



This is a post-peer-review, pre-copyedit version of an article published in Tree Genetics and Genomes. The final authenticated version is available online at: <https://doi.org/10.1007/s11295-021-01526-7>

Springer Nature terms of use for archived accepted manuscripts (AMs) of subscription articles at: <https://www.springernature.com/gp/open-research/policies/accepted-manuscript-terms>

Document downloaded from:



1 **Comparison of selection methods for the establishment of a core collection**
2 **using SSR markers for hazelnut (*Corylus avellana* L.) accessions from European**
3 **germplasm repositories**

4
5 Paolo Boccacci ¹⁾, Maria Aramini ²⁾, Matthew Ordidge ³⁾, Theo J.L. van Hintum ⁴⁾, Daniela Torello
6 Marinoni ⁵⁾, Nadia Valentini ⁵⁾, Jean-Paul Sarraquigne ⁶⁾, Anita Solar ⁷⁾, Mercè Rovira ⁸⁾, Loretta
7 Bacchetta ²⁾, Roberto Botta ⁵⁾

8
9 ¹⁾ Institute for Sustainable Plant Protection - National Research Council (IPSP-CNR). Strada delle Cacce, 73 -
10 10135 Torino, Italy

11 ²⁾ ENEA - Agenzia nazionale per le nuove tecnologie, l'energia e lo sviluppo economico sostenibile, Via
12 Anguillarese 301 - 00123 S.M. di Galeria (RM), Italy

13 ³⁾ School of Agriculture, Policy and Development, University of Reading, Whiteknights, RG6 6AR, Reading,
14 United Kingdom

15 ⁴⁾ Centre for Genetic Resources, the Netherlands, Wageningen Plant Research, P.O. Box 16, 6700 AA,
16 Wageningen, the Netherlands

17 ⁵⁾ Department of Agricultural, Forestry and Food Science - University of Torino, Largo Paolo Braccini, 2 - 10095
18 Grugliasco (TO), Italy

19 ⁶⁾ Association Nationale des Producteurs de Noisette (ANPN), 47290 Cancon, France

20 ⁷⁾ University of Ljubljana, Biotechnical Faculty, Department of Agronomy, Jamnikarjeva 101, 1000 Ljubljana,
21 Slovenia

22 ⁸⁾ Institut de Recerca i Tecnologia Agroalimentàries (IRTA), Centre Mas Bové, Ctra. Reus-El Morell, km 3.8 -
23 43120 Constantí (Tarragona), Spain

24

25 **Corresponding author:** Paolo Boccacci, e-mail: paolo.boccacci@ipsp.cnr.it

26

27

28 **ORCID ID of the authors**

29 Paolo Boccacci: 0000-0001-8574-0478; Matthew Ordidge: 0000-0003-0115-5218; Theo van Hintum: 0000-0003-

30 4953-4700; Daniela Torello Marinoni: 0000-0002-3679-4813; Nadia Valentini: 0000-0002-8820-9006; Merce

31 Rovira: 0000-0002-0540-3752; Loretta Bacchetta: 0000-0001-8878-4054; Anita Solar: 0000-0001-9755-4998;

32 Roberto Botta: 0000-0002-1952-8775.

1 **Abstract**

2 Hazelnut (*Corylus avellana* L.) is one of the most important tree nut crops in Europe. Accessions are
3 conserved in twelve European *ex situ* germplasm repositories, located in countries where hazelnut
4 production occurs. In this work, we used ten single sequence repeat (SSR) markers as the basis to
5 establish a core collection representative of the hazelnut genetic diversity conserved in different
6 European collections. A total of 480 accessions, 430 from *ex situ* collections and 50 landraces
7 maintained *on-farm*, were used. SSR analysis identified 181 true-to-type genotypes, that represented our
8 whole hazelnut germplasm collection (WHGC). Four approaches (MSTRAT, Power Core, Core Hunter
9 single- and multi-strategy) based on maximization (M) strategy were used to determine the best
10 sampling method. Core Hunter multi-strategy, optimizing simultaneously both allele coverage (Cv) and
11 Cavalli-Sforza and Edward (Dce) distance with equal weight, outperformed the others and was selected
12 as the best approach. The final core collection (Cv-Dce30) comprised 30 entries (16.6%). It recovered
13 all the SSR alleles with the minimum number of accessions and preserved parameter variations when
14 compared to WHGC. Entries covered all six gene pools obtained from the population structure analysis
15 of WHGC, further confirming the representativeness of Cv-Dce30. Our findings contribute towards
16 improving both the conservation and management of European hazelnut genetic resources and could be
17 used to optimize future research by identifying a minimum number of accessions on which to focus.

18

19

20 **Key words:** Filbert; Microsatellite; *Ex situ* and *in situ* conservation; Germplasm management; Plant
21 genetic resources

22

1 Introduction

2 The European hazelnut (*Corylus avellana* L.) is one of the most important tree nut crops in terms of
3 worldwide production (averaging 939,927 tons per annum in 2015-2019). The Black Sea countries
4 account for most of the average annual world production (data 2015-2019): Turkey (606,409 tons),
5 Azerbaijan (43,584 tons), and Georgia (25,440 tons). Other important producers are Italy (116,945 tons),
6 the USA (36,652 tons), Iran (15,583 tons), France (11,994 tons), and Spain (10,364 tons) followed by
7 Chile, Poland, Serbia, Kyrgyzstan, and Uzbekistan (FAOSTAT 2021). World production is based
8 entirely on cultivars selected over many centuries from local wild populations (Thompson et al. 1996)
9 and about 500 cultivars have been described in the literature and are available from one or more *ex situ*
10 germplasm repositories (Köksal 2000; Botta et al. 2019). However, only about 20 cultivars are widely
11 grown and another 30 are considered promising for breeding (Botta et al. 2019). Collections consist
12 primarily of cultivated forms of *C. avellana* and are mainly located in countries where this production
13 occurs. A total of 510 hazelnut accessions, corresponding to 222 cultivars and 58 selections, are
14 conserved in 12 European collection fields: four in Italy, three in Portugal, two in Spain, and one each
15 in Slovenia, France, and Greece (Bacchetta et al. 2015; Botta et al. 2019). More than 700 *Corylus*
16 accessions are preserved in the major world hazelnut collection located in Oregon (USA) (Hummer
17 2001), while a collection containing 20 registered cultivars and more than 400 accessions collected from
18 the Black Sea coast is in Turkey (Öztürk et al. 2017). *In situ* conservation strategies have been applied
19 only recently, after *on-farm* explorations conducted in southern Europe (Ferreira et al. 2010; Boccacci
20 et al. 2013).

21 Germplasm collections ensure the long-term conservation of genetic resources and provide easy
22 access to plant breeders, researchers, and other users. The management and use of large germplasm
23 collections requires significant economic costs for routine tasks, such as conservation, regeneration,
24 duplication, documentation, and evaluation. Moreover, collections invariably contain duplicate and
25 redundant accessions that may invalidate both the efficiency of the conservation and the effectiveness
26 of germplasm evaluation and use (van Hintum et al. 2000). Consequently, the long-term conservation
27 of collections can be endangered. Thus, a core collection concept was introduced in the 1980s to define

1 a limited set of accessions from the whole collection that represents, with a minimum of repetitiveness,
2 the genetic diversity of a crop species and its wild relatives (van Hintum et al. 2000). Most core
3 collections developed include 5-20% of the accessions present in the collection, capturing 70-90% of
4 the diversity, and without redundant entries (van Hintum et al. 2000). Core collections do not replace
5 the whole collections from which these are obtained, however they can optimize the characterization
6 and evaluation efforts by focusing on a subset of accessions (van Hintum et al. 2000). Recognizing these
7 objectives, core collections were recommended by the global plan of action for the conservation and
8 sustainable utilization of plant genetic resources for food and agriculture as a necessary activity to
9 progress the use of genetic resources (FAO 1996).

10 The development of core collections is traditionally based on passport data or phenotypic traits
11 that are often unreliable and incomplete or influenced by environmental factors, respectively. DNA
12 markers, such as simple sequence repeats (SSRs) and single-nucleotide polymorphisms (SNPs), are the
13 tool of choice for the development of core collections. They can accurately represent the genetic
14 diversity of the entire collection and have no problems related to incomplete data and environmental
15 interactions, typically linked to passport and phenotypic markers. SSRs generally show a high level of
16 polymorphism than SNPs, leading to population-specific alleles that are useful for revealing population
17 structure. Nevertheless, SSRs are usually developed in small numbers for one species, and they may not
18 reflect the genome-wide genetic diversity respect to SNPs (Bernard et al. 2020). The latter are much
19 more frequent in the genomes and many SNPs can be identified using high-throughput genomics tools,
20 allowing to develop panels of markers useful for genetic diversity and fine mapping. Thus, SSR and
21 SNP markers bring different views of the population structure, and their other characteristics are
22 reported by Guichoux et al. (2011). In hazelnut, a total of 718 SSR markers have been developed and
23 more than 430 of them were used for the development of a reference linkage map (Mehlenbacher 2018;
24 Botta et al. 2019). They have been used to fingerprint cultivars, to identify duplicate accessions and
25 parents, to study genetic diversity in cultivated and wild populations, and in association mapping studies
26 (Mehlenbacher 2018; Botta et al. 2019). On the contrary, SNPs have been used in hazelnut only recently
27 and to develop two high-density genetic maps (Botta et al. 2019).

1 Different strategies and bioinformatic tools have been proposed to construct core collections
2 (Schoen and Brown 1993; Marita et al. 2000; Franco et al. 2005) and these have been compared in
3 annual (e.g., Franco et al. 2006) and perennial species (e.g., Escribano et al. 2008). Both studies
4 concluded that the maximization (M) strategy (Schoen and Brown 1993), which maximizes the number
5 of alleles, is highly suitable for constructing core collections. Several algorithms based on the M-strategy
6 have been developed and implemented in different software, such as MSTRAT (Gouesnard et al. 2001),
7 PowerCore (Kim et al. 2007), and Core Hunter (Thachuk et al. 2009; De Beukelaer et al. 2012).

8 Many studies concerning the construction of core sets have been performed in annual species.
9 Nevertheless, the benefits of developing core collections are perhaps most evident in woody perennial
10 species, usually maintained as clones in collection fields. These are due to higher management costs per
11 accession than those needed to maintain seed germplasm (Escribano et al. 2008). Development of core
12 collections using SSR markers has been performed in several fruit tree species, such as apple (Liang et
13 al. 2015; Lassois et al. 2016), apricot (Wang et al. 2011), carob tree (Di Guardo et al. 2019), chestnut
14 (Pereira-Lorenzo et al. 2017), fig (Balas et al. 2014), grape (Le Cunff et al. 2008; Štajner et al. 2014),
15 hazelnut (Öztürk et al. 2017), olive (Belaj et al. 2012; Díez et al. 2012; El Bakkali et al. 2013), pear
16 (Miranda et al. 2010; Liu et al. 2015), and walnut (Bernard et al. 2020).

17 The main objective of this work was to develop a core collection representative of the hazelnut
18 genetic diversity conserved in different *ex situ* and *in situ* European germplasm repositories. For that
19 purpose, first we used different M-strategy approaches to build and select the respective best subset
20 based on ten SSR markers. In a second step, the diversity parameters of each subset were compared to
21 select the final core collection. A quality evaluation of each sampling method was also performed
22 following the model proposed by Odong et al. (2013). Finally, the population structure and relatedness
23 among genotypes were also investigated.

24 25 **Material and methods**

26 **Plant material and microsatellite genotyping**

1 A total of 410 hazelnut accessions were collected from nine different *ex situ* germplasm repositories
2 located in six Countries: UK, Spain, France, Italy, Slovenia, and USA. Moreover, 6 landrace accessions
3 were also collected from an *on-farm* survey in the Nuoro province (Sardinia, Italy) (Online Resource 1).

4 Total genomic DNA was extracted from 0.20 g of young leaves or immature catkins using the
5 modified procedure of Thomas et al. (1993). A total of 10 SSR loci selected by Boccacci and Botta
6 (2010) were analysed: CaT-B107, CaT-B501, CaT-B502, CaT-B503, CaTB504, CaT-B505, CaT-B507,
7 CaT-B508 (Boccacci et al. 2005), CaC-B020, and CaC-B028 (Bassil et al. 2005). PCR amplifications
8 were performed in a volume of 15 µl containing 40 ng DNA, 0.5 U Taq-DNA polymerase (Bioline,
9 Meridian Bioscience, OH, USA), 3 µl 5x PCR buffer (Bioline, Meridian Bioscience), 2.2 mM MgCl₂,
10 200 µM dNTPs, and 0.5 µM of each primer. The PCR conditions were: a first denaturation step at 95
11 °C for 9 min, followed by 26 cycles of denaturation (30 s at 95 °C), annealing (45 s at 55 °C and 50 °C
12 for CaT-B502), and extension (90 s at 72 °C). The final elongation step was carried out at 72 °C for 30
13 min.

14 Total genomic DNA was extracted from the UK samples using the Nucleospin[®] Plant II kit
15 (Macherey-Nagel GmbH & Co. KG, Deuren, Germany) according to manufacturer's instructions. SSR
16 loci were the same as above, but loci were amplified in two multiplex reactions: i) MP1: CaT B107,
17 CaT B501, CaT B502, CaT B504, and CaC B028; ii) MP2: CaT B503, CaT B505, CaT B507, CaT
18 B508, and CaC B020. PCR amplifications were performed in a volume of 11 µl containing 10 ng DNA
19 and using the Type-it Microsatellite PCR kit (QIAGEN, MD, USA) according to the manufacturer's
20 protocol. The PCR conditions were: a first denaturation step at 95 °C for 5 min, followed by 35 cycles
21 of denaturation (30 s at 95 °C), annealing (45 s at 55 °C decreasing by 0.5 °C per cycle for the first 10
22 cycles), and extension (60 s at 72 °C). The final elongation step was carried out at 72 °C for 15 min.

23 Amplification products were analysed using an ABI-PRISM 3130 Genetic Analyzer capillary
24 electrophoresis instrument; UK samples were analysed using an ABI 3730xl capillary electrophoresis
25 instrument (both, Applied Biosystems, Foster City, CA, USA). Results were processed with
26 GeneMapper software (Applied Biosystems), and alleles were designated by their size in base pairs

1 using a GeneScan-500 LIZ standard (Applied Biosystems). UK data were aligned to the main dataset
2 by applying a simple conversion based on a series of overlapping cultivars between the two datasets.

3 Data obtained at the same SSR loci reported by Boccacci et al. (2013) for 17 reference cultivars
4 from the Hazelnut Research Institute (HRI) at Giresun (Turkey), 3 reference cultivars from the National
5 Agricultural Research Foundation - Pomology Institute (NAGREF-PI) at Naoussa (Greece), and 44
6 landraces surveyed *on-farm* in southern Europe were also added. Thus, microsatellite data from a total
7 of 480 accessions were processed using the software Identity 4.0 (Wagner and Sefc 1999) to calculate
8 the total probability of identity (Paetkau et al. 1995) and to identify samples with identical SSR
9 genotype. When two or more accessions had identical SSR genotype, only one was retained for further
10 analysis.

11

12 **Construction of the core collections by different M-strategies**

13 Three different approaches based on the maximization (M) strategy were used to build core collections:

- 14 • The standard M-strategy described by Schoen and Brown (1993) was employed as implemented in
15 MSTRAT (Gouesnard et al. 2001). Nei's diversity index (Nei 1987) was used as diversity criterion.
16 A total of five subsets with 10, 20, 30, 40 and 50 entries, respectively, were developed, together with
17 the optimized subset selected by the software algorithm. For each sampling size, 100 independent
18 replicates and 200 iterations were generated and the replicates that maintained the highest number of
19 alleles and genetic diversity scores were selected;
- 20 • The advanced M-strategy proposed by Kim et al. (2007) was carried out as implemented in
21 PowerCore v. 1.0;
- 22 • The advanced stochastic local search (SLS) algorithm, replica exchange Monte Carlo, developed by
23 Thachuk et al. (2009) and implemented in Core Hunter II (De Beukelaer et al. 2012). The software
24 can select core subsets using different allocation strategies by optimizing one genetic parameter or
25 many parameters simultaneously. By maximizing only genetic distance parameters the software
26 selects the most genetically distant accessions, whereas by optimizing diversity index accessions are
27 selected with the highest allelic variability. In this study, six allocation strategies were used: i)

1 optimizing each of the following measures independently: average Cavalli-Sforza and Edwards
2 (Dce) and Modified Rogers (Mr) as genetic distance parameter, expected proportion of heterozygous
3 loci (He) and Shannon's diversity index (Sh) as allelic diversity indices, and allele coverage (Cv); ii)
4 optimizing simultaneously both Cv and Dce (Cv-Dce) with equal weight assigned to each parameter.
5 Indeed, when a weight of 50% was assigned to Cv and 50% to Dce, all observed alleles were captured
6 in the sampled subset (Online Resource 2, Fig. S1). For each strategy, five subsets with 10, 20, 30,
7 40, and 50 entries, respectively, were developed.

8

9 **Characterization and validation of the representativeness of the core collections**

10 In order to evaluate the ability of each sampling strategy in capturing the diversity and representativeness
11 in the sampled subsets, as compared to the whole germplasm collection, different parameters were
12 considered: i) no significant differences in the number of alleles (A), genetic diversity (GD), observed
13 heterozygosity (Ho), and polymorphism information content (PIC), computed by the ANOVA analysis
14 with the SPSS software (IBM, Armonk, NY, USA). *Post hoc* Dunnett's test was used to compare the
15 means of diversity estimates from different core subsets with the entire collection used as control; ii)
16 number of loci with significantly different allele frequencies (Fr). Each locus was analysed
17 independently, comparing the frequency of each allele between the entire collection and each core subset
18 by the chi-squared test (Escribano et al. 2008). A, GD, Ho, PIC, and Fr were calculated using
19 PowerMarker v.3.25 (Liu and Muse 2005).

20 After selecting the best subset for each sampling strategy based on the above parameters, the
21 representativeness of the subsets was validated against the criteria proposed by Escribano et al. (2008):
22 i) capture all the alleles present in the original collection; ii) show no significant differences in frequency
23 distribution of alleles in at least 95% of the loci from that of the whole collection; iii) show no significant
24 differences in diversity indices, GD and Ho, between the core and the whole collection.

25

26 **Comparison and quality of sampling strategies**

1 Once the most representative core subsets were determined for each strategy, the effectiveness of each
2 sampling method was evaluated following the criteria reported by Thachuk et al. (2009) which expected
3 the best subset to have the: i) highest average genetic distance between accessions; ii) highest allele
4 richness; iii) lowest proportion of non-informative alleles and, equivalently, the highest allele coverage.
5 In order to assess this, Modified Rogers (MR) and Cavalli-Sforza and Edwards (CE) genetic distances,
6 Shannon's diversity index (SH), expected proportion of heterozygous loci (HE), number of effective
7 alleles (NE), proportion of non-informative alleles (PN), and allele coverage (CV) were calculated for
8 each core collection and compared with respect to the entire collection. Each parameter was optimized
9 independently by performing 20 runs using Core Hunter II (De Beukelaer et al. 2012).

10 The quality of each sampling method was also determined against two criteria proposed by Odong
11 et al. (2013):

- 12 • Average distance between each accession in the whole collection and the nearest entry in the core
13 collection (A-NE), a criterion to indicate the representativeness of a core collection. If the A-NE
14 realized value is low, there is always an entry close to each accession;
- 15 • Average distance between each entry in the core collection and the nearest neighbouring entry in the
16 core collection (E-NE), a criterion to indicate to what extent the entries are spaced in the diversity
17 space, represented by the whole collection. If the E-NE realized value is high, the entries cover the
18 entire space, and each entry should be as different as possible from each other.

19 In order to create a baseline to evaluate these criteria, 1,000 random subsets (rA-NE and rE-NE)
20 were generated of the desired sizes. The criteria were determined for each random set and the standard
21 deviation of the results per size were calculated. To indicate the potential value of the criteria (pA-NE
22 and pE-NE), an optimisation was done for both criteria for all core sizes. All calculations were
23 performed following the genetic distance optimisation (GDOpt) procedure described by Odong et al.
24 (2011).

25

26 **Genetic structure analysis**

1 The genetic structure analysis and the relatedness among genotypes of the final core collection
2 were performed within our whole hazelnut germplasm collection (WHGC), composed by true-to-type
3 genotypes obtained from the Identity analysis. The population structure was explored using
4 STRUCTURE v. 2.3.4 (Pritchard et al. 2000), a model-based Bayesian clustering method, assigning
5 individuals to subpopulations with no *a priori* grouping assumptions. The admixture model was applied,
6 and allele frequencies were assumed to be correlated. A burn-in period of 1,000,000 generations and
7 2,000,000 Markov chain Monte Carlo replications were used. STRUCTURE was run 10 independent
8 times for each K value ranging from 1 to 20. The most likely K value was determined using the ΔK
9 method (Evanno et al. 2005), as implemented in CLUMPAK (Kopelman et al. 2015). The resulting
10 matrices of estimated group membership coefficients (Q) were permuted using the *Greedy* algorithm
11 implemented in CLUMPP (Jakobsson and Rosenberg 2007) and bar plots were drawn using
12 STRUCTURE PLOT v 2.0 (Ramasamy et al. 2014). Genotypes with probability of membership $\geq 80\%$
13 ($Q \geq 0.8$) were assigned to the same group, while those with intermediate admixture coefficients ($Q <$
14 0.8) in any group were classified as “admixed” and were clustered in a separate mosaic group (M). The
15 genetic relationships among genotypes were also investigated using the weighted Neighbor-Joining (NJ)
16 dendrogram and the principal coordinate analysis (PCoA) implemented in DARwin v. 6.0 (Perrier and
17 Jacquemoud-Collet 2006). The NJ tree and the two-dimensional PCoA scatterplot were both constructed
18 based on Dice dissimilarity scores (10,000 bootstraps).

19

20 **Results**

21 **Sets of synonyms**

22 The genetic profiles of 480 accessions across 10 SSR loci were analysed using the Identity software to
23 identify duplicates, synonyms, and mistakes. Among them, 430 accessions are conserved in 11
24 international *ex situ* germplasm repositories, while 50 accessions are landraces maintained *on-farm* in
25 five southern European countries (Portugal, Spain, Italy, Slovenia, and Greece). Moreover, only 106
26 accessions were identified with a unique cultivar name, while the remaining 374 were groups of two or
27 more accessions labelled with the same name.

1 SSR analysis identified a total of 181 genotypes showing a unique profile, with a total probability
2 of identity of 1.85×10^{-12} (Online Resource 1). By comparison, 252 accessions (52.5% of the total)
3 were deemed to be duplicates and 47 accessions (9.7% of the total) were classified as planting or labeling
4 mistakes. Among the accessions deemed to be duplicates, a total of 18 synonym groups were identified
5 (Online Resource 1). Each set grouped accessions with similar nut and husk morphology and most of
6 them were already reported in the literature (Bocacci et al. 2006, 2008, 2013; Gökirmak et al. 2009;
7 Gürcan et al. 2010; Bacchetta et al. 2015). Nevertheless, some new synonyms were also identified: i)
8 the German accessions ‘Kurzhuellige Zellernuss’, ‘Minna's Zellernuss’, ‘Volle Zellernuss’, and
9 ‘Gunslebenert Zellernuss’ showed the same SSR profile of the cultivars ‘Barr's Zellernuss’, ‘Gustav's
10 Zellernuss’, ‘Merveille de Bollwiller’ (syn. ‘Hall’s Giant’), and ‘Gunslebert’, respectively; ii) the
11 English accession ‘Inghilterra’ was genetically identical to ‘Bandnuss’ (syn. ‘Apolda’); iii) the Spanish
12 accessions ‘Closca molla’ and ‘Punxenc’ revealed the same genetic profile as ‘Comun Alava’ and ‘Pere
13 Mas’, respectively; iv) the local cultivars ‘Negret primerenc’ and ‘Negret primerenc cort’ showed the
14 same microsatellite profile as ‘Negret’.

15

16 **Development of core collections and comparison to the whole collection**

17 A total of 181 true-to-type genotypes, representing our whole hazelnut germplasm collection (WHGC),
18 were used to construct core collections by means of three different approaches based on the M-strategy.
19 The performance of each sampling strategy for assembling core collections was evaluated over a range
20 of putative core subset (sample) sizes. Thus, a total of five subsets with 10 (5.5%), 20 (11.0%), 30
21 (16.6%), 40 (22.1%), and 50 (27.6%) entries, respectively, were developed, except for the PC strategy
22 where only one subset can be obtained.

23 The results of the variability parameters (A, GD, Ho, and PIC) obtained from a total of 37 subsets
24 compared with the initial collection are reported in Table 1. No SSR loci with significantly different
25 allele frequencies (Fr) were observed and thus, the criterion of no significant differences ($P < 0.05$) in
26 at least 95% of the loci was met in all subsets. The characterization of the WHGC showed 118
27 amplification fragments (A) with a mean GD, Ho and PIC of 0.79, 0.80 and 0.76, respectively. Among

1 the subsets obtained by the MSTRAT (MS) strategy, only MS10 showed a significant difference ($P <$
2 0.05) in the number of alleles (A). MS50 was deemed the best subset capturing all the alleles present in
3 the WHGC, while the optimized MS19 subset given by the algorithm captured a total of 99 alleles
4 (83.9%). The core collection obtained by the Power Core (PC) strategy, representing a full coverage of
5 all the alleles existing in the WHGC, comprised 53 entries (29.3%) and no significant differences were
6 observed with the entire collection. Among the subsets obtained with the Core Hunter (CH) single-
7 strategy, optimizing the Dce, Mr, Cv, He, and Sh indices independently, significantly higher values (P
8 < 0.05) were detected for the number of alleles (A) at Dce10, Dce20, Mr10, Mr20, Mr30, Cv10, He10,
9 and Sh10 subsets, for Ho in all Mr subsets, and for GD and PIC in all Dce, He and Sh subsets. Thus,
10 only by optimizing the Cv index was it possible to build a core collection that respected the criteria
11 proposed by Escribano et al. (2008) and the best subset was recovered with a minimum of 30 entries
12 (Cv30). The CH multi-strategy, where both Cv and Dce were optimizing simultaneously with equal
13 weight (50%), only Cv-Dce30 respected the criteria proposed by Escribano et al. (2008) and was
14 selected as the best subset.

15

16 **Selection of the final core collection**

17 A total of four core collections were selected from each sampling method: MS50 from MSTRAT (MS),
18 PC53 from Power Core (PC), Cv30 and Cv-Dce30 from Core Hunter (CH) single- and multi-strategy,
19 respectively.

20 In Table 2 all sampling strategies are compared with the whole collection (WHGC) and are listed
21 the mean values of the independents runs for each of the following parameters: the genetic distances
22 MR and CE, the genetic diversity indices SH, HE, and NE, and of the auxiliary values PN and CV
23 (Thachuk et al. 2009). All core subsets showed higher average genetic distance between entries and
24 higher allelic richness than the WHGC. Moreover, all sampling strategies were optimal in minimizing
25 PN and maximizing CV (0.0 and 100.0, respectively). Among them, the CH multi-strategy (Cv-Dce30)
26 showed slightly higher values at CE, SH, and HE and the highest NE value.

1 In Table 3 are reported the A-NE and E-NE quality parameters for each strategy: the realized
2 values (A-NE and E-NE), the potential optimal value (pA-NE and pE-NE), the average value from 1,000
3 random sets and the corresponding standard deviation (rA-NE and rE-NE). All four realized values for
4 the A-NE criterion were higher than their respective potential optimal values but were not different from
5 those of the random sets. Therefore, no strategies improved the A-NE criterion when compared to the
6 random sets. On the contrary, all four E-NE values were considerably higher than random values and
7 considerably less than potential values. In proportion, none of them reached more than 80% (PC53 and
8 Cv-Dce30) of the maximum achievable, while the random sets reached 67-69% of the maximum
9 achievable. Thus, all four strategies did improve the E-NE criterion, as compared to a random set, and
10 the CH multi-strategy (Cv-Dce30) outperformed the others.

11 Thus, the subset to form the final core collection of our whole collection was obtained by the CH
12 multi-strategy (Cv-Dce30) and was composed by 30 entries (16.6%); the relationship among frequencies
13 of alleles between this subset and the WHGC was very highly correlated ($R^2=0.93$) (Online Resource 2,
14 Fig. S2).

15

16 **Genetic population structure**

17 The estimation of ΔK (Online Resource 2, Fig. S3) from the analysis of 181 unique genotypes revealed
18 the highest value for $K = 3$ ($\Delta K = 246.82$), but high values were also obtained for $K = 2$ ($\Delta K = 205.22$)
19 and $K = 5$ ($\Delta K = 184.24$). In $K = 2$, genotypes were grouped in two gene pools (Fig. 1): one composed
20 mainly by cultivars from the Central Europe (CEU) and the British Islands (BI), and another composed
21 by cultivars from the Iberian Peninsula (IbeP), the Italian Peninsula (ItaP), and the Balkans-Black Sea
22 (BBS). A total of 63 cultivars were not clearly placed in these groups ($Q < 0.8$) and were classified as
23 admixed. In $K = 3$ were observed three groups composed mainly by cultivars from CEU and BI, ItaP,
24 and BBS, respectively (Fig. 1). Cultivars from the Iberian Peninsula (IbeP) were widespread within all
25 three groups, while 62 cultivars were classified as admixed. In $K = 5$ genotypes were classified into five
26 groups (Fig. 1). Cultivars from CEU and BI were placed in two separate groups: Q1 was composed by
27 9 cultivars from CEU and 2 accessions of unknown origin ('Mogulnuss' and 'Pallagrossa'), while Q2

1 included 8 cultivars from BI, 2 from CEU, and 3 accessions of unknown origin ('Apolda', 'Bearn', and
2 'Sodlinger'). 18 cultivars from IbeP showed the tendency to constitute a separate group (Q3), together
3 with 4 cultivars and 3 landraces from ItaP, 1 landrace from BBS, and 1 accession of unknown origin
4 ('Comen'). Q4 grouped 13 cultivars and 10 landraces from ItaP, 4 cultivars and 1 landrace from IbeP,
5 and 1 cultivar from CEU. Q5 clustered 20 cultivars from BBS, 5 landraces from ItaP, and 2 accessions
6 of unknown origin ('Fructo rubro' and 'Jann's'). A total of 74 genotypes were classified as admixed (Q
7 < 0.8) and were deemed "mosaics" (M group).

8 NJ dendrogram and PCoA scatterplot showed a clustering of the 181 genotypes similar to that
9 obtained from STRUCTURE analysis. In the NJ dendrogram (Fig. 2), genotypes were grouped in three
10 main clusters (I, II, and III), corresponding to $K = 3$, that showed a substructure similar to that observed
11 in $K = 5$: CEU (Q1) and BI (Q2) constituted two distinct subgroups into cluster I; IbeP (Q3) and ItaP
12 (Q4) were separated in several subgroups into cluster II; BBS (Q5) corresponded to the cluster III.
13 Admixed genotypes (M) were distributed in all three main clusters. In the PCoA scatterplot (Fig. 3), the
14 projection of the genotypes on a two-dimensional plane defined by the first two PCs (15.85 % of the
15 cumulative variation) showed: i) a separation between groups CEU and BI (right half of the graph) and
16 groups IbeP, ItaP, and BBS (left half of the graph), as in $K = 2$; ii) a separation between group ItaP (top
17 left), group BBS (lower left), and groups CEU and BI (right half), as in $K = 3$; iii) a general tendency to
18 separate each Q group obtained with the $K = 5$ stratification. CEU (Q1) was placed in the upper right of
19 the graph, while BI (Q2) was positioned in the right half. ItaP (Q4) was placed in the upper left, while
20 IbeP (Q3) and BBS (Q5) were located separately in the lower left. M genotypes were scattered in all
21 four parts of the graph.

22 Considering the WHGC population structure obtained, the genotypes included in the Cv-Dce30
23 core collection covered all six groups: 1 from Q1 (9.1%), 2 from Q2 (15.4%), 4 from Q3 (15%), 7 from
24 Q4 (24.1%), 5 from Q5 (18.5%), and 11 from M (15%).

25

26 **Discussion**

27 **Building the core collection**

1 Mislabeling and duplication are important challenges for germplasm conservation. In addition,
2 the existence of synonyms is a characteristic challenge in vegetatively propagated woody perennial
3 species. Thus, SSR markers have become very valuable tools in the management of *ex situ* and *in situ*
4 hazelnut collections. In our work, as duplicates were indicated the genotypes that showed the same SSR
5 profile, and they were 52.5 % of the total number of accessions analysed. Among them were reported
6 two types of duplicates: i) accessions labelled with the same name and collected from different collection
7 fields. In this first case was possible to define the true-to-type SSR genotype of most cultivars by
8 comparing the profiles obtained from several accessions or identify some mislabeling among them; and
9 ii) accessions labelled with a different name, conserved in different collection fields or in the same
10 collection field. In this second case, the results allowed us to identify some mislabeling (9.7% of the
11 total number of accessions analysed) due to planting or labeling mistakes, as happened to Bassil et al.
12 (2009) during a backup of the USDA collection. Moreover, it was also possible to confirm several sets
13 of synonyms reported in the literature (Boccacci et al. 2006, 2008, 2013; Gökirmak et al. 2009; Gürcan
14 et al. 2010; Bacchetta et al. 2015) and identify new ones. Among them, the local cultivars ‘Negret
15 primerenc’ and ‘Negret primerenc cort’ revealed the same genetic profile of ‘Negret’, although are more
16 productive and their fruits mature earlier (Rovira et al. 2017). They represented a possible case of clonal
17 mutation, and a similar result was observed between ‘Tonda di Biglini’ and ‘Tonda Gentile delle
18 Langhe’ by Valentini et al. (2014). Consequently, to construct our core collection it was important to
19 identify mislabeling, duplicates, and synonyms from the whole hazelnut germplasm collection
20 (WHGC), to delete a significant source of redundancy and build the core collection only from true-to-
21 type genotypes.

22 The main strategies used to construct core collections from molecular marker data can be
23 classified into two groups. The first methods are based on genetic distance, with or without stratified
24 sampling techniques, that cluster the accessions and then select entries from each cluster using different
25 allocation approaches (van Hintum et al. 2000; De Beukelaer et al. 2012). The second methods are based
26 on the M-strategy that construct cores with high allelic richness, maximizing the number of observed
27 alleles at each marker locus (Schoen and Brown 1993). M-methods are the only approaches that recover

1 all the alleles of the whole collection, including rare alleles, and keep the original allele frequencies at
2 each marker, favouring smaller subsets, reducing redundancy, and capturing most of the genetic
3 diversity (Marita et al. 2000; Gouesnard et al. 2001). In fruit and nut tree genera, such as *Annona*, *Ficus*
4 and *Castanea*, the M-strategy was the most efficient method to develop core collections, outperforming
5 other strategies (Escribano et al. 2008; Balas et al. 2014; Pereira-Lorenzo et al. 2017); and for this reason
6 was largely used by many other authors (Le Cunff et al. 2008; Miranda et al. 2010; Belaj et al. 2012;
7 Díez et al. 2012; El Bakkali et al. 2013; Štajner et al. 2014; Liang et al. 2015; Liu et al. 2015; Öztürk et
8 al. 2017; Bernard et al. 2020). Nevertheless, the choice of the most appropriate evaluation measures
9 depends on the purpose of the core collection. Methods based on the allele representativeness (genetic
10 distances) are preferred by plant breeders, while methods based on the allele richness, including rare
11 and localized alleles, interest taxonomists and geneticists (Marita et al. 2000).

12 The M-strategy, based on MSTRAT (Gouesnard et al. 2001), PowerCore (Kim et al. 2007), and
13 Core Hunter (Thachuk et al. 2009, De Beukelaer et al. 2012) algorithms, was used to develop our core
14 collections (Table 1). The best subsets obtained from each sampling strategies (MS50, PC53, Cv30, and
15 Cv-Dce30) captured all the alleles with the minimum number of accessions, without significant
16 differences in allele frequencies and preserving the parameter variations when compared to the WHGC
17 (Table 2). The core subsets obtained from the Core Hunter (CH) simple- and multi-strategies showed
18 the minimum number of entries (30 accessions) compared to those obtained from MSTRAT (50
19 accessions, MS50) and PowerCore (53 accessions, PC53). As reported by Thachuk et al. (2009), our
20 results confirmed that the CH strategy was able to select significantly smaller core subsets that retain all
21 unique alleles within a whole collection and a similar result was also observed in olive by Díez et al.
22 (2012).

23 The main advantage of the Core Hunter software is its ability to build core collections using
24 different allocation strategies by optimizing one parameter or many parameters simultaneously.
25 Generally, core subsets optimized using multiple criteria perform worse than these obtained using
26 individual measures (Thachuk et al. 2009; Díez et al. 2012). Nevertheless, our core subset obtained from
27 CH multi-strategy (Cv-Dce30) showed slightly higher values at MR, CE, SH, and HE and the highest

1 NE value, compared to that obtained from CH single-strategy (Cv30) (Table 2). Thus, CH multi-strategy
2 approach, optimizing Cv and Dce simultaneously with equal weight, was selected as the best strategy to
3 build our final core collection. It satisfied both the breeders' and geneticists'/taxonomists' perspectives
4 despite including only 16.6 % of the WHCG genotypes. This value sits within the 5–20 % proposed by
5 van Hintum et al. (2000) and is lower than the 19 % obtained from the Turkish national hazelnut
6 collection (Öztürk et al. 2017). Other studies on fruit tree crops reported a minimum requirement of 4
7 % inclusion in grape (Le Cunff et al. 2008; Štajner et al. 2014), a common range from 13 % in fig (Balas
8 et al. 2014) to 15-19 % in olive (Belaj et al. 2012; Díez et al. 2012; El Bakkali et al. 2013), and a
9 maximum of 28.6 % in apple (Liang et al. 2015) and 29.8 % in chestnut (Pereira-Lorenzo et al. 2017).

10 According to Ondong et al. (2013), the criterion of choice for evaluating the quality of core
11 collections should be determined by the objectives or type of the core collection. Thus, they proposed
12 two genetic distance-based criteria, A-NE and E-NE, for evaluating the quality of two important types
13 of core collections, respectively: i) a core collection (CC-I) where each entry represents one (itself) or
14 more accessions of the whole collection, in order to maximize the representativeness of genetic diversity
15 (A-NE); ii) a core collection (CC-X) where the diversity of the traits of the entries is maximized, in
16 order to represent the total genetic diversity (E-NE). Using these criteria, our results indicated that all
17 the best subsets obtained from each sampling strategy (MS50, PC53, Cv30, and Cv-Dce30) aimed at
18 covering the range of the genetic diversity, rather than representing the accessions from WHGC. All
19 values for the A-NE criterion were not significantly different from that of a random set. On the contrary,
20 all four strategies improved the E-NE criterion, as compared to a random set, and the CH multi-strategy
21 (Cv-Dce30) outperformed the others (Table 3). Since the objective of our final core collection was the
22 maximisation of the allelic richness, including rare and localized alleles, the E-NE criterion was the
23 most appropriate to evaluate the quality of the four core collections obtained. Nevertheless, the fact that
24 the size of the core collections varied made E-NE comparison difficult (Table 3). In the case of the Cv30
25 and Cv-Dce30 subsets (30 entries), the CH multi-strategy outperformed the CH single-strategy in all
26 aspects of the E-NE criterion, but not by a large factor (80% vs 77% of the maximum achievable).
27 Comparing these two subsets with the larger core collections MