. J. S.

Incomplete Amaranth domestication

References

- 559
- 560 Abbo S, van Oss RP, Gopher A, Saranga Y, Ofner I, Peleg Z (2014) Plant domestication ver-
561 sus crop evolution: a conceptual framework for cereals and grain legumes. *Trends in Plant*
562 *Science*, **19**, 351 – 360.
- 563 Adhikary D, Pratt DB (2015) Morphologic and taxonomic analysis of the weedy and cultivated
564 *Amaranthus hybridus* species complex. *Systematic Botany*, **40**, 604–610.
- 565 Alexander DH, Novembre J, Lange K (2009) Fast model-based estimation of ancestry in unre-
566 lated individuals. *Genome Research*, **19**, 1655–64.
- 567 Arnold B, Corbett-Detig RB, Hartl D, Bomblies K (2013) RADseq underestimates diversity and
568 introduces genealogical biases due to nonrandom haplotype sampling. *Molecular Ecology*,
569 **22**, 3179–90.
- 570 Arreguez GA, Martínez JG, Ponessa G (2013) *Amaranthus hybridus* L. ssp. *hybridus* in an ar-
571 chaeological site from the initial mid-Holocene in the Southern Argentinian Puna. *Quaternary*
572 *International*, **307**, 81–85.
- 573 Baohua S, Xuejie Z (2002) Chromosome numbers of 14 species in *Amaranthus* from China.
574 *Acta Phytotaxonomica Sinica*, **40**, 428–432.
- 575 Begun DJ, Holloway AK, Stevens K, *et al* (2007) Population genomics: whole-genome analysis
576 of polymorphism and divergence in *Drosophila simulans*. *PLoS Biology*, **5**, e310.
- 577 Besnard G, Rubio de Casas R (2015) Single vs multiple independent olive domestications: the
578 jury is (still) out. *New Phytologist*.
- 579 Brenner DM, Baltensperger DD, Kulakow PA, *et al* (2010) Genetic Resources and Breeding of
580 *Amaranthus*. In *Plant Breeding Reviews*, pp. 227–285. John Wiley & Sons, Inc.
- 581 Bryant D, Moulton V (2004) Neighbor-net: an agglomerative method for the construction of
582 phylogenetic networks. *Molecular Biology and Evolution*, **21**, 255–65.
- 583 Catchen J, Hohenlohe PA, Bassham S, Amores A, Cresko WA (2013) Stacks: an analysis tool
584 set for population genomics. *Molecular Ecology*, **22**, 3124–40.
- 585 Catchen JM, Amores A, Hohenlohe P, Cresko W, Postlethwait JH (2011) Stacks: building and
586 genotyping loci *de novo* from short-read sequences. *G3 (Bethesda, Md.)*, **1**, 171–82.
- 587 Coons MP (1978) The status of *Amaranthus hybridus* L. in South America. *Cienc. Nat.*, **18**.
- 588 Cornille A, Gladieux P, Smulders MJM, *et al* (2012) New insight into the history of domesti-
589 cated apple: Secondary contribution of the european wild apple to the genome of cultivated
590 varieties. *PLoS Genetics*, **8**, e1002703.
- 591 Costea M, DeMason D (2001) Stem morphology and anatomy in *Amaranthus* L. (*Amaran-*
592 *thaceae*), taxonomic significance. *Journal of the Torrey Botanical Society*, **128**, 254–281.
- 593 Cruaud A, Gautier M, Galan M, *et al* (2014) Empirical assessment of RAD sequencing for
594 interspecific phylogeny. *Molecular Biology and Evolution*, **31**, 1272–4.
- 595 Danecek P, Auton A, Abecasis G, *et al* (2011) The variant call format and VCFtools. *Bioinform-*
596 *atics (Oxford, England)*, **27**, 2156–8.

