

Analysis of the genetic diversity and structure across a wide range of germplasm reveals prominent gene flow in apple at the European level

Article

Accepted Version

Urrestarazu, J., Denance, C., Ravon, E., Guyader, A., Guinsel, R., Feugey, L., Poncet, C., Lateur, M., Houben, P., Ordidge, M. ORCID: <https://orcid.org/0000-0003-0115-5218>, Fernandez-Fernandez, F., Evans, K. M., Paprstein, F., Sedlak, J., Nybom, H., Garkava-Gustavsson, L., Miranda, C., Gassmann, J., Kellerhals, M., Suprun, I., Pikunova, A. V., Krasova, N. G., Tortutaeva, E., Dondini, L., Tartarini, S., Laurens, F. and Durel, C. E. (2016) Analysis of the genetic diversity and structure across a wide range of germplasm reveals prominent gene flow in apple at the European level. *BMC Plant Biology*, 160 (1). 130. ISSN 1471-2229 doi: <https://doi.org/10.1186/s12870-016-0818-0> Available at <https://centaur.reading.ac.uk/65747/>

It is advisable to refer to the publisher's version if you intend to cite from the work. See [Guidance on citing](#).

To link to this article DOI: <http://dx.doi.org/10.1186/s12870-016-0818-0>

Publisher: BioMed Central

All outputs in CentAUR are protected by Intellectual Property Rights law, including copyright law. Copyright and IPR is retained by the creators or other copyright holders. Terms and conditions for use of this material are defined in the [End User Agreement](#).

www.reading.ac.uk/centaur

CentAUR

Central Archive at the University of Reading

Reading's research outputs online

1 **Analysis of the genetic diversity and structure across a wide range of** 2 **germplasm reveals prominent gene flow in apple at the European level**

3
4 Jorge Urrestarazu ^{1, 2, 11 *}, Caroline Denancé ^{1 *}, Elisa Ravon ¹, Arnaud Guyader ¹, Rémi Guisnel
5 ¹, Laurence Feugey ¹, Charles Poncet ³, Marc Lateur ⁴, Patrick Houben ⁴, Matthew Ordidge ⁵,
6 Felicidad Fernandez-Fernandez ⁶, Kate M. Evans ⁷, Frantisek Paprstein ⁸, Jiri Sedlak ⁸, Hilde
7 Nybom ⁹, Larisa Garkava-Gustavsson ¹⁰, Carlos Miranda ¹¹, Jennifer Gassmann ¹², Markus
8 Kellerhals ¹², Ivan Suprun ¹³, Anna V. Pikunova ¹⁴, Nina G. Krasova ¹⁴, Elnura Torutaeva ¹⁵,
9 Luca Dondini ², Stefano Tartarini ², François Laurens ¹, Charles-Eric Durel ^{1, c}

10
11 *These authors contributed equally to this work

12 ^c Corresponding author: Charles-Eric Durel (charles-eric.durel@angers.inra.fr)

13
14 ¹ Institut de Recherche en Horticulture et Semences – UMR 1345, INRA, SFR 4207 QUASAV,
15 42 rue Georges Morel, 49071 Beaucouzé cedex, France

16 ² Department of Agricultural Sciences, University of Bologna, Viale Giuseppe Fanin 44, 40127
17 Bologna, Italy

18 ³ Plateforme Gentyane, INRA UMR1095 Genetics, Diversity and Ecophysiology of Cereals,
19 63100 Clermont-Ferrand, France

20 ⁴ CRA-W, Centre Wallon de Recherches Agronomiques, Plant Breeding & Biodiversity,
21 Bâtiment Emile Marchal, Rue de Liroux, 4 - 5030 Gembloux, Belgium

22 ⁵ University of Reading, School of Agriculture, Policy and Development, Whiteknights,
23 Reading RG6 6AR, United Kingdom

24 ⁶ NIAB EMR, East Malling Research, East Malling, Kent, ME19 6BJ, United Kingdom

25 ⁷ Washington State University Tree Fruit Research and Extension Center, 1100 N Western Ave,
26 Wenatchee WA 98801, United States

27 ⁸ RBIPH, Research and Breeding Institute of Pomology Holovousy Ltd., 508 01 Horice, Czech
28 Republic

29 ⁹ Swedish University of Agricultural Sciences, Department of Plant Breeding, Balsgård,
30 Fjälkestadvägen 459, 291 94 Kristianstad, Sweden

31 ¹⁰ Swedish University of Agricultural Sciences, Department of Plant Breeding, Box 101, 230
32 53 Alnarp, Sweden

33 ¹¹ Public University of Navarre (UPNA), Campus Arrosadia, 31006 Pamplona, Spain

34 ¹² Agroscope, Institute for Plant Production Sciences IPS, Schloss 1, P.O. Box, 8820 Wädenswil,
35 Switzerland

36 ¹³ NCRRIH&V, North Caucasian Regional Research Institute of Horticulture and Viticulture,
37 39, 40-letiya Pobedy street, Krasnodar, 350901, Russian Federation

38 ¹⁴ VNIISPK, The All Russian Research Institute of Fruit Crop Breeding, 302530, p/o Zhilina,
39 Orel district, Russian Federation

40 ¹⁵ Kyrgyz National Agrarian University, 68 Mederova Street, 720005, Bishkek, Kyrgyzstan

43 Jorge Urrestarazu ^{1,2,11} *; jorge.urrestarazu@unavarra.es
44 Caroline Denancé ¹ *; caroline.denance@angers.inra.fr
45 Elisa Ravon¹; elisa.ravon@angers.inra.fr
46 Arnaud Guyader ¹; Arnaud.Guyader@angers.inra.fr
47 Rémi Guisnel ¹; Remi.Guisnel@angers.inra.fr
48 Laurence Feugey ¹; Laurence.Feugey@angers.inra.fr
49 Charles Poncet ³; charles.poncet@clermont.inra.fr
50 Marc Lateur ⁴; lateur@cra.wallonie.be
51 Patrick Houben ⁴; p.houben@cra.wallonie.be
52 Matthew Ordidge ⁵; m.ordidge@reading.ac.uk
53 Felicidad Fernandez-Fernandez ⁶; Felicidad.Fernandez@emr.ac.uk
54 Kate M. Evans ⁷; kate_evans@wsu.edu
55 Frantisek Paprstein ⁸; fp@vsuo.cz
56 Jiri Sedlak ⁸; sedlak@vsuo.cz
57 Hilde Nybom ⁹; Hilde.Nybom@slu.se
58 Larisa Garkava-Gustavsson ¹⁰; Larisa.Gustavsson@slu.se
59 Carlos Miranda ¹¹; carlos.miranda@unavarra.es
60 Jennifer Gassmann ¹²; jennifer.gassmann@agroscope.admin.ch
61 Markus Kellerhals ¹²; markus.kellerhals@agroscope.admin.ch
62 Ivan Suprun ¹³; ivan-sn@rambler.ru
63 Anna V. Pikunova ¹⁴; pikuanna84@mail.ru
64 Nina G. Krasova ¹⁴; n.krasova@yandex.ru
65 ElnuraTorutaeva ¹⁵; elnura.torutaeva@mail.ru
66 Luca Dondini ² ; luca.dondini@unibo.it
67 Stefano Tartarini ²; stefano.tartarini@unibo.it
68 François Laurens ¹; Francois.laurens@angers.inra.fr
69 Charles-EricDurel. ^{1,c}; charles-eric.durel@angers.inra.fr
70

71 **Abstract**
72

73 **Background:** The amount and structure of genetic diversity in dessert apple germplasm
74 conserved at a European level is mostly unknown, since all diversity studies conducted in
75 Europe until now have been performed on regional or national collections. Here, we applied a
76 common set of 16 SSR markers to genotype more than 2,400 accessions across 14 collections
77 representing three broad European geographic regions (North+East, West and South) with the
78 aim to analyze the extent, distribution and structure of variation in the apple genetic resources in
79 Europe.

80

81 **Results:** A Bayesian model-based clustering approach showed that diversity was organized in
82 three groups, although these were only moderately differentiated ($F_{ST}=0.031$). A nested
83 Bayesian clustering approach allowed identification of subgroups which revealed internal
84 patterns of substructure within the groups, allowing a finer delineation of the variation into eight
85 subgroups ($F_{ST}=0.044$). The first level of stratification revealed an asymmetric division of the
86 germplasm among the three groups, and a clear association was found with the geographical
87 regions of origin of the cultivars. The substructure revealed clear partitioning of genetic groups
88 among countries, but also interesting associations between subgroups and breeding purposes of
89 recent cultivars or particular usage such as cider production. Additional parentage analyses
90 allowed us to identify both putative parents of more than 40 old and/or local cultivars giving
91 interesting insights in the pedigree of some emblematic cultivars.

92

93 **Conclusions:** The variation found at group and sub-group levels may reflect a combination of
94 historical processes of migration/selection and adaptive factors to diverse agricultural
95 environments that, together with genetic drift, have resulted in extensive genetic variation but
96 limited population structure. The European dessert apple germplasm represents an important
97 source of genetic diversity with a strong historical and patrimonial value. The present work thus
98 constitutes a decisive step in the field of conservation genetics. Moreover, the obtained data can
99 be used for defining a European apple core collection useful for further identification of
100 genomic regions associated with commercially important horticultural traits in apple through
101 genome-wide association studies.

102

103 **Keywords:** *Malus x domestica* Borkh., genetic resources, population structure, variability, SSR
104 markers, differentiation, parentage analysis.

105

106 **Background**

107 Cultivated apple (*Malus x domestica* Borkh.) is one of the most important fruit crops
108 grown in temperate zones and the most important in the *Rosaceae* family [1]. Although there
109 are more than 10,000 documented apple cultivars worldwide and the apple production area is
110 widespread geographically, the global production is dominated by relatively few cultivars, many
111 of which are closely related [2, 3]. Moreover, in the last century, despite the existence of a large
112 number of apple breeding programs worldwide, only a few well-adapted genotypes (e.g., ‘Red
113 Delicious’, ‘Golden Delicious’, ‘Jonathan’, ‘McIntosh’ or ‘Cox’s Orange Pippin’) were
114 extensively used in apple breeding to release new varieties with desirable traits [2, 4, 5]. The
115 additional release of clonal selections of the most popular and widely grown varieties has
116 further contributed towards the uniformity of commercial apple orchards [6–8]. The gradual
117 replacement of the traditional and locally well-adapted cultivars by a few wide-spread modern
118 varieties has led to a dramatic loss of genetic diversity in the orchards and may also hamper
119 future plant breeding.

120 The recognition of this situation has encouraged the establishment of action towards the
121 preservation of apple genetic resources worldwide. Multiple apple collections are presently
122 maintained in Europe, preserving mainly old cultivars which have been grown traditionally in
123 their respective regions, but also other cultivars with diverse geographic origins introduced a
124 long time ago, that represent elite selections from before the time of formal breeding. Most of
125 these existing collections were established before molecular identification became available,
126 and in the absence of marker data, the criteria used in the past for selecting the germplasm to be
127 preserved in collections focused mainly on morphology (pomology), eco-geography and/or
128 passport information [9]. The effectiveness of these conservation approaches depends upon the
129 criteria used for selecting germplasm and it has been suggested that genetic diversity may not
130 always be optimal in these, or equivalent collections in other crops [10, 11], and therefore,
131 unintended internal redundancies are expected. Assessment of the genetic diversity in fruit tree
132 species is nowadays mainly performed by marker genotyping techniques [12]. Molecular

133 markers have therefore become an indispensable tool in the management of germplasm
134 collections, and their use is widely applied in characterization to assist and complement
135 phenotypic assessments and to re-examine the composition of the collections [11, 13–16]. The
136 use of molecular markers has not only important implications with regard to the efficiency of
137 the management of the genetic resources, but constitutes a key instrument to evaluate diversity,
138 to elucidate the underlying genetic structure of the germplasm and to quantify relatedness and
139 differentiation between populations among other multiple applications [17–20]. Such
140 knowledge is of high relevance since the conservation of plant genetic resources only fulfills its
141 full potential when they are used effectively, which requires knowledge of the extent and
142 structure of the variation occurring within the material preserved [21].

143 Until now, the studies of diversity and genetic structure conducted in European apple have
144 been based on the analyses of material from limited geographic areas (mostly nation-scale) [11,
145 14, 22–26]. By contrast, the extent and structure of the apple genetic diversity conserved at a
146 European level have remained largely unknown. The main obstacle is the different sets of SSR
147 markers used in the different European collections preventing an overall comparison [27]. Thus,
148 in the frame of the EU-FruitBreedomics project [28] a single set of 16 SSR markers was used in
149 a very broad set of apple germplasm (~2440 accessions, mostly of dessert use) preserved in
150 collections located in eleven countries and representing three broad European geographical
151 regions (North+East, West and South) in order to determine the diversity in apple collections at
152 a European scale, to evaluate gene flow in cultivated apple across Europe, as well as to elucidate
153 the stratification of germplasm into population subdivisions and finally, to perform parentage
154 analysis. This is the largest study of apple genetic resources at the pan-European level.

155 **Results**

156 **SSR polymorphism – identification and redundancy**

157 Among the 2,446 accessions, ten accessions did not show clear PCR amplifications and
158 were discarded from the analysis. Pairwise comparison of multilocus profiles revealed 219
159 groups of redundancies (Additional file 1), leading to the removal of 405 redundant accessions

160 before further analyses (16% of redundancy). The number of accessions in each of these
161 identical SSR profile groups varied from two to nine. The cumulative probability of identity
162 (P_{ID}) was extremely low: $P_{ID} = 1.3 \times 10^{-22}$, thus highlighting the low risk of erroneous attribution
163 of accessions to duplicate groups. Redundancies were found both within and between
164 collections, leading to the confirmation of numerous previously documented synonyms (e.g.,
165 ‘Papirovka’ and ‘White Transparent’, ‘London Pippin’ and ‘Calville du Roi’, or ‘Président van
166 Dievoet’ and ‘Cabarette’) and allowing the putative identification of numerous unknown
167 synonyms or mutant groups (e.g., ‘Gloria Mundi’ = ‘Mela Zamboni’ = ‘Audiena de Oroz’ =
168 ‘Belle Louronnaise’, ‘Court-Pendu Plat/Doux/Gris’ = ‘Krátkostopka královská’, ‘Reinette de
169 Champagne’ = ‘Maestro Sagarra’ or ‘Reinette Simirenko’ = ‘Renetta Walder’ = ‘Burdinche’).
170 Redundancy groups also supported the notion of several national/local name translations such as
171 the English cultivar ‘Cornish Gilliflower’ translated into ‘Cornwallské hřebíčkové’ (i.e.,
172 ‘Cornish clove’), or ‘White Transparent’ and ‘Skleněné žluté’ (i.e., ‘yellow glass’) in Czech and
173 ‘Transparente Blanca’ in Spanish, the Russian cultivar ‘Korichnoe polosatoe’ translated into
174 ‘Kaneläpple’ in Swedish (i.e., ‘cinnamon apple’), or the cultivar ‘La Paix’ translated into
175 ‘Matčino’ (i.e., ‘Mother’, a synonym of ‘La Paix’) in Czech. Several cases of homonymy (i.e.,
176 accessions with the same name but different SSR profiles) were also found, e.g., three different
177 SSR profiles for the same accession names ‘Pomme Citron’ or ‘Charles Ross’. Data allowed
178 identifying some obvious labeling errors, e.g., X2698 ‘Court Pendu Plat’ which was shown to
179 be the rootstock ‘MM106’, or CRAW-0362 ‘Transparente de Croncels’ which was found likely
180 to actually be ‘Filippa’ (Additional file 1). Following these observations, the apple germplasm
181 dataset was reduced to 2,031 unique genotypes (i.e., exhibiting distinct SSR profiles). Among
182 these individuals, 162 (8% of the different genotypes) were removed since they had a putative
183 triploid profile, while another ten were discarded because of too much missing SSR data, or
184 because further identified as rootstock or outliers in a preliminary Principal Coordinate
185 Analysis. The final number of unique diploid genotypes further analyzed was therefore 1,859.
186 Using passport data and other accessible information, it was possible to attribute geographical
187 regions of origin (either for three broad designated European regions or, when possible, specific

188 countries) for a large part of the unique genotypes. Roughly 89% (1,653) of these genotypes
189 could be geographically assigned, with 261, 1,074 and 318 genotypes assigned to
190 Northern+Eastern, Western and Southern historical regions of origin, respectively (Additional
191 file 1). In brief, the Northern+Eastern region was composed of germplasm originating in Nordic
192 European countries plus Russia, the Western region was composed of germplasm originating in
193 Western and Central European countries and the Southern region was composed of germplasm
194 from Spain and Italy (see Methods for more details). The remaining 11% consisted of either
195 genotypes lacking passport information or genotypes with contradictory information in passport
196 data from different origins. Similarly, the specific country of origin could be attributed to 1,550
197 genotypes out of the 1,653 geographically assigned (Additional file 1). It is important to note
198 that the European region or country of origin assigned to a genotype was independent from the
199 location of the collection where the sampled accession was maintained, since many collections
200 contained accessions from various origins.

201

202 **Genetic diversity across and within European regional groups**

203 The 16 SSR markers amplified a total of 369 alleles across the 1,859 apple accessions used
204 for diversity analysis, ranging from 17 (CH02c09 and CH05f06) to 35 (CH02c06) alleles per
205 locus. The average number of alleles per locus was 23.06, whereas the mean effective number
206 of alleles per locus was 6.59 (Table 1). High average number of alleles per locus and almost
207 identical mean effective number of alleles per locus were noted for the three geographical
208 regions of origin of the germplasm. Allelic richness was normalized to the smallest group (i.e.,
209 North+East) to avoid a group size-dependent bias of results. Overall, the results obtained for the
210 material of the three designated regions of origin suggested the existence of a high and
211 relatively homogeneous allelic diversity across Europe (Table 1). Within the 369 alleles
212 identified in the overall set (i.e., across Europe), 73.4% and 52.0% were found at frequencies
213 below 5% and 1%, respectively (Table 1; data not shown for 1%). A similar proportion of rare

214 alleles was obtained for the material from the three designated geographical regions of origin,
215 with the exception of alleles detected at a frequency $< 1\%$ with Northern+Eastern and Southern
216 European origins, for which slightly lower percentages were identified ($\approx 38\%$). Almost
217 identical mean H_e values were obtained for the overall dataset (0.83) and for the germplasm
218 from each of the three geographical groups (Table 1). Cross-comparison of the allelic
219 composition for the accessions classified into geographic categories showed that 221 out of the
220 362 alleles (seven alleles appeared only in accessions that could not be classified into
221 geographic groups) were detected in all three geographical groups, 59 alleles (16.3%) were
222 identified in two geographic groups only, whereas 82 alleles (22.6%) were specifically found
223 only in one geographic group (i.e., private alleles). At the national level (i.e., countries of origin
224 of the unique genotypes), some countries exhibited a higher rate of private alleles than others:
225 especially, genotypes assigned to Switzerland, Italy and Russia harboured 15, 14 and 14 private
226 alleles (respectively), genotypes from Spain and France harboured 7 private alleles each,
227 whereas genotypes from the Netherlands, Belgium, Great Britain or Sweden had a maximum of
228 one private allele. The pattern of distribution of the frequent alleles (frequency > 0.05) between
229 Southern, Northern+Eastern and Western germplasm was analyzed for each locus separately
230 using Chi² tests. Highly significant differences in the allelic distributions ($P < 0.001$) were found
231 between all the geographic groups for all markers except for the CH-Vf1 locus when comparing
232 Southern and Western germplasm (data not shown).

233

234 **Genetic structure and differentiation**

235 A Bayesian model-based clustering method was applied to the 1,859 unique diploid
236 genotypes in order to elucidate the underlying genetic structure at a European scale. The
237 analysis of Evanno's ΔK statistic indicated unambiguously $K=3$ as the most likely level of
238 population stratification (Fig. 1 a1). The mean proportion of ancestry of the genotypes to the
239 inferred groups was 0.81. Using the threshold of $qI \geq 0.80$ to define strong assignments to
240 groups, 1,175 genotypes (63%) were identified as strongly associated to a group. This

241 partitioning level corresponded to an asymmetric division of the material into three groups: K1
242 composed of 506 genotypes, K2 containing 401 genotypes, and K3, the largest group,
243 comprising 952 genotypes. Diversity estimates revealed high levels of allelic variation within
244 each group, with allelic richness ranging between 16.0 (K3) and 18.6 (K1) (Table 2). Genetic
245 discrimination between the three groups was confirmed through a multivariate Principal
246 Coordinate Analysis (PCoA) (Fig. 2). In the bi-dimensional plot, K1 was located mostly to the
247 left of the Y axis, and K2 mostly below the X axis, while K3 occurred to the right of the Y axis
248 and mostly above the X axis. A Neighbor-joining tree also showed three different main clusters
249 (Fig. 3), supporting the identification of the three groups by the Bayesian method.

250 The genetic differentiation between the three designated geographic regions of origin was
251 low ($F_{ST} = 0.021$, $P < 0.001$, Table 3), suggesting a weak genetic structure for this crop at a
252 European scale in terms of geographical origin. The level of genetic differentiation between the
253 three groups inferred by Structure was only slightly higher ($F_{ST} = 0.031$, $P < 0.001$). The largest
254 differentiation between pairs of groups was found between Northern+Eastern and Southern
255 germplasm ($F_{ST} = 0.042$, $P < 0.001$), whereas much lower F_{ST} values were found between the
256 Western and each of the Northern+Eastern ($F_{ST} = 0.023$, $P < 0.001$) and Southern ($F_{ST} = 0.015$,
257 $P < 0.001$) materials.

258

259 The relationship between membership of accessions within the three groups defined by
260 Structure and their geographical regions of origin was also analyzed. 80% and 75% of the
261 accessions from Northern+Eastern and Southern Europe clustered in K2 and K1 respectively.
262 The relationship between the material with Western European origin and the third group (K3)
263 was less evident (63%), but still visible by comparison (Fig. 1b). Although the genetic
264 differentiation revealed between the three groups defined by Structure was not very high, the
265 existence of a relationship between the grouping by geographical regions of origin of the
266 accessions and the three inferred groups is noteworthy. Furthermore, when considering the
267 specific country of origin attributed to the cultivars, the distribution within the three Structure-

268 defined groups appears to follow a clear gradient from North(East) to South of Europe (Fig. 4);
269 the cultivars from Northern Europe and Russia were mainly assigned to the K2 group and the
270 Spanish and Italian cultivars were mainly assigned to the K1 group, with intermediate patterns
271 found for those countries located at the interfaces of the broad regions.

272

273 Within the admixed accessions (i.e., $qI < 0.8$) for which the geographical regions of origin
274 (Northern+Eastern, Southern and Western) was known, we defined a membership coefficient
275 threshold ($qI < 0.55$) with the aim of identifying genotypes unambiguously in *admixis*, in order to
276 examine whether a supplemental relationship could be found between geographical region and
277 grouping by Structure for the admixed material. For the unambiguously admixed material (i.e.,
278 $qI < 0.55$) of Southern European origin, the average proportion of ancestry (qI) was 0.45 to K1
279 (the group mostly associated with material from Southern Europe), followed by 0.42 to K3 and
280 0.13 to K2, the groups mostly composed by material from Western and Northern+Eastern
281 European origins, respectively (data not shown); a slightly less pronounced, but complementary,
282 pattern was observed for the unambiguously admixed germplasm (i.e., $qI < 0.55$) of
283 Northern+Eastern Europe with average proportions of ancestry of 0.43, 0.35 and 0.22 to K2, K3
284 and K1, respectively. For the unambiguously admixed material (i.e., $qI < 0.55$) of Western origin
285 the average proportion of ancestry to each of these three groups was almost identical
286 (approximately 1/3). This result was in line with the lower F_{ST} values found between the groups
287 K1 / K3 ($F_{ST} = 0.024$, $P < 0.001$) in comparison with the slightly higher differentiation between
288 the groups K1 / K2 ($F_{ST} = 0.039$, $P < 0.001$) and K2 / K3 ($F_{ST} = 0.036$, $P < 0.001$). The dispersion
289 of the three groups in the PCoA plot was also in agreement with these results, showing the
290 highest overlap between K1 and K3 followed by K2 and K3.

291

292 **Nested-Bayesian clustering approach: substructuring of the diversity**

293 In order to investigate the substructuring of the diversity within each of the three groups
294 identified in the initial analysis we used a nested application of the Structure software. To do
295 this, the three groups were analyzed independently. To evaluate the strength of the hypothetical

296 subdivisions (i.e., subgroups) within each group, simulations for each K value were examined,
297 paying attention to the internal consistency between the runs, the mean proportion of ancestry of
298 accessions within each subgroup, and the proportion of accessions unequivocally assigned ($qI \geq$
299 0.80).

300 The analysis of the relationships between K and ΔK for K1 suggested a probable
301 subdivision of this material into three subgroups and the assignment of genotypes was well
302 correlated between runs. The average proportion of ancestry for the accessions clustered in the
303 three subgroups of K1 was 0.75, with 44% of the accessions showing strong assignments. Two
304 subgroups for K2 and three for K3 were similarly established. In both cases, the assignment of
305 genotypes was well correlated between runs, and almost identical average proportions of
306 ancestry to those for the subgroups of K1 were obtained with slightly higher proportions of
307 strongly assigned accessions (47% and 50% respectively). Secondary peaks at other K values
308 were also explored but these subdivisions had less statistical support (data not shown).
309 Therefore, we adopted eight subgroups as the most suitable partitioning degree of substructuring
310 (Fig. 1 a2). For these eight subgroups the affinity of almost half of the individuals (47%) to their
311 respective subgroups was strong and the assignment of *admixed* accessions was consistent
312 between runs. The examination of the eight subgroups showed considerable differences in size,
313 ranging from 148 (K1.3) to 415 (K3.3) genotypes, and variable proportion of accessions
314 strongly assigned to the inferred subgroups (Table 2). K3.2 was the subgroup with the highest
315 proportion of strongly assigned genotypes (57%), whereas K1.3 had the highest proportion of
316 admixed accessions. The proportion of accessions unambiguously assigned for the remaining
317 six subgroups ranged from 41% to 54%, whereas the mean proportion of ancestry for the
318 accessions clustered in each one of the eight subgroups was very stable (≈ 0.75).

319 The analysis of the relationship between the different subgroups and the putative countries
320 of origin of the germplasm indicated potentially interesting correlations, especially for groups
321 K1 and K3. About 70% of the subgroup K1.2 consisted of germplasm originating from Spain.
322 Similarly, 46% of the subgroup K1.1 and 50% of the subgroup K1.3 consisted of germplasm

323 originating from Switzerland and Italy, respectively (Additional file 2); the latter subgroup was
324 also composed of a further 39% of the cultivars with a French origin and interestingly, a
325 significant proportion of these were attributed to Southeastern France (data not shown). The
326 disentangling of the substructuring pattern therefore allowed not only the dissection of the
327 internal distribution of the diversity within group K1, but also the detection of three subgroups
328 strongly associated with some particular countries of origin. With respect to the collections from
329 the Northern+Eastern part of Europe (Sweden, Finland and Russia), no clear differentiation of
330 the germplasm in the two subgroups of K2 was observed (Additional file 2). For the subgroup
331 K3.1, about half of the germplasm consisted of cultivars from either the United Kingdom or
332 France. All of the 40 cultivars selected in the French collection as being recently bred, clustered
333 in a single small subgroup (K3.2) which was mostly composed of English, US and, perhaps
334 more surprisingly, Spanish cultivars. Major standard cultivars such as ‘Golden Delicious’, ‘Red
335 Delicious’, ‘Jonathan’ and ‘Ingrid Marie’ were also assigned to this subgroup, as well as ‘Cox’s
336 Orange Pippin’ and ‘James Grieve’. Interestingly, most of the 40 cider apple cultivars (87%)
337 were assigned to one subgroup (K3.3) which was mostly composed of French, English, and
338 Swiss cultivars. The other standard cultivars were assigned to the latter subgroup and to
339 subgroup K3.1.

340

341 Genetic diversity estimates were calculated for all the subgroups obtained by the nested
342 Bayesian model-based clustering (Table 2). While H_e ranged from 0.76 (K1.3 and K3.2) to 0.84
343 (K1.1), indicating a high level of heterozygosity contained in all the subgroups, the percentage
344 of alleles represented in each one of the eight subgroups was very variable, ranging from 46%
345 (K3.2) to 76% (K3.1). Some private alleles were identified in all subgroups except for K3.2.
346 They were most abundant in K1.1, but a considerable number of them were found also in K3.3
347 and K2.1. Most of the private alleles (approx. 72%) were also unique as they were identified in
348 only one accession. To properly evaluate the allelic diversity between the eight subgroups, we
349 applied a rarefaction approach to compensate for the differences in subgroup size. The allelic

350 richness obtained for the eight subgroups supported the previous results, confirming the highest
351 diversity in K1.1 and the lowest diversity in K3.2.

352 Estimates of genetic differentiation showed that only 3.7% (K1) and 3.4% (K2) accounted
353 for variation among subgroups within groups (Table 3). The genetic differentiation between the
354 subgroups into which K3 was subdivided was considerably lower (Table 3). Considering the
355 eight subgroups obtained by the overall Nested Bayesian approach, the results showed that
356 variation among subgroups accounted for 4.4% of the total variation. Regarding the F_{ST} pairwise
357 tests between subgroups (Table 4), irrespective of whether they belonged to the same group or
358 not, the highest F_{ST} corresponded to the pair K1.3 / K2.2 ($F_{ST}=0.087$, $P<0.001$), followed by
359 K1.3 / K2.1 ($F_{ST}=0.077$, $P<0.001$), and the lowest to the pairs K1.1 / K3.3 ($F_{ST}=0.016$, $P<0.001$)
360 and K3.1 / K3.3 ($F_{ST}=0.023$, $P<0.001$).

361 **Parentage reconstruction**

362 Two-parents-offspring relationships within the 1,859 diploid genotypes were explored
363 using CERVUS software. A total of 46 putative trios (offspring and two inferred parents) were
364 identified with high (95%) confidence level. These consisted of two already documented trios,
365 ('Calville Rouge du Mont Dore' and 'Belle de Mleiev' and their parents; [23]), as well as
366 another 10 recent and 34 old cultivars (Table 5). The two parents of the 10 modern cultivars, for
367 which full parentage was already documented were correctly inferred (e.g., 'Heta', 'Jaspi' and
368 'Pirkko' = 'Lobo' x 'Huvitus', 'Pirja' = 'Huvitus' x 'Melba', or 'Mio' = 'Worcester Pearmain'
369 x 'Oranie'). In most cases, the two parents of the older cultivars were not known and thus newly
370 inferred (Table 5). Inferred parentage was found for old cultivars from various European
371 countries (6x for Italy; 4x for Great Britain, Switzerland, Czech Republic, and Sweden; 3x for
372 Germany and Spain; 2x for Belgium). Perhaps unsurprisingly, some accessions were more
373 frequently inferred as parents, such as the two French cultivars 'Reine des Reinettes' (= 'King
374 of the Pippins') or 'Transparente de Croncels' which were each identified three times.
375 Geographic convergence of parentage was frequently observed (e.g., 'Kramforsäpple' =
376 'Sävstaholm' x 'Åkerö', all three from Sweden; 'Beauty of Moray' = 'Keswick Codlin' x

377 'Stirling Castle', all three from Great Britain; 'Roja de Guipuzcoa' = 'Urte Sagarra' x 'Maxel
378 Gorri', all three from Spain; or 'Scodellino' = 'Abbondanza' x 'Decio', all three from Italy).
379 But hybridizations between cultivars from distant countries were also observed (e.g., 'Rotwiler'
380 presumably from Switzerland = 'King of the Pippins' x 'Alexander' from France and Ukraine,
381 respectively; or 'Godelieve Hegmans' from Belgium = 'Red Astrakan' x 'Transparente de
382 Cronicels' from Russia and France, respectively). It should be noted that the female and male
383 status of the inferred parents could not be specified from the available SSR markers.

384 **Discussion**

385

386 **Identification and redundancy**

387 The exchange of genotyping data between research units has increased considerably in
388 recent years, with the aim to investigate the extent and distribution of diversity for specific
389 crops at a wide geographic scale. In this study, the application of a common set of 16 SSR
390 markers on a wide set of dessert apple cultivars distributed across three broad European regions
391 allowed the detection of redundant accessions and duplicated genotypes between and within
392 collections, and the description of the structuration of a significant part of the European apple
393 diversity. Cross-comparison of SSR data in attempts to combine datasets from multiple sources
394 has often been problematic due to challenges in harmonizing the allelic sizes between different
395 laboratories [18, 29, 30]. By combining existing data over numerous shared reference
396 accessions in our collections with the re-genotyping of a subset of the accessions, we were able
397 to strongly secure the SSR allele adjustment over sites. This dataset represents a highly valuable
398 resource for the comparison of apple germplasm collections throughout Europe and the rest of
399 the world. Taking into consideration the rich allelic diversity present in the European apple
400 germplasm, it would be useful to identify a relatively small set of varieties that offer a good
401 representation of the allelic variability identified in this germplasm to act as an internal control
402 (i.e., a reference set) between laboratories for future use.

403 Interestingly, duplicate groups involving accessions from different collections underlined
404 some putative drift in the cultivar denomination. Some good examples were 'Pott's seedling'

405 and ‘Pottovo’ (FBUNQ14), or ‘Signe Tillish’ and ‘Signatillis’ (FBUNQ34). In addition, ‘sports’
406 are often given derivative names (e.g., ‘Crimson Peasgood’ as a sport of ‘Peasgood’s Nonsuch’)
407 but the current analysis was not set up to distinguish between clones and ‘sports’ of cultivars
408 with potential morphological differences. Many likely errors in denomination of genotypes were
409 also detected when multiple representatives of a given cultivar were detected within a group, but
410 a single supposed representative was obviously outside of the group and was often associated
411 with representatives of a different cultivar. For example, ‘Drap d’Or’ and ‘Chailleux’
412 (FBUNQ92) are known to be synonyms used in France for the same cultivar, and accession
413 DCA_D35 ‘Drap Dore’, which was found to belong to the group FBUNQ50, was most likely a
414 denomination error since almost all other members of this group were ‘Winter Banana’. In other
415 cases, accessions with uncertain denomination could be resolved, such as CRAW-1858
416 ‘Reinette Baumann?’ (FBUNQ21) and accession CRAW-1108 ‘Peasgood Nonsuch?’
417 (FBUNQ51) for which the molecular analyses confirmed that they were most likely ‘true-to-
418 type’ cultivars. The question of ‘trueness-to-type’ is a major issue in apple germplasm
419 management where extensive budwood exchange between regions and countries has occurred
420 for centuries. Indeed, an erroneously denominated accession can be transmitted from collection
421 to collection for years, such that a large number of representatives within a duplicate group (as
422 per the present study) should not always be considered definitive proof of the trueness-to-type
423 of accessions but this objective evidence is extremely valuable in highlighting issues to resolve.
424 Since genebank curators have often collected material of old cultivars from private gardens or
425 from tree pasture orchards, unidentified or misidentified material can later be detected either by
426 classical phenotypic characters and/or by using genetic markers. As an example, this study
427 showed that an old so called local cultivar ‘Madame Colard’ (CRAW-0365 – FBUNQ72),
428 described to have been raised in 1910 by the nurseryman Joseph Colart at Bastogne (Belgium),
429 exhibited the same SSR profile as the old English apple cultivar called ‘Royal Jubilee’ (UK-
430 NFC 2000085) raised already in 1888. Further comparison with historical descriptions could
431 conclude that they are the same cultivar. Additional insights from the passport data of
432 accessions would be needed to help in tracing the transmission of the material from collection to

433 collection and pomological characterization will be required to compare accessions to published
434 descriptions of the variety. This will remain a task for the curators of collections, in order to
435 improve curation of germplasm in a coordinated way.

436 It is important to note that the criteria used to select the accessions at the country-level
437 were not always the same. For instance, the INRA and UNIBO material corresponded to former
438 “core collections” built to encompass a large variability not restricted to the national/local
439 accessions [23, 24]. Conversely, the UK-NFC and FRUCTUS material was restricted to older
440 diploid accessions considered to derive from UK and Switzerland, respectively. A similar,
441 despite less stringent situation was applied also for CRA-W, RBIPH, SLU, and the Spanish
442 accessions (UPNA, UDL and EEAD). For MTT, NCRRI, VNIISPK, and KNAU, the national
443 representativeness was more limited and strictly restricted to accessions considered to be
444 emblematic landrace cultivars. The germplasm was thus somewhat heterogeneous in nature, but
445 still allowed a broad examination of the European dessert apple diversity. In the future, it will
446 be useful to enlarge the dataset to include additional accessions from the collections considered
447 here as well as other European collections [11, 31] or collections from other regions worldwide
448 [32–34] to provide a wider perspective on genetic resource conservation of apple worldwide.

449 **Genetic diversity**

450 The high level of diversity and heterozygosity in apple germplasm at a European level
451 agreed with previous results obtained at collection-scale in several European countries, e.g.,
452 Italy [24], Spain [26], France [23], Sweden [22], Czech Republic [25] or Switzerland [14]. The
453 large diversity found is consistent with the weak bottleneck effect reported in connection with
454 the domestication of this species [35–37]. Probably a combination of factors are involved: i)
455 vegetative propagation methods that have been adopted since ancient times favoring the
456 dispersal of cultivars across geographic regions [38, 39], ii) forced allogamy due to the self-
457 incompatibility system of *Malus × domestica* [40], iii) multiple hybridization events at each
458 geographical region combined with human activities, e.g., selection and breeding [36, 37] and,
459 iv) diversifying selection associated with adaptive criteria for the subsistence in diverse

460 agricultural environments [41, 42]. Interestingly, the distribution of private SSR alleles over the
461 countries of origin of the unique genotypes was somewhat unbalanced at the European level
462 with much higher occurrences in genotypes assigned to Switzerland, Italy or Russia than in
463 genotypes originating from Northern-Western Europe. Whilst these findings should be
464 considered with caution because of possible biases linked to the initial sampling or to the size
465 differences of the genotype sets, this study underlines that accessions originating from Southern
466 Europe and Russia could be expected to bring original genetic diversity into modern breeding
467 programs especially for traits related to more extreme climate adaptation. Overall, the highly
468 diverse germplasm studied here contains much more genetic variation than do modern apple
469 cultivars, many of which having been selected for optimal performance within a narrow range
470 of environmental conditions [5, 37, 42].

471 **Coordinated actions: a key point for better knowledge of the resources conserved**

472 This large-scale analysis in apple germplasm constitutes a good example of the efficiency
473 and value of coordinated international actions to enhance the knowledge of diversity conserved
474 at a European level. The results obtained offer a valuable step to undertake actions to coordinate
475 European resources towards optimizing the management of apple germplasm across Europe in
476 line with the aspirations of the European Cooperative program on Plant Genetic Resources
477 (ECPGR). The results also offer a potential starting point that may open new opportunities for
478 apple breeding in the near future. All breeding advances are built upon the diversity available,
479 and a key role of the germplasm collections is to help safeguard natural forms of genetic
480 variation and to make them accessible to plant biologists, breeders, and other key users [15].
481 The extensive germplasm evaluated in this study consisted mainly of old and/or locally grown
482 accessions across Europe, many of which remain underutilized in cultivation or breeding
483 programs. The preservation of traditional cultivars in living germplasm collections must be
484 regarded as an invaluable reservoir of insufficiently explored genetic diversity that may become
485 useful for apple breeding in a near future, and the establishment of coordinated genetic data is
486 hoped to increase the accessibility of this material to breeding programs. From the perspective

487 of modern-day fruit production, most of these old varieties would now be considered as obsolete
488 since they are not particularly well-adapted to current agricultural practices and marketing.
489 Nevertheless, this material should be considered as a reservoir of potentially interesting genes to
490 be used for further improvement. This is particularly relevant in a crop like apple, for which the
491 current production is highly dependent on a very limited number of cultivars with a narrow
492 genetic basis for the bulk of current production [5]. As an example, it can be mentioned that
493 50% of the commercially marketed apple production in the European Union consists of only
494 four cultivars, ‘Golden Delicious’, ‘Gala’, ‘Idared’ and ‘Red Delicious’ [43]. The low diversity
495 of the subset of elite cultivars used for commercial production during recent decades is likely to
496 result in a bottleneck hampering future genetic improvement [37]. The recognition of this
497 situation should encourage the establishment of coordinated actions across different levels
498 (regional, national and international scales) to define strategies for the efficient conservation of
499 the genetic resources of this species.

500 **Genetic structure: major divisions and substructuring of the diversity**

501 The attribution of country of origin to traditional cultivars can be a matter of endless
502 debate, especially for those dating back two-three centuries or more. Initial descriptions in
503 pomologies and booklets can be subject to errors in denomination confused by historical
504 distribution and renaming, resulting in synonymy, as well as the re-use of old names for more
505 recent findings or misidentifications. This is less problematic for the better known old cultivars
506 as many of them have been widely documented and monitored over years in several countries.
507 However, for local cultivars and/or landraces where less information is available, the correct
508 attribution can be complicated, especially between neighboring countries. It is also worthy of
509 note that the ‘country of origin’ relies on a political construct, which can be prone to significant
510 change within the potential lifetime of many varieties of apple (and other long lived perennial
511 crops). Therefore, we first used a conservative approach and discussed our findings in terms of
512 three broad European regions of origin. Then, we analyzed the structuration at a country-scale,
513 but noting that the exact attribution of a given country to a genotype was not always

514 unanimously agreed so that this finer level of analysis should be considered with an element of
515 caution.

516 Using a Bayesian model-based clustering method we were able to initially discern the
517 existence of three robust groups reflecting major divisions of the germplasm. These groups were
518 linked with the three geographical regions of origin, although differentiated only to a low
519 degree. This would reflect a situation whereby the cultivars from a given region were more
520 frequently derived from crosses between parental cultivars from the same region than from
521 cultivars from elsewhere. Nevertheless, the migration of the plant material associated to human
522 movement together with hundreds of years of empirical selection may have caused a significant
523 gene flow across Europe. This is clearly indicated by the low genetic differentiation between
524 groups and has shaped the overall pattern of genetic diversity. A spatially and temporally
525 dynamic process where seeds and mainly graftwood were exchanged between geographically
526 distinct populations has contributed to the increase of the genetic diversity in each area through
527 unintentional gene flow or human-mediated intentional crosses [35, 36, 44]. The background
528 common to other long lived tree fruits, including factors such as multiple origins of cultivated
529 populations, ongoing crop-wild species gene flow and clonal distribution of genotypes together
530 with the features associated with fruit tree species (lengthy juvenile phase, extensive
531 outcrossing, widespread hybridization or mechanisms to avoid selfing) has defined the way they
532 evolve in nature and resulted in extensive population genetic variation, but limited population
533 structure [44]. A possible cause of divergence between the three identified groups could be the
534 differential adaptation to distinct environmental conditions as are the case between Southern,
535 Western and Northern+Eastern Europe. A similar situation was postulated for grapevine
536 cultivars where the genetic structure appeared to be strongly shaped by geographic origin and
537 intentional selection [13]. But since selection causes differentiation in particular regions of the
538 genome on which selection pressure is acting [45], another likely cause of the population
539 structure is genetic drift (i.e., changes in allelic frequencies caused by chance events) as also

540 shown in e.g., apricot [46]. Together with selection, migration and drift can shape the local
541 adaptation of species [47].

542 Although there may have been some mistakes in attributing cultivars to country of origin,
543 the genetic makeup of the cultivars at the European level clearly appeared to show a North-East
544 to South gradient. Interestingly, some countries exhibit intermediate marker data patterns in
545 consistency with their intermediate geographic positions. This was clearly manifested at the
546 national scale for the German and Czech cultivars which were shared between K2 and K3
547 groups. Similarly, the French and Swiss cultivars were shared between K1 and K3 groups. By
548 contrast, cultivars from Southern Europe (Spain and Italy), from Northwestern Europe (United
549 Kingdom and Ireland, Belgium, the Netherlands), and from North+Eastern Europe (Sweden and
550 Finland) and Russia were mostly assigned to a single group (K1, K3, and K2, respectively). For
551 the admixed germplasm from Southern and Northern+Eastern European geographical regions of
552 origin, a certain degree of introgression with the Western germplasm was also indicated in
553 contrast to the low contribution of the Northern+Eastern germplasm into the Southern
554 germplasm and viceversa. Thus, in agreement with the correspondence between clustering and
555 regions or countries of origin of the germplasm, the geographical proximity appears to align
556 with the patterning observed in the admixed accessions.

557 In cases demonstrating the presence of a significant hierarchical population structure as this
558 study suggests, this method preferentially detects the uppermost level of structure [26, 48–50].
559 As a consequence, when large datasets in species with a complex background are analyzed, it is
560 possible for an underlying substructure to remain undetected within the major divisions of the
561 germplasm. In this context, the “*nested (or two-steps) Structure*” clustering method has been
562 shown to be an efficient tool to delineate further levels of substructure in both apple and other
563 plant species [10, 24, 26, 49–52]. In this study, the three groups inferred from the first round of
564 Structure analysis were used as the starting point for revealing internal substructuring. Eight
565 subgroups were identified with remarkable differences in both allelic composition and richness,
566 as well as a considerable number of private alleles associated to particular subgroups.

567 Nevertheless, the relationship between the placement of the genotypes in the subgroups and
568 their country of origin varied considerably between subgroups in contraposition to the clearer
569 and more consistent clustering trend within the three groups. As discussed earlier, this
570 stratification may reflect historical processes of selection and adaptation to local conditions that
571 might suggest a “*fine-delineation*” of the intra-variation within each main geographical region
572 of origin. This is most probably the case for the K1.1 subgroup which mainly consists of
573 Spanish cultivars and could reflect a process of both local adaption and isolation by distance
574 related to the Pyrenean barrier. For the K1.3 subgroup, local adaptation to the Southern region
575 could be inferred together with a potential for more intense commercial exchange between Italy
576 and Southern France. For other subgroups, the relationship with particular countries or small
577 regions was not obvious, but some interesting associations between subgroups of group K3 and
578 recent cultivars and some of their founders or particular usage (cider apple cultivars) could be
579 noticed.

580 **Relatedness and family relationships**

581 The previously reported parentage of 10 recent cultivars was correctly inferred in all cases.
582 These results served as a control and validated the parentage assignment obtained with the
583 CERVUS software [53] indicating that the number and informativeness of SSR markers were
584 sufficient at least for these cultivars. The 16 SSR markers were nevertheless limited in their
585 ability to infer parentages, and additional cases might have been detected with a larger number
586 of SSR markers. In a recent paper [54], it was suggested that the number of 27 SSR loci used in
587 that study was a minimum to be utilized for full parentage reconstruction. Basically, the LOD
588 score tests used in the CERVUS software are computed according to the SSR allelic
589 frequencies, and thus, parentages involving common alleles are more difficult to detect. By
590 contrast, parentages involving low frequency and rare alleles are more easily detected. On that
591 basis, it is worthy to note that the more frequently detected parents (i.e., ‘Reine des Reinettes’ =
592 ‘King of the Pippins’, and ‘Transparente de Croncels’) are possibly representing a biased view
593 of the frequently involved parents, as they most probably carry rare or low frequency alleles in

594 at least some SSR loci. Putative parents present in the dataset but carrying more common alleles
595 may have been hidden because of the statistical limits of their detection with 16 SSR markers. A
596 similar situation was observed by [23] with the frequent appearance of ‘Reine des Reinettes’ as
597 a parent of four old cultivars out of 28, using 21 SSR markers. In the near future, medium and
598 high density SNP arrays [55–57] will provide much more power to infer parentages.

599 The parentage of some old cultivars was either confirmed, in the case of ‘Ernst Bosch’ =
600 ‘Ananas Reinette’ x ‘Mank's Codlin’ (synonym: ‘Evino’) or augmented, in the case of ‘Ben’s
601 Red’ = ‘Devonshire Quarrenden’ x ‘Box Apple’ (Table 5) where the second parent was initially
602 hypothesized to be ‘Farleigh Pippin’ [58]. Distances between the geographic origins of the
603 inferred parents (when known), ranged from crosses between geographically close cultivars to
604 crosses between very distant cultivars, reflecting the large gene flow across Europe caused by,
605 e.g., extensive exchange of budwood over centuries.

606 Some traditional folklore about the origination of old apple cultivars could be either
607 substantiated or refuted by the SSR-based parentage information. As one example, the old
608 Swedish cultivar ‘Förlovningsäpple’ is said to derive from a locally acquired seed in Northern
609 Sweden where only a few cold-hardy apples can be grown. The two unknown parents were here
610 inferred to be the Swiss cultivar ‘Heuapfel’ and the wide-spread cultivar ‘Saint Germain’
611 (X1646) also known as ‘Vitgylling’ in Sweden, a name used for a group of more or less similar,
612 white-fruited, early-ripening and winter-hardy cultivars. Interestingly, the ‘Vitgylling’ accession
613 included in the present study (BAL072) did not have the same SSR profile as ‘Saint Germain’,
614 but they share one allele for all 16 SSR loci and may therefore be related. In two other cases,
615 traditional Swedish folklore indicated that a sailor brought an exotic seed to the island of
616 Gotland and to Kramfors in Northern Sweden, respectively, resulting in ‘Stenkyrke’ and
617 ‘Kramforsäpple’. For ‘Stenkyrke’, one parent is the Swedish ‘Fullerö’ and the second is the
618 German cultivar ‘Danziger Kantapfel’ which has been much grown in Sweden. The origin of
619 ‘Stenkyrke’ is thus probably much more local than anticipated. Similarly, the surmised

620 American sailor origin of the seed giving rise to ‘Kramforsäpple’ is refuted by the fact that the
621 parents of this cultivar are the Swedish ‘Sävstaholm’ and ‘Åkerö’.

622 It is important to keep in mind that trueness-to-type of the accessions is not guaranteed,
623 thus the labeling of the offspring or the parents can be erroneous in some cases. Conversely, the
624 inferred parentages are robustly established so that the genetic relationships between the
625 accessions are valid independently of their names. Crosses between the two inferred parents
626 could be performed to reproduce the cross which gave birth to the offspring cultivar, especially
627 if genetic analysis of some particular traits of the latter genotype indicates an interesting
628 application in plant breeding.

629 **Conclusions**

630 The analysis of a large and representative set of *Malus x domestica* genotypes indicated
631 that apple germplasm diversity reflects its origination within three main geographic regions of
632 Europe, and that a weak genetic structure exists at the European level. This structuring of
633 genetic variation in European dessert apple is caused by evolutionary processes relevant to the
634 domestication of perennial fruit species with factors such as gene flow created by, e.g., ancient
635 roads of commerce across the continent, other human activities like intentional selection and
636 later breeding, and genetic drift. The remarkable differences in the allelic variation found at
637 group and subgroup levels of germplasm stratification constitute a strong indication of that the
638 diversity is hierarchically organized into three *genepools*, with consistent evidence of a pattern
639 of internal substructure. The potential value for modern fruit production is mostly unknown
640 since a majority of the accessions are poorly evaluated from an agronomic point of view. Thus,
641 phenotypic data obtained with standardized methods is required to determine the commercial
642 potential of the preserved material and to enable its use in new crosses to increase the genetic
643 basis of the cultivated apple.

644 The integration of data for collections from different European geographic regions using
645 standardized methods will undoubtedly form an important step in developing the European

646 strategy for conservation of apple germplasm and constitute the starting point to define a
647 European “apple core collection”. This will constitute a decisive step in the field of conservation
648 genetics, and may also have direct implications on the improvement of our understanding of the
649 species, including i) the identification of genomic regions associated with commercially
650 important horticultural traits, ii) the discovery of new germplasm features that may be taken
651 advantage of for efficient breeding and iii) the analysis of genotype x environmental interactions
652 for studying the stability of the most economically important traits for this species.

653 **Methods**

654 **Plant material**

655 Apple germplasm collections from nine European countries, plus Western part of Russia
656 and Kyrgyzstan, were available for this study (Additional file 1): France (INRA, Institut
657 National de la Recherche Agronomique, 399 accessions), Italy (UNIBO, University of Bologna,
658 216 acc.), Belgium (CRA-W, Centre Wallon de Recherche Agronomique, 408 acc.), Czech
659 Republic (RBIPH, Research and Breeding Institute of Pomology Holovousy, 263 acc.), United
660 Kingdom (UK-NFC, University of Reading, 310 acc.), Sweden (SLU, Swedish University of
661 Agricultural Sciences, 199 acc.), Finland (MTT Agrifood Research, 50 acc.), Spain (UPNA,
662 Public University of Navarre, UDL, University of Lleida, and EEAD, Aula Dei Experimental
663 Station, 269 acc.), Switzerland (FRUCTUS, Agroscope, 237 acc.), Russia (NCRRIHV, North
664 Caucasian Regional Research Institute of Horticulture and Viticulture, and VNIISPK, The All
665 Russian Research Institute of Horticultural Breeding, 83 acc.), Kyrgyzstan (KNAU, Kyrgyz
666 National Agrarian University, 12 acc.). In all countries, the accessions were mostly chosen as
667 old local/national dessert cultivars (registered or at least known before 1950), but 12 standard
668 dessert cultivars were also included to strengthen comparisons between collections, namely
669 ‘Golden Delicious’, ‘Red Delicious’, ‘McIntosh’, ‘Rome Beauty’, ‘Granny Smith’, ‘Jonathan’,
670 ‘Winter Banana’, ‘Ingrid Marie’, ‘Ananas Reinette’, ‘Reinette de Champagne’, ‘Discovery’ and
671 ‘Alkmene’. Moreover, 40 old cider apple cultivars and 40 recently-bred dessert cultivars were
672 sampled in the INRA collection in order to investigate particular patterns. Altogether, 2,446

673 accessions were thus considered (Additional file 1). Available collections were somewhat
674 heterogeneous in nature as some of them corresponded to already established core collections
675 (INRA and UNIBO) whereas others were selected for the present study thanks to available SSR
676 marker data (UK-NFC and FRUCTUS, see below), or were chosen as a subset of mainly local
677 cultivars (CRA-W, RBIPH, SLU, MTT, UPNA, UDL, EEAD, NCRRIHV, VNIISPK and
678 KNAU). Cultivars that were known to be triploid or duplicated were avoided since this analysis
679 was performed with an aim to subsequently use a major part of the material in a Genome Wide
680 Association Study to be carried out within the EU FruitBreedomics project [28].

681 **SSR genotyping**

682 A set of 16 SSR markers developed by different groups [59–62] was used to genotype the
683 2,446 accessions (Additional file 3). These SSR markers are distributed over 15 out of the 17
684 apple linkage groups, and 15 of them are included in a former list recommended by the ECPGR
685 *Malus/Pyrus* working group [63]. The 16th marker of this list, NZ05g08, was replaced by the
686 marker CH-Vf1 because the former showed either complex scoring pattern or low level of
687 polymorphism in previous studies [23, 26]. SSR marker data were fully available for the
688 collection from INRA [23]. SSR data were available (i.e., for some, but not all of the 16 SSR
689 markers) for collections from UK-NFC [64], FRUCTUS [14], UPNA, UDL and EEAD [26],
690 and UNIBO [24], so that only the missing SSR marker data were generated in the present study.
691 Fully new SSR datasets were generated for collections from CRA-W, RBIPH, SLU, MTT,
692 NCRRIHV, VNIISPK, and KNAU.

693 Forward primers were labeled with four different fluorescent dyes (6-FAM, VIC, NED, or
694 PET) in order to be combined into four different multiplexed reactions (Additional file 3).
695 Polymerase chain reactions (PCR) for the four multiplex PCRs were performed in a final
696 volume of 11 μ L using 10 ng of DNA template, 0.18 μ M of each primer (with the exception of
697 some markers as described in Additional file 3), and 1 \times PCR Master mix of QIAGEN kit
698 multiplex PCR (Qiagen, Hilden, Germany). PCR cycling conditions were as follows: pre-
699 incubation for 15 min at 94°C, followed by 4 cycles using a touchdown amplification program

700 with an annealing temperature reduced by 1°C per cycle from 60°C to 55°C, followed by 34
701 cycles, each consisting of 30 s denaturing at 94°C, 90 s annealing at 55°C, and 60 s elongation
702 at 72°C, the last cycle ending with a final 15-min extension at 72°C. SSR amplification products
703 were analyzed with an ABI3730 XL sequencing system (Applied Biosystems, Foster City, CA,
704 USA). Fragment analysis and sizing were carried out using GeneMapper v.4.0 software
705 (Applied Biosystems, Foster City, CA, USA); chromatograms were independently read by two
706 operators. When SSR marker data were already available and obtained at different sites, SSR
707 allele sizes were carefully adjusted between collections, both by use of reference accessions
708 known to be in common between collections and by re-genotyping a subset of each collection
709 with the full set of 16 SSR markers to confirm the allele adjustment.

710

711 **Diversity assessments**

712 The multilocus SSR profiles were compared pairwise in order to establish the genetic
713 uniqueness of each accession. Accessions were considered as duplicates if they had identical
714 SSR fingerprints, or if they had one allelic difference for a maximum of two SSR loci thus
715 making room for some genotyping errors and/or spontaneous SSR mutations. On this basis,
716 redundant profiles were removed from the dataset to avoid bias in genetic analyses and
717 duplicate groups were labeled with unique group ID codes (FBUNQ codes). An accession was
718 declared as a putative triploid when at least three of the 16 SSR loci exhibited three distinct
719 alleles. Analyses of descriptive diversity statistics were conducted at locus level. For each SSR
720 marker, SPAGeDi v.1.3 software [65] was used to estimate the number of alleles (N_A), the
721 number of alleles with a frequency below 5% (N_B), the number of effective alleles (N_E), and the
722 observed (H_o) and expected (H_e) heterozygosity. The probability of identity (P_{ID}) was
723 calculated as follows [66]:

$$724 \quad P_{ID} = \sum p_i^4 + \sum \sum (2p_i p_j)^2$$

725 where p_i and p_j are the frequencies of the i^{th} and j^{th} alleles and $i \neq j$. The cumulative P_{ID} over the
726 16 SSR was computed as the product of the P_{ID} of each individual marker.

727

728 **Determination of the geographical regions of origin of the unique genotypes**

729 Using passport data along with reviewing published records with a focus on old literature
730 (national compilations/varietal catalogues/reports) and specialized websites we were able to
731 discern the geographical regions of origin for a large part of the unique genotypes analyzed.
732 This was further helped by the resolution of identified duplicates and comparison of accessions
733 against additional SSR data of the whole UK-NFC apple collection kindly made available from
734 the UK-NFC database [64] and of the whole FRUCTUS collection kindly made available by
735 Agroscope [14]. We first decided to define three broad historical European regions of origin of
736 the germplasm according to geographical proximity and traditional agricultural relations
737 between them: North+East (Sweden, Norway, Finland, Denmark, Baltic countries, plus Russia,
738 Ukraine and Kyrgyzstan), West (Ireland, United Kingdom, France, Belgium, the Netherlands,
739 Switzerland, Germany, Czech Republic) and South (Spain and Italy). When available, countries
740 of origin of the cultivars were also documented although, this information should be considered
741 with caution since the information on the countries of origin was not always fully consistent
742 within duplicates groups.

743

744 **Analysis of the genetic structure**

745 The software Structure v.2.3.4 [67] was used to estimate the number of hypothetical
746 subpopulations (K) and to quantify the proportion of ancestry of each genotype to the inferred
747 subpopulations. No prior information about the geographical origin of the accessions was
748 considered in the analysis. Ten independent runs were carried out for K values ranging from two
749 to 10 using 500,000 Markov Chain Monte Carlo (MCMC) iterations after a burn-in of 200,000
750 steps assuming an admixture model and allelic frequencies correlated. In order to assess the best
751 K value supported for our dataset, the ΔK method [68] was used through the Structure harvester
752 v.0.6.93 website [69] to examine the rate of change in successive posterior probabilities over the
753 range of K values. When the results described above suggested additional substructuring of the
754 diversity in subgroups, a second-level (nested) application of the Structure clustering method

755 was carried out analyzing separately each of the K major groups previously obtained [10, 24,
756 26, 50, 51]. Genotypes were assigned to the group (or sub-group) for which they showed the
757 highest membership coefficient, considering an accession strongly assigned to each partitioning
758 level if its proportion of ancestry (qI) was ≥ 0.80 [70–72]; otherwise they were considered as
759 “admixed”. The placement of genotypes on groups (or sub-groups) was determined using
760 CLUMPP v.1.1 [73], which evaluates the similarity of outcomes between population structure
761 runs. CLUMPP output was used directly as input for Distruct v1.1 [74] in order to graphically
762 display the results.

763 To validate the genetic structure revealed by the Bayesian model-based clustering two
764 complementary approaches using the Darwin software package v6.0.10 [75] were considered: i)
765 an unweighted neighbor-joining tree constructed based on dissimilarities between the unique
766 genotypes (using a Simple Matching coefficient), and ii) a multivariate Principal Coordinate
767 Analysis (PCoA).

768

769 **Genetic differentiation**

770 Population differentiation was estimated by analyses of molecular variance (AMOVA)
771 through Genodive [76] under two scenarios: i) three broad European geographic regions of
772 origin of the material (North+East, West and South); and ii) the major groups (and sub-groups)
773 defined by Structure. Pairwise F_{ST} estimates for the different partitioning levels considered in
774 each case were also obtained using Genodive [76]. Descriptive statistics were calculated for the
775 material clustered according to geographical regions of origin as well as for each group (or sub-
776 group) identified by the Bayesian model-based clustering method, including H_o and H_e , number
777 of total alleles, number of private alleles, i.e., those only found in one (sub)division level, and
778 number of unique alleles, i.e., those only detected in one unique accession. The software
779 FSTAT v.2.9.3.2 [77] was applied to compute the allelic richness after scaling down to the
780 smallest partitioning level in the different scenarios considered.

781

782 **Parentage reconstruction and relatedness between the accessions**

783 On the basis of the SSR profiles of the unique genotypes, accessions were analyzed to infer
784 possible parent-offspring relationships using Cervus v.3.0 software [53]. In order to reveal only
785 robust parentages, we limited the study to the inferences of ‘two-parents offspring’ relationships
786 and did not consider inferences of ‘one-parent offspring’ relationships where the lacking parent
787 offers more flexibility but more speculative assignments as well, especially with only 16 SSR
788 markers. Two criteria were considered to establish strict parentage relationships: i) a confidence
789 level of the LOD score and ii) the Delta LOD value (defined as the difference in LOD scores
790 between the first and second most likely two-candidate parents inferred) both higher than 95%.
791 Finally, an additional constraint was added to strengthen the results by limiting the maximum
792 number of tolerated locus mismatches to only one in any inferred two-parents offspring trio,
793 assuming that such a slight difference may be attributable to possible scoring errors, occurrence
794 of null alleles or occasional mutational events [54, 78].

795

796 **Declarations**

797 **List of abbreviations**

798 **SSR:** Simple Sequence Repeat

799 **PCR:** Polymerase Chain Reaction

800 **LOD:** Logarithm of Odds ratio

801 **F_{ST}:** fixation index ‘F-statistics’

802 **ECPGR:** European Cooperative Programme for Plant Genetic Resources European.

803

804 **Ethics approval and consent to participate**

805 Not applicable.

806

807 **Consent for publication**

808 Not applicable.

809

810 **Competing interests**

811 The authors declare that they have no competing interests.

812

813 **Funding**

814 This work has been partly funded under the EU seventh Framework Programme by the
815 FruitBreedomics project N°265582: “Integrated approach for increasing breeding efficiency in
816 fruit tree crops”. The views expressed in this work are the sole responsibility of the authors and
817 do not necessarily reflect the views of the European Commission. Genotyping of the Spanish
818 collection was partially funded by INIA, Instituto Nacional de Investigación y Tecnología
819 Agraria y Alimentaria (project grant no RF2004-008-C03-00). Genotyping of the Swiss
820 collection was funded by the Swiss Federal Office for Agriculture. Genotyping of the French
821 collection was funded by the FRB, ‘Fondation pour la Recherche sur la Biodiversité’. Initial
822 genotyping of the UK-NFC material was funded by the UK Depart for the Environment Food
823 and Rural Affairs (Defra), grant GC0140. Providing of VNIISPK material (DNA isolation and
824 delivery) have been done with support of Russian Scientific Fund, Project 14-1600127. JU has
825 been partially supported by an Early Stage Research Fellowship of the Institute of Advanced
826 Studies (University of Bologna).

827

828 **Authors’ contributions**

829 JU and CD carried out the statistical analyses. CD and ER carried out the molecular analysis of
830 the accessions not previously genotyped and CD performed the alignment of the SSR profiles of
831 the whole dataset. CP coordinated the fragment analyses of the accessions not previously
832 genotyped. ML, PH, MO, FP, JS, HN, LGG, CM, JG, IS, AVP, LD, and ST contributed in the
833 classification of the plant material in the three broad European geographic regions and national
834 origins, and in the interpretation of the SSR duplicated groups, synonyms and parentages
835 inferred. AG, RG, LF, ML, PH, MO, FFF, KME, FP, JS, HN, LGG, CM, JG, MK, IS, AVP,
836 NGK, ET, LD and ST provided plant material, DNA or SSR profiles of accessions preserved in
837 the studied germplasm. CED conceived and coordinated the study. FL coordinated the EU

838 FruitBreedomics project. JU and CED wrote the manuscript. ML, MO, HN, LGG, CM, MK, IS,
839 AVP, and ST critically reviewed the manuscript. All authors read and approved the final
840 manuscript.

841

842 **Authors' information**

843 Jorge Urrestarazu and Caroline Denancé share first authorship.

844

845 **Availability of supporting data**

846 The dataset supporting the conclusions of this article will be available in the Genome Database
847 for Rosaceae (GDR) (<https://www.rosaceae.org/>).

848

849 **Acknowledgments**

850 The staff at the ANAN genotyping platform of the SFR 149 QUASAV (Angers, France), at the
851 GENTYANE genotyping platform (INRA, Clermont-Ferrand, France) and at the Unité
852 Expérimentale Horticole (INRA, Angers, France) are warmly acknowledged for their help in
853 producing the genotyping data and for maintaining the apple germplasm collection,
854 respectively. The INRA MIGALE bioinformatic platform is also greatly acknowledged for
855 having given support to intense data analyses. Natural Resources Institute, Finland, is gratefully
856 acknowledged for providing access to 50 Finnish apple cultivars in their germplasm collections.
857 UDL, University of Lleida, and EEAD, Aula Dei Experimental Station, are gratefully
858 acknowledged for providing access to apple cultivars in their germplasm collections.
859 Acknowledgements to Slepkov S. from Maykop experimental station of Vavilov's N.I. Vavilov
860 Institute of Plant Industry (MOSVIR) for making available plant material of Russian apple
861 cultivars included in the present analysis.

862

863 **References**

- 864 1. Food and Agriculture Organization of the United Nations. FAO statistics database on the
865 World Wide Web <http://faostat.fao.org/site/567/DesktopDefault.aspx?PageID0567#anchor>.
866 Accessed 27 July 2015.
867
- 868 2. Hokanson SC, Lamboy WF, Szewc-McFadden AK, McFerson JR. Microsatellite (SSR)
869 variation in a collection of *Malus* (apple) species and hybrids. *Euphytica*. 2001;118:281–
870 94.
871
- 872 3. Janick J, Moore JN (1996) Fruit breeding. Volume I: tree and tropical fruits. New York:
873 Wiley; 1996.
874
- 875 4. Laurens F, Durel CE, Patocchi A, Peil A, Salvi S, Tartarini S, Velasco R, van de Weg WE.
876 Review on apple genetics and breeding programs and presentation of a new initiative of a
877 news European initiative to increase fruit breeding efficiency. *J Fruit Sci*. 2010;27:102–7.
878
- 879 5. Noiton DAM, Alspach PA. Founding clones, inbreeding, coancestry, and status number of
880 modern apple cultivars. *J Am Soc Hortic Sci*. 1996;121:773–82.
881
- 882 6. Brooks RM, Olmo HP. Register of new fruit and nut varieties list 35. *HortScience*.
883 1991;26:951–78.
884
- 885 7. Brooks RM, Olmo HP. Register of new fruit and nut varieties list 36. *HortScience*.
886 1994;29:942–69.
887
- 888 8. Brooks RM, Olmo HP. Register of fruit and nut varieties. Alexandria: ASHS; 1997.
889
- 890 9. de Vicente. The evolving role of genebanks in the fast-developing field of molecular
891 genetics. Ed. Issues in genetic resources No XI, August 2004. Rome: International Plant
892 Genetic Resources Institute-IPGRI; 2004.
893
- 894 10. Jing RC, Vershinin A, Grzebyta J, Shaw P, Smykal P, Marshall D, Ambrose MJ, Ellis
895 THN, Flavell AJ. The genetic diversity and evolution of field pea (*Pisum*) studied by high
896 throughput retrotransposon based insertion polymorphism (RBIP) marker analysis. *BMC*
897 *Evol Biol*. 2010;10:44.
898
- 899 11. van Treuren R, Kemp H, Ernsting G, Jongejans B, Houtman H, Visser L. Microsatellite
900 genotyping of apple (*Malus x domestica* Borkh.) genetic resources in the Netherlands:
901 application in collection management and variety identification. *Genet Resour Crop Evol*.
902 2010;57:853–65.
903
- 904 12. Nybom H, Weising K, Rotter B. DNA fingerprinting in botany: past, present, future.
905 *Investig Genet*. 2014;5:1.
906
- 907 13. Bacilieri R, Lacombe T, Le Cunff L, Di Vecchi-Staraz M, Laucou V, Genna B, Péros JP,
908 This P, Boursiquot JM. Genetic structure in cultivated grapevines is linked to geography
909 and human selection. *BMC Plant Biol*. 2013;13:25.
910
- 911 14. Bühlmann A, Gassmann J, Ingenfeld A, Hunziker K, Kellerhals M, Frey JE. Molecular
912 characterization of the Swiss fruit genetic resources. *Erwerbs-Obstbau*. 2015;57:29–34.
913
- 914 15. McCouch SR, McNally KL, Wang W, Sackville Hamilton R. Genomics of gene banks: a
915 case study in rice. *Am J Bot*. 2012;99:407–23.
916

- 917 16. Wünsch A, Hormaza JI. Cultivar identification and genetic fingerprinting of temperate fruit
918 tree species using DNA markers. *Euphytica*. 2002;125:59–67.
919
- 920 17. Ellegren H. Microsatellites: simple sequences with complex evolution. *Nat Rev Genet*.
921 2004;5:435–45.
922
- 923 18. Guichoux E, Lagache L, Wagner S, Chaumeil P, Léger P, Lepais O, Lepoittevin C,
924 Malausa T, Revardel E, Salin F, Petit RJ. Current trends in microsatellite genotyping. *Mol*
925 *Ecol Resour*. 2011;11:591–611.
926
- 927 19. Kalia RK, Rai MK, Kalia S, Singh R, Dhawan AK. Microsatellite markers: an overview of
928 the recent progress in plants. *Euphytica*. 2011;177:309–34.
929
- 930 20. Sunnucks P. Efficient genetic markers for population biology. *Trends Ecol Evol*.
931 2000;15:199–03.
932
- 933 21. Urrestarazu J, Royo JB, Santesteban LG, Miranda CM. Evaluating the influence of the
934 microsatellite marker set on the genetic structure inferred in *Pyrus communis* L. *PLoS One*
935 2015;10:e0138417.
936
- 937 22. Garkava-Gustavsson L, Kolodinska Brantestam A, Sehic J, Nybom H. Molecular
938 characterisation of indigenous Swedish apple cultivars based on SSR and S-allele analysis.
939 *Hereditas*. 2008;145:99–112.
940
- 941 23. Lassois L, Denancé C, Ravon E, Guyader A, Guisnel R, Hibrand-Saint-Oyant L, Poncet C,
942 Lasserre-Zuber P, Feugey L, Durel CE. Genetic diversity, population structure, parentage
943 analysis and construction of core collections in the French apple germplasm based on SSR
944 markers. *Plant Mol Biol Rep*. 2015; doi:10.1007/s11105-015-0966-7.
945
- 946 24. Liang W, Dondini L, De Franceschi P, Paris R, Sansavini S, Tartarini S. Genetic diversity,
947 population structure and construction of a core collection of apple cultivars from Italian
948 germplasm. *Plant Mol Biol Rep*. 2015;33:458–73.
949
- 950 25. Patzak J, Paprštejn F, Henychová A, Sedlák J. Comparison of genetic diversity structure
951 analyses of SSR molecular marker data within apple (*Malus × domestica*) genetic
952 resources. *Genome*. 2012;55:647–65.
953
- 954 26. Urrestarazu J, Miranda C, Santesteban LG, Royo JB. Genetic diversity and structure of
955 local apple cultivars from Northeastern Spain assessed by microsatellite markers. *Tree*
956 *Genet Genomes*. 2012;8:1163–80.
957
- 958 27. Sehic J, Garkava-Gustavsson L, Nybom H. More harmonization needed for DNA-based
959 identification of apple germplasm. *Acta Hort*. 2013;976: 277–83.
960
- 961 28. Laurens F, Aranzana MJ, Arus P, Bassi D, Bonany J, Corelli L, Davey M, Durel CE,
962 Guerra W, Pascal T, Patocchi A, Peace C, Peil A, Quilot-Turion B, Rees J, Troillard V,
963 Stella A, Troggio M, Velasco R, White A, Gao ZS, Van de Weg WE. The new EU project
964 Fruitbreedomics: an integrated approach for increasing breeding efficiency in fruit tree
965 crops. *Plant and Animal Genome XX*, 14-18 January 2012, San Diego, CA (poster)
966
- 967 29. Morin PA, Manaster C, Mesnick SL, Holland R. Normalization and binning of historical
968 and multi-source microsatellite data: overcoming the problems of allele size shift with
969 ALLELOGRAM. *Mol Ecol Resour*. 2009;9:1451–5.
970

- 971 30. Putman AI, Carbone I. Challenges in analysis and interpretation of microsatellite data for
972 population genetic studies. *Ecol Evol.* 2014;4:4399–28.
973
- 974 31. Gasi F, Simon S, Pojskic N, Kurtovic M, Pejic I. Genetic assessment of apple germplasm
975 in Bosnia and Herzegovina using microsatellite and morphologic markers. *Sci Hortic.*
976 2010;126:164–71.
977
- 978 32. Gharghani A, Zamani Z, Talaie A, Oraguzie NC, Fatahi R, Hajnajari H, Wiedow C,
979 Gardiner SE. Genetic identity and relationships of Iranian apple (*Malus x domestica*
980 Borkh.) cultivars and landraces, wild *Malus* species and representative old apple cultivars
981 based on simple sequence repeat (SSR) marker analysis. *Genet Resour Crop Evol.*
982 2009;56:829–42.
983
- 984 33. Gao Y, Liu F, Wang K, Wang D, Gong X, Liu L, Richards CM, Henk AD, Volk GM.
985 Genetic diversity of *Malus* cultivars and wild relatives in the Chinese National Repository
986 of Apple Germplasm Resources. *Tree Genet Genomes.* 2015; doi:10.1007/s11295-015-
987 0913-7.
988
- 989 34. Gross BL, Volk GM, Richards CM, Forsline CL, Fazio G, Chao CT. Identification of
990 “duplicate” accessions within the USDA-ARS National Plant Germplasm System *Malus*
991 Collection. *J Am Soc Hortic Sci.* 2012;137:333–42.
992
- 993 35. Cornille A, Gladioux P, Smulders MJM, Roldán-Ruiz I, Laurens F, Le Cam B, Nerseyan
994 A, Clavel J, Olonova M, Feugey L, Gabrielyan I, Zhang XG, Tenailon MI, Giraud. New
995 insight into the history of domesticated apple: secondary contribution of the European wild
996 apple to the genome of cultivated varieties. *PLoS Genet.* 2012;8:e1002703.
997
- 998 36. Cornille A, Giraud T, Smulders MJM, Roldán-Ruiz I, Gladioux P. The domestication and
999 evolutionary ecology of apples. *Trends Genet.* 2014;30:57–65.
1000
- 1001 37. Gross BL, Henk AD, Richards CM, Fazio G, Volk GM. Genetic diversity in *Malus x*
1002 *domestica* (*Rosaceae*) through time in response to domestication. *Am J Bot.*
1003 2014;101:1770–9.
1004
- 1005 38. Hartmann HT, Kester DE, Davies FT, Geneve RL. Plant propagation: principles and
1006 practices. 7th ed. Upper Saddle River, NJ: Prentice Hall; 2002.
1007
- 1008 39. Zohary D, Hopf D. Domestication of plants in the Old World: the origin and spread of
1009 cultivated plants in West Asia, Europe and the Nile Valley. Oxford: Oxford University
1010 Press; 2000.
1011
- 1012 40. De Franceschi P, Dondini L, Sanzol J. Molecular bases and evolutionary dynamics of self-
1013 incompatibility in the *Pyrinae* (*Rosaceae*). *J Exp Bot.* 2012;63:4015–32.
1014
- 1015 41. Knutson L, Stoner AK. Biotic diversity and germplasm preservation: global imperative.
1016 Kluwer Academic Publishers; 1989.
1017
- 1018 42. McCouch S. Diversifying selection in plant breeding. *PLoS Biol.* 2004;2:e347.
1019
- 1020 43. World Apple and Pear Association. European apple and pear crop forecast. Brussels:
1021 World Apple and Pear Association; 2013.
1022
- 1023 44. Miller AJ, Gross BL. Forest to field: perennial fruit crop domestication. *Am J Bot.*
1024 2011;98:1389–1414.
1025

- 1026 45. Biswas S, Akey JM. Genomic insights into positive selection. Trends Genet. 2006;22 :437–
1027 46.
1028
- 1029 46. Bourguiba H, Audergon JM, Krichen L, Trifi-Farah N, Mamouni A, Trabelsi S, D'Onofrio
1030 C, Asma BM, Santoni S, Khadari B. Loss of genetic diversity as a signature of apricot
1031 domestication and diffusion into the Mediterranean Basin. BMC Plant Biol. 2012;12:49.
1032
- 1033 47. Blanquart F, Gandon S, Nuismer SL. 2012. The effects of migration and drift on local
1034 adaptation to a heterogeneous environment. J Evol Biol. 2012;25:1351–63.
1035
1036
- 1037 48. Emanuelli F, Lorenzi S, Grzeskowiak L, Catalano V, Stefanini M, Troglio M, Myles S,
1038 Martinez-Zapater JM, Zyprian E, Moreira FM, Grando MS. Genetic diversity and
1039 population structure assessed by SSR and SNP markers in a large germplasm collection of
1040 grape. BMC Plant Biol. 2013;13:39.
1041
- 1042 49. Lia VV, Poggio L, Confalonieri VA. Microsatellite variation in maize landraces from
1043 Northwestern Argentina: genetic diversity, population structure and racial affiliations.
1044 Theor Appl Genet. 2009;119:1053–67.
1045
- 1046 50. Li XW, Meng XQ, Jia HJ, Yu ML, Ma RJ, Wang LR, Cao K, Shen ZJ, Niu L, Tian JB,
1047 Chen MJ, Xie M, Arus P, Gao ZS, Aranzana MJ. Peach genetic resources: diversity,
1048 population structure and linkage disequilibrium. BMC Genet. 2013;14:84.
1049
- 1050 51. Jacobs MMJ, Smulders MJM, van den Berg RG, Vosman B. What's in a name; genetic
1051 structure in *Solanum* section *Petota* studied using population-genetic tools. BMC Evol
1052 Biol. 2011;11:42.
1053
- 1054 52. Jing R, Ambrose MA, Knox MR, Smykal P, Hybl M, Ramos A, Caminero C, Burstin J,
1055 Duc G, van Soest LJM, Swiecicki WK, Pereira MG, Vishnyakova M, Davenport GF,
1056 Flavell AJ, Ellis THN. Genetic diversity in European *Pisum* germplasm collections Theor
1057 Appl Genet. 2012;125:365–80.
1058
- 1059 53. Kalinowski ST, Taper ML, Marshall TC. Revising how the computer program cervus
1060 accommodates genotyping error increases success in paternity assignment Mol Ecol.
1061 2007;16:1099–2006.
1062
- 1063 54. Salvi S, Micheletti D, Magnago P, Fontanari M, Viola R, Pindo M, Velasco R. One-step
1064 reconstruction of multi-generation pedigree networks in apple (*Malus × domestica* Borkh.)
1065 and the parentage of Golden Delicious. Mol Breed. 2014;34:511–24.
1066
- 1067 55. Bianco L, Cestaro A, Sargent DJ, Banchi E, Derdak S, Di Guardo M, Salvi S, Jansen J,
1068 Viola R, Gut I, et al. Development and validation of a 20K single nucleotide polymorphism
1069 (SNP) whole genome genotyping array for apple (*Malus x domestica* Borkh). PloS ONE.
1070 2014;9:e110377.
1071
- 1072 56. Bianco L, Cestaro A, Linsmith G, Muranty H, Denancé C, Théron A, Poncet C, Micheletti
1073 D, Kersshbamer E, Di Pierro EA, Larger S, Pindo M, van de Weg WE, Davassi A, Laurens
1074 F, Velasco R, Durel CE, Troglio M. Development and validation of the
1075 Axiom®Apple480K SNP genotyping array. 2016; submitted.
1076
- 1077 57. Chagné D, Crowhurst RN, Troglio M, Davey MW, Gilmore B, Lawley C, Vanderzande S,
1078 Hellens RP, Kumar S, Cestaro A et al. Genome-wide SNP detection, validation, and
1079 development of an 8K SNP array for apple. PloS ONE. 2012;7:e31745.

1080
1081 58. Morgan J, Richards A, Dowle E. The new book of Apples: the definitive guide to apples,
1082 including over 2000 varieties. London: Ebury Press; 2002
1083
1084 59. Hokanson SC, Szewc-McFadden AK, Lamboy WF, McFerson JR. Microsatellite (SSR)
1085 markers reveal genetic identities, genetic diversity and relationships in a *Malus x domestica*
1086 Borkh. core subset collection. Theor Appl Genet. 1998;97:671–83.
1087
1088 60. Liebhard R, Gianfranceschi L, Koller B, Ryder CD, Tarchini R, van de Weg E, Gessler C.
1089 Development and characterisation of 140 new microsatellites in apple (*Malus x domestica*
1090 Borkh.). Mol Breed. 2002;10:217–41.
1091
1092 61. Silfverberg-Dilworth E, Matasci CL, van de Weg WE, van Kaauwen MPW, Walser M,
1093 Kodde LP, Soglio V, Gianfranceschi L, Durel CE, Costa F, Yamamoto T, Koller B,
1094 Gessler C, Patocchi A. Microsatellite markers spanning the apple (*Malus x domestica*
1095 Borkh.) genome. Tree Genet Genomes 2006;2:202–24.
1096
1097 62. Vinatzer BA, Patocchi A, Tartarini S, Gianfranceschi L, Sansavini S, Gessler C. Isolation
1098 of two microsatellite markers from BAC clones of the Vf scab resistance region and
1099 molecular characterization of scab-resistant accessions in *Malus* germplasm. Plant Breed.
1100 2004;123:321–6.
1101
1102 63. Evans KM, Fernández F, Govan C. Harmonising fingerprinting protocols to allow
1103 comparisons between germplasm collections - *Pyrus*. Acta Hortic. 2009;814:103–6.
1104
1105 64. Fernández-Fernández F. Fingerprinting the National apple and pear collections. Final
1106 report of DEFRA research project GC0140.
1107 <http://randd.defra.gov.uk/Document.aspx?Document=GC0140SID5FingerprintingFernandez.pdf>;
1108 2010. 1–18.
1109
1110 65. Hardy OJ, Vekemans X. SPAGEDi: a versatile computer program to analyse spatial
1111 genetic structure at the individual or population levels. Mol Ecol Notes. 2002;2:618–20.
1112
1113 66. Waits LP, Luikart G, Taberlet P. Estimating the probability of identity among genotypes in
1114 natural populations: cautions and guidelines. Mol Ecol. 2001;10:249–56.
1115
1116 67. Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus
1117 genotype data. Genetics. 2000;155:945–59.
1118
1119 68. Evanno G, Regnaut S, Goudet J. Detecting the number of clusters of individuals using the
1120 software STRUCTURE: a simulation study. Mol Ecol. 2005;14:2611–20.
1121
1122 69. Earl DA, vonHoldt BM. STRUCTURE HARVESTER: a website and program for
1123 visualizing STRUCTURE output and implementing the Evanno method. Cons Genet
1124 Resour. 2012;4:359–61.
1125
1126 70. Breton C, Pinatel C, Médail F, Bonhomme F, Bervillé A. Comparison between classical
1127 and Bayesian methods to investigate the history of olive cultivars using SSR-
1128 polymorphisms. Plant Sci. 2008;175:524–32.
1129
1130 71. Miranda C, Urrestarazu J, Santesteban LG, Royo JB, Urbina V. Genetic diversity and
1131 structure in a collection of ancient Spanish pear cultivars assessed by microsatellite
1132 markers. J Am Soc Hortic Sci. 2010;135:428–37.
1133

- 1134 72. Marra FP, Caruso T, Costa F, Di Vaio C, Mafrica R, Marchese A. Genetic relationships,
1135 structure and parentage simulation among the olive tree (*Olea europaea* L. subsp.
1136 europaea) cultivated in Southern Italy revealed by SSR markers. *Tree Genet Genomes*.
1137 2013;9:961–73.
1138
- 1139 73. Jakobsson M, Rosenberg NA. CLUMPP: a cluster matching and permutation program for
1140 dealing with label switching and multimodality in analysis of population structure.
1141 *Bioinformatics*. 2007;23:801–6.
1142
- 1143 74. Rosenberg NA. DISTRUCT: a program for the graphical display of population structure.
1144 *Mol Ecol Notes*. 2004;4:137–8.
1145
- 1146 75. Perrier X, Jacquemoud-Collet JP. DARwin Software <http://darwin.cirad.fr/darwin>; 2006.
1147
- 1148 76. Meirmans PG, van Tienderen PH. GENOTYPE and GENODIVE: two programs for the
1149 analysis of genetic diversity of asexual organisms. *Mol Ecol Notes*. 2004;4:792–4.
1150
- 1151 77. Goudet J. FSTAT 2.9.3.2. <http://www2.unil.ch/popgen/softwares/fstat.htm>; 2002.
1152
- 1153 78. Lacombe T, Boursiquot JM, Laucou V, Di Vecchi-Staraz M, Péros JP, This P. Large-scale
1154 parentage analysis in an extended set of grapevine cultivars (*Vitis vinifera* L.). *Theor Appl*
1155 *Genet*. 2013;126:401–14.
1156
1157
1158

1159 **Tables**

1160

1161 **Table 1** Average measures of genetic diversity at two different levels: overall set of accessions
 1162 and according to the three geographical regions of origin (North+Eastern, Southern and
 1163 Western). Number of alleles per locus (NA), number of rare alleles (NB), effective number of
 1164 alleles (NE), allelic richness (AR), and observed (Ho) and expected (He) heterozygosity are
 1165 included

Material	NA	NB ^a	NE	AR ^b	Ho	He
Overall set (1859 genotypes)	23.06	16.94	6.59	-	0.81	0.83
European regions of origin						
Northern+Eastern Europe	16.75	10.87	6.24	16.57	0.83	0.82
Southern Europe	17.50	11.87	6.29	16.95	0.81	0.82
Western Europe	20.31	13.94	6.18	16.36	0.81	0.82

1166 ^a Rare alleles were considered if they appeared in a frequency below 5%

1167 ^b For the geographical European regions of origin, allelic richness was computed after normalization according to the smallest
 1168 population size (i.e., Northern+Eastern Europe)

1169

1170 **Table 2** Descriptive information for each of the three major groups and eight subgroups of
 1171 genotypes identified by the Bayesian model-based clustering method. Summary statistics
 1172 include the partitioning of number of individuals in each group, expected heterozygosity (He),
 1173 total, private, unique, and average number of alleles (A). Allelic richness is scaled to the
 1174 smallest group (K2; N=401) or subgroup (K1.3; N=148)

Group/Subgroup	Number of genotypes in the group/subgroup		He	Number of alleles				Allelic richness
	Number Genotypes	Frequency of genotypes with $qI \geq 0.8$		Total	Private	Unique	A	
K1	506	60%	0.823	307	34	16	19.19	18.63
K2	401	57%	0.816	287	23	15	17.94	17.76
K3	952	67%	0.801	294	22	14	18.36	15.99
K1.1	209	42%	0.842	282	17	12	17.63	16.38
K1.2	149	54%	0.789	215	3	1	13.44	13.20
K1.3	148	36%	0.761	228	6	3	14.25	13.86
K2.1	244	48%	0.818	268	14	11	16.75	14.73
K2.2	157	53%	0.778	211	5	4	13.19	12.67
K3.1	375	41%	0.775	242	7	6	15.13	12.32
K3.2	162	57%	0.760	171	0	0	10.69	10.31
K3.3	415	51%	0.809	255	14	8	15.94	13.43

1175

1176 **Table 3** Analysis of molecular variance (AMOVA) based on the 16 SSR loci of the apple germplasm evaluated in this study corresponding to three regions of
 1177 origin (Northern+Eastern, Southern and Western Europe) and groups and subgroups defined by Structure analysis

1178

Populations	df ^a		Variance components (%)		
	W ^b	A	W	A	<i>p</i> value
3, geographic origins	1653	2	97.9	2.1	0.001
3, groups defined by Structure	1859	2	96.9	3.1	0.001
3, subgroups of K1	506	2	96.3	3.7	0.001
2, subgroups of K2	401	1	96.6	3.4	0.001
3, subgroups of K3	952	2	97.3	2.7	0.001
8, subgroups (K1+ K2+ K3)	1859	7	95.6	4.4	0.001

^adf: degrees of freedom, ^bW: within populations, ^cA: among populations

1179

1180 **Table 4** Pairwise estimates of F_{ST} among the eight subgroups obtained by the nested Bayesian
1181 clustering approach.

Subgroup	K1.1	K1.2	K1.3	K2.1	K2.2	K3.1	K3.2	K3.3
K1.1	—							
K1.2	0.030	—						
K1.3	0.035	0.051	—					
K2.1	0.028	0.067	0.077	—				
K2.2	0.049	0.076	0.087	0.035	—			
K3.1	0.034	0.051	0.061	0.061	0.055	—		
K3.2	0.051	0.065	0.070	0.070	0.058	0.029	—	
K3.3	0.016	0.042	0.060	0.038	0.051	0.023	0.038	—

1182 All the estimates were highly significant ($P < 0.001$)

1183

1184 **Table 5** Full parentages of 46 apple cultivars inferred within the set of the 1859 apple unique accessions using 16 SSR markers with their accession codes,
 1185 accession names (AcceNumber), their duplicate codes according to the SSR profile (FBUNQ) and their putative country of origin (OriginHist)

Offspring ID	Accename	FBUNQ	OriginHist	First candidate ID	Accename	FBUNQ	OriginHist	Second candidate ID	Accename	FBUNQ	OriginHist	Status ^j
X1618	Calville Rouge du Mont Dore	963	FRA	BAL086	Alexander	30	UKR	DCA_I05	Mele Ubriache ^a	361	FRA	doc.
X1846	Belle de Mleiev	1563	-	X0557	Mc Intosh	508	CAN	1957218	King of the Pippins	37	FRA	doc.
BAL035	Heta	1774	FIN	CRAW-0433	Lobo	788	CAN	FIN09	Huvitus	4922	FIN	recent
BAL039	Jaspi	1776	FIN	CRAW-0433	Lobo	788	CAN	FIN09	Huvitus	4922	FIN	recent
FIN18	Pirkko	4930	FIN	CRAW-0433	Lobo	788	CAN	FIN09	Huvitus	4922	FIN	recent
BAL010	Rödluvan	107	SWE	CRAW-0433	Lobo	788	CAN	BAL023	Barchatnoje	1768	RUS	recent
BAL109	Arona	1819	LVA	CRAW-0433	Lobo	788	CAN	BAL112	Iedzenu	1822	LVA	recent
BAL176	Nyckelby	1861	SWE?	CRAW-0433	Lobo	788	CAN	1957188	Cox's Pomona	2033	GBR	recent?
BAL059	Pirja	444	FIN	FIN09	Huvitus	4922	FIN	CRAW-0836	Melba	167	CAN	recent
FIN43	Pirkkala	4949	FIN	BAL042	Kaneläpple	512	RUS	FIN14	Lavia	4926	FIN	recent
BAL154	Mio	543	SWE	CZ_G2D_0045	Worcester parména	550	GBR	BAL056	Oranie	48	SWE	recent
BAL052	Oberle	1784	CAN	BAL027	Early Red Bird	236	CAN	CRAW-0266	Stark Earliest	468	USA	old
BAL091	Förlovningsäpple	1804	SWE	CHE0893	Heuapfel	1248	CHE	X1646	Saint Germain	31	-	old
BAL167	Valldaäpple	1853	SWE	CHE0893	Heuapfel	1248	CHE	BAL179	Göteborgs Flickäpple	1863	SWE	old
BAL099	Kramforsäpple	1811	SWE	BAL161	Sävstaholm	573	SWE	BAL195	Åkerö	308	SWE	old
BAL158	Stenkyrke	463	SWE	BAL171	Fullerö	1857	SWE	CZ_LJ_0045	Malinové podzimní ^b	722	POL	old
FIN07	Finne	4920	FIN	BAL161	Sävstaholm	573	SWE	FIN08	Grenman	4921	FIN	old
1942035	Beauty of Moray	1925	GBR	2000053	Keswick Codlin	1438	GBR	2000090	Stirling Castle	2103	GBR	old
1951242	Brighton	2011	NZL?	X4915	Red Dougherty	939	NZL	CZ_LC_0411	Hlaváčkovo ^c	23	USA	old
1957208	Ben's Red	2035	GBR	CRAW-0020	Devonshire Quarrenden	622	GBR	1955077	Box Apple	2025	GBR	old
1965004	Fred Webb	2054	GBR	1946088	Winter Marigold	324	GBR	1957181	Gascoyne's Scarlet	45	GBR	old
2000083	Rivers' Early Peach	2099	GBR	2000051	Irish Peach	2093	IRL	BAL169	Aspa	1855	SWE	old
BMN0011	Roja de Guipuzcoa	3854	ESP	BMN0017	Urte Sagarra	956	ESP	BMN0171	Maxel Gorri	3896	ESP	old
BMZ016	Cella	3935	ESP	BMN0022	Erreka Sagarra	957	ESP	X5102	Bisquet	535	FRA	old
BMN0070	Madotz-01	3869	ESP	1957218	King of the Pippins	37	FRA	X7201	Transparente de Croncels	62	FRA	old
CHE1322	Rotwiler	1271	CHE?	1957218	King of the Pippins	37	FRA	BAL086	Alexander	30	UKR	old

CHE1788	Roseneggler	3718	CHE	1957218	King of the Pippins	37	FRA	CZ_BoN_0429	Trat. Laze	2284	CZE	old
CHE0032	Ernst Bosch	1003	DEU	1947074	Ananas Reinette	69	NLD	CZ_GF_0415	Evino ^d	7	GBR	old
CHE0168	Eibner	3258	CHE	CRAW-0836	Melba	167	CAN	CZ_BoN_0424	Trevinské červené ^e	71	USA	old
CHE1390	Klefeler	3589	CHE	KRAS123	Papirovka	25	RUS	X7199	Rose de Berne	83	CHE	old
CRAW-0226	Laubain n°1	2126	BEL	CRAW-0086	Bismarck	3	AUS	CZ_GS_0478	Ušlechtilé žluté ^f	90	GBR	old
CRAW-0105	Godelieve Hegmans	2116	BEL	BAL175	Röd Astrakan	82	RUS	X7201	Transparente de Croncels	62	FRA	old
CZ_BB_0442	Nathusiovo	2268	DEU	CZ_GL_0464	Bláhovo Libovické	2311	CZE	X7201	Transparente de Croncels	62	FRA	old
CZ_BB_0434	Panenské veliké	2265	CZE	CZ_GP_0469	Panenské české	1529	CZE	X1344	Reinette de Landsberg	61	DEU	old
CZ_GK_0412	Proche	2308	CZE	CRAW-0425	Calville Rouge d'Automne	13	FRA	X1344	Reinette de Landsberg	61	DEU	old
CZ_BoN_0421	Moravcovo	2283	CZE	CZ_GP_0469	Panenské české	1529	CZE	CHE0269	Pomme Bülleöpfungel	1377	-	old
CZ_GL_0456	Bláhův poklad	694	-	CZ_GL_0464	Bláhovo Libovické	2311	CZE	CZ_GG_0438	Malinové hornokrajské ^g	47	NLD	old
CZ_GP_0473	Petr Broich	2321	DEU	1957175	Annie Elizabeth	15	GBR	2000075	Peasgood's Nonsuch	51	GBR	old
CZ_BB_0458	Šarlatová parména	2269	CZE	CZ_GG_0442	Malinové holovouské	452	CZE	X8233	Petite Madeleine	24	-	old
CZ_BB_0466	Podzvičinské ^h	231	-	X0691	Boiken	108	DEU	X1071	Reinette de Caux	629	NLD	old
DCA_017	S.Giuseppe	1646	ITA	DCA_090	Abbondanza	327	ITA	DCA_C44	Rambour Frank (MI)	493	FRA	old
DCA_H03	Scodellino	1642	ITA	DCA_090	Abbondanza	327	ITA	DCA_E52	Decio	397	ITA	old
DCA_E72	Gelato Cola	330	ITA	DCA_E69	Gelato (CT)	780	-	DCA_F74	Limoncella (TN) ⁱ	708	ITA	old
DCA_H62	Liscio di Cumiana	1713	ITA	DCA_H29	Carla	114	-	DCA_C21	Renetta di Grenoble	263	ITA	old
DCA_I96	Ros Magior	1658	ITA	DCA_I80	Rus d' Muslot	321	-	X1115	Rome Beauty	334	USA	old
DCA_F47	Mela Golden Simile di Villa Collemantina	1692	ITA	DCA_A20	Rosa Mantovana (TN)	101	ITA	CRAW-0025	Yellow Bellflower	77	USA	old

1186

1187 ^a DCA_I05 'Mele Ubriache' duplicate with 'Calville Rouge d'Hiver' [23]

1188 ^b based on 11 SSR [64] the accession CZ_LJ_0045 'Malinové podzimní' was shown to be duplicated with 'Danziger Kantapfel'

1189 ^c based on 11 SSR [64] the accession CZ_LC_0411 'Hlaváčkovo' duplicate with 'Nothern Spy'

1190 ^d based on 11 SSR [64] the accession CZ_GF_0415 'Evino' duplicate with 'Mank's Codlin'

1191 ^e based on 11 SSR [64] the accession CZ_BoN_0424 'Trevinské červené' duplicate with 'King David'

1192 ^f based on 11 SSR [64] the accession CZ_GS_0478 'Ušlechtilé žluté' duplicate with 'Golden Noble'

1193 ^g based on 11 SSR [64] the accession CZ_GG_0438 'Malinové hornokrajské' duplicate with 'Framboise'

1194 ^h based on 11 SSR [64] and on 13 SSR [14] the accession CZ_BB_0466 'Podzvičinské' duplicate with 'Altlander Pfannkuchenapfel' and 'Thurgauer Kent'

1195 ⁱ based on 11 SSR [64] the accession DCA_F74 'Limoncella' (TN) duplicate with 'Cola'

1196 ^j recent or old cultivars ; doc. = inferred parentage already documented in [23]

1197

1198 **Caption for Figures**

1199

1200 **Figure 1** Graphical display of the results of the Structure analyses. a1) Proportions of ancestry
1201 of 1859 unique diploid apple genotypes for $K=3$ groups inferred with Structure v.2.3.4 software
1202 [67]. Each genotype is represented by a vertical bar partitioned into $K=3$ segments representing
1203 the estimated membership fraction in three groups. The three groups are depicted using the
1204 following color codes: Red = group K1; Blue = group K2; Green = group K3. a2) Proportions of
1205 ancestry of the same 1859 genotypes following a nested Structure analysis within each
1206 previously defined group. For K1 and K3 three subgroups are shown and for K2 two subgroups
1207 are shown. Each genotype is represented by a vertical bar partitioned into $K=2$ or 3 subgroups
1208 representing the estimated membership fraction in each subgroup. Genotypes are presented in
1209 the same order in a1. The subgroups are depicted using the following color codes: light Pink =
1210 K1.1; Purple = K1.2; dark Pink = K1.3; light Blue = K2.1; dark Blue = K2.2; fluorescent Green
1211 = K3.1; dark Green = K3.2; light Green = K3.3. b) Proportions of ancestry of 1653 unique
1212 diploid apple genotypes with known European region of origin for $K=3$ groups inferred with the
1213 same Structure analysis as in a. The genotypes are sorted according to their European region of
1214 origin (North+East, West, and South).

1215

1216 **Figure 2** Scatter plot of the Principal Coordinate Analysis (PCoA) of the 1859 apple accessions
1217 based on the 16 SSR data. The three groups are depicted using the following color codes: Red =
1218 group K1; Blue = group K2; Green = group K3.

1219

1220 **Figure 3** Neighbor-joining dendrogram based on simple matching dissimilarity matrix
1221 calculated from the dataset of 16 SSR markers for the 1859 genotypes clustered in the three
1222 groups revealed by the Bayesian model-based clustering method. The three groups are depicted
1223 using the following color codes: Red = group K1; Blue = group K2; Green = group K3.

1224

1225 **Figure 4** Genetic composition of the groups of cultivars clustered by country of origin for $K=3$
1226 groups inferred with Structure. For the detailed country list, see Additional file 1. The pies
1227 represent the proportion of each group in each country: color codes are as per Figure 1a1.
1228
1229

1230 **Caption for Additional files**

1231

1232 **Additional file 1** (.xls file) List of the 2446 accessions considered in the present study with their
1233 accession code (AcceNumber), name (AcceName), the name of the providing collection
1234 (Collection), their duplicate code according to the SSR profile (FBUNQ, see text), their ploidy
1235 level (Ploidy) determined according to the occurrences of three alleles per locus (see text), their
1236 status (Analyzed) as analyzed or not-analyzed within the duplicate group (when adequate), their
1237 documented European geographic region of origin (Eur_reg_orig), their putative country of
1238 origin (Country_orig), their group assignment (Group) inferred by the Structure analysis with
1239 the highest proportion of ancestry (*qI*max), and their subgroup assignment (Subgroup) inferred
1240 by the nested Structure analysis with the highest proportion of ancestry (*qI*max nested). In the
1241 ‘*qI*max’ and ‘*qI*max nested’ columns, a bold number indicates that the highest subgroup
1242 proportion of ancestry (*qI*) is equal to or greater than 0.8. The proportions of ancestry for each
1243 of the 3 groups (*qI* K1, *qI* K2, *qI* K3) and for either the 3 (*qI* K1.1, *qI* K1.2, *qI* K1.3), the 2 (*qI*
1244 K2.1, *qI* K2.2), or the 3 (*qI* K3.1, *qI* K3.2, *qI* K3.3) subgroups of groups K1, K2, or K3
1245 (respectively) are then given, the latter subgroups corresponding to the group exhibiting the
1246 highest proportion of ancestry (*qI*max).

1247 In the ‘Analyzed’ column:

- 1248 - ‘A’ indicates an accession that has been considered in the statistical analyses;
- 1249 - ‘E’ indicates an accession that has been excluded from the statistical analyses (mostly
1250 because another duplicated accession has been retained; in that case, the group or
1251 subgroup membership and the *qI* max probability have been imputed according to the
1252 analyzed duplicated accession);
- 1253 - ‘E (SSR)’ indicates an accession that has been excluded from the statistical analyses
1254 because of low number of SSR marker data (< 12 SSR);
- 1255 - ‘E (Ext.)’ indicates an accession that has been excluded from the statistical analyses
1256 because of its status as outlier in a preliminary Principal Coordinate Analysis;

1257 - 'E (Rs)' indicates an accession that has been excluded from the statistical analyses
1258 because of its status as a rootstock identified using the SSR profile (e.g., M9, MM106
1259 or MM111 instead of the expected accession).

1260 In the 'AcceNumber' column, the various colors are only attributed to distinguish the various
1261 collections under study. In the 'FBUNQ' column, the water-green color is attributed to the even
1262 numbers to ease the duplicate group visualization.

1263

1264 **Additional file 2** (.TIFF file) Genetic composition of cultivars clustered by country of origin for
1265 the eight subgroups inferred with Structure. For the detailed country list, see Additional file 1.
1266 The pies represent the proportion of each subgroup in each country: color codes are as per
1267 Figure 1a2.

1268

1269 **Additional file 3** (.xls file) Characteristics of the 16 SSR markers used in this study with
1270 indication of the corresponding multiplex and dye.

1271

1272 Footnotes:

1273 ^a [61]; ^b [60]; ^c [59]; ^d [62]; ^e Primer concentration within a given multiplex has been adjusted to
1274 get more homogeneous SSR marker amplification intensities.