Incomplete Amaranth domestication

- 597 Doebley JF, Gaut BS, Smith BD (2006) The molecular genetics of crop domestication. *Cell*,
598 **127**, 1309–21.
- 599 Dolezel J, Bartos J, Voglmayr H, Greilhuber J (2003) Nuclear DNA content and genome size of
600 trout and human. *Cytometry. Part A : The Journal of the International Society for Analytical*
601 *Cytology*, **51**, 127–8; author reply 129.
- 602 Dolezel J, Sgorbati S, Lucretti S (1992) Comparison of three DNA fliiorocliromes for flow cyto-
603 metric estimatioe of nuclear DNA conteet in plants. *Physiologia Plantarum*, **85**, 625–631.
- 604 Elshire RJ, Glaubitz JC, Sun Q, *et al* (2011) A robust, simple genotyping-by-sequencing (gbs)
605 approach for high diversity species. *PLoS ONE*, **6**, e19379.
- 606 Gepts P (2014) The contribution of genetic and genomic approaches to plant domestication
607 studies. *Current Opinion in Plant Biology*, **18**, 51–9.
- 608 Hake S, Ross-Ibarra J (2015) Genetic, evolutionary and plant breeding insights from the do-
609 mestication of maize. *eLife*, **4**, 1–8.
- 610 Huang YF, Poland JA, Wight CP, Jackson EW, Tinker NA (2014) Using genotyping-by-
611 sequencing (GBS) for genomic discovery in cultivated oat. *PLoS ONE*, **9**, e102448.
- 612 Hufford MB, Lubinsky P, Pyhäjärvi T, Devengenzo MT, Ellstrand NC, Ross-Ibarra J (2013) The
613 genomic signature of crop-wild introgression in maize. *PLoS Genetics*, **9**, e1003477.
- 614 Huson DH, Bryant D (2006) Application of phylogenetic networks in evolutionary studies.
615 *Molecular Biology and Evolution*, **23**, 254–67.
- 616 Jain S, Hauptil H, Vaidya K (1982) Outcrossing rate in grain amaranths. *Journal of Heredity*,
617 **73**, 71–72.
- 618 Jimenez FR, Maughan PJ, Alvarez A, *et al* (2013) Assessment of genetic diversity in Peru-
619 vian amaranth (*Amaranthus caudatus* and *A. hybridus*) germplasm using single nucleotide
620 polymorphism markers. *Crop Science*, **53**, 532.
- 621 Jombart T, Ahmed I (2011) adegenet 1.3-1: new tools for the analysis of genome-wide SNP
622 data. *Bioinformatics (Oxford, England)*, **27**, 3070–1.
- 623 Jombart T, Devillard S, Balloux F (2010) Discriminant analysis of principal components: a new
624 method for the analysis of genetically structured populations. *BMC Genetics*, **11**, 94.
- 625 Kassa MT, Penmetsa RV, Carrasquilla-Garcia N, *et al* (2012) Genetic patterns of domestication
626 in pigeonpea (*Cajanus cajan* L. millsp.) and wild cajanus relatives. *PLoS ONE*, **7**.
- 627 Kauffman CS, Weber LE (1990) Grain Amaranth. In *Advances in new crops* (edited by J Janick,
628 JE Simon), pp. 127–139. Timber Press, Portland.
- 629 Khaing AA, Moe KT, Chung JW, Baek HJ, Park YJ (2013) Genetic diversity and population
630 structure of the selected core set in *Amaranthus* using SSR markers. *Plant Breeding*, **132**,
631 165–173.
- 632 Kietlinski KD, Jimenez F, Jellen EN, Maughan PJ, Smith SM, Pratt DB (2014) Relationships
633 between the weedy (*Amaranthaceae*) and the grain amaranths. *Crop Science*, **54**, 220.
- 634 Kulakow P, Hauptli H, Jain S (1985) Genetics of grain amaranths I. mendelian analysis of six
635 color characteristics. *Journal of Heredity*, **76**, 27–30.

Incomplete Amaranth domestication

- 636 Lenser T, Theißen G (2013) Molecular mechanisms involved in convergent crop domestication.
637 *Trends in Plant Science*, pp. 1–11.
- 638 Londo JP, Chiang YC, Hung KH, Chiang TY, Schaal BA (2006) Phylogeography of asian wild
639 rice, *Oryza rufipogon*, reveals multiple independent domestications of cultivated rice, *Oryza*
640 *sativa*. *Proceedings of the National Academy of Sciences*, **103**, 9578–9583.
- 641 Mallory MA, Hall RV, McNabb AR, Pratt DB, Jellen EN, Maughan PJ (2008) Development and
642 characterization of microsatellite markers for the grain amaranths. *Crop Science*, **48**, 1098.
- 643 Mandel JR, Dechaine JM, Marek LF, Burke JM (2011) Genetic diversity and population struc-
644 ture in cultivated sunflower and a comparison to its wild progenitor, *Helianthus annuus* L.
645 *Theoretical and Applied Genetics*, **123**, 693–704.
- 646 Mastretta-Yanes A, Arrigo N, Alvarez N, Jorgensen TH, Piñero D, Emerson BC (2015) Re-
647 striction site-associated dna sequencing, genotyping error estimation and de novo assembly
648 optimization for population genetic inference. *Molecular Ecology Resources*, **15**, 28–41.
- 649 Maughan P, Smith S, Fairbanks D, Jellen E (2011) Development, characterization, and linkage
650 mapping of single nucleotide polymorphisms in the grain amaranths (*Amaranthus sp.*). *The*
651 *Plant Genome Journal*, **4**, 92.
- 652 Maughan PJ, Yourstone SM, Jellen EN, Udall Ja (2009) SNP discovery via genomic reduction,
653 barcoding, and 454-pyrosequencing in amaranth. *The Plant Genome Journal*, **2**, 260.
- 654 Meyer RS, DuVal AE, Jensen HR (2012) Patterns and processes in crop domestication: an
655 historical review and quantitative analysis of 203 global food crops. *The New Phytologist*,
656 **196**, 29–48.
- 657 Morris GP, Ramu P, Deshpande SP, *et al* (2013) Population genomic and genome-wide asso-
658 ciation studies of agroclimatic traits in sorghum. *Proceedings of the National Academy of*
659 *Sciences*, **110**, 453–458.
- 660 Nabholz B, Sarah G, Sabot F, *et al* (2014) Transcriptome population genomics reveals severe
661 bottleneck and domestication cost in the african rice (*Oryza glaberrima*). *Molecular Ecology*,
662 **23**, 2210–27.
- 663 Olsen KM, Wendel JF (2013) A bountiful harvest: genomic insights into crop domestication
664 phenotypes. *Annual Review of Plant Biology*, **64**, 47–70.
- 665 Papa R, Acosta J, Delgado-Salinas a, Gepts P (2005) A genome-wide analysis of differentiation
666 between wild and domesticated *Phaseolus vulgaris* from Mesoamerica. *Theoretical and*
667 *Applied Genetics*, **111**, 1147–1158.
- 668 Paradis E, Claude J, Strimmer K (2004) APE: Analyses of phylogenetics and evolution in R
669 language. *Bioinformatics*, **20**, 289–290.
- 670 Pickrell JK, Pritchard JK (2012) Inference of population splits and mixtures from genome-wide
671 allele frequency data. *PLoS Genetics*, **8**.
- 672 Poets AM, Fang Z, Clegg MT, Morrell PL (2015) Barley landraces are characterized by geo-
673 graphically heterogeneous genomic origins. *Genome Biology*, **16**, 173.
- 674 Poland JA, Brown PJ, Sorrells ME, Jannink JL (2012) Development of high-density genetic
675 maps for barley and wheat using a novel two-enzyme genotyping-by-sequencing approach.
676 *PLoS ONE*, **7**, e32253.

Incomplete Amaranth domestication

- 677 Rayburn AL, McCloskey R, Tatum TC, Bollero GA, Jeschke MR, Tranel PJ (2005) Genome size
678 analysis of weedy species. *Crop Science*, **45**, 2557.
- 679 Saghai-Marooif MA, Soliman KM, Jorgensen RA, Allard RW (1984) Ribosomal DNA spacer-
680 length polymorphisms in barley: mendelian inheritance, chromosomal location, and popu-
681 lation dynamics. *Proceedings of the National Academy of Sciences of the United States of*
682 *America*, **81**, 8014–8018.
- 683 Sang T, Li J (2013) Molecular genetic basis of the domestication syndrome in cereals. In
684 *Cereal Genomics II* (edited by PK Gupta, RK Varshney), pp. 319–340. Springer Netherlands,
685 Dordrecht.
- 686 Sauer J (1967) The grain amaranths and their relatives: a revised taxonomic and geographic
687 survey. *Annals of the Missouri Botanical Garden*, **54**, 103–137.
- 688 Stetter MG, Schmid KJ (2016) Phylogenetic relationships and genome size evolution within the
689 genus *Amaranthus* indicate the ancestors of an ancient crop. *bioRxiv*.
- 690 Stetter MG, Zeitler L, Steinhaus A, Kroener K, Biljecki M, Schmid KJ (2016) Crossing meth-
691 ods and cultivation conditions for rapid production of segregating populations in three grain
692 amaranth species. *Frontiers in Plant Science*, **7**, 816.
- 693 Team RC (2014) *R: A Language and Environment for Statistical Computing*. R Foundation for
694 Statistical Computing, Vienna, Austria.
- 695 Trucco F, Jeschke MR, Rayburn AL, Tranel PJ (2005) *Amaranthus hybridus* can be pollinated
696 frequently by *A. tuberculatus* under field conditions. *Heredity*, **94**, 64–70.
- 697 Weir B, Cockerham C (1984) Estimating F-statistics for the analysis of population structure.
698 *Evolution*, **38**, 1358–1370.

699 Data Accessibility

700 The original genomic data are available on the European Nucleic Archive (ENA) under the
701 study accession number PRJEB13013. Scripts and phenotypic raw data are available under
702 <http://dx.doi.org/10.5061/dryad.m5kk3> on Dryad (<http://datadryad.org/>).

703 Author Contributions

704 M.G.S. and K.J.S. designed research; M.G.S. and K.J.S. performed research; T.M. contributed
705 analytic tools; M.G.S. analyzed data; and M.G.S. and K.J.S. wrote the paper.

706 Conflict of interest

707 The authors declare no conflict of interest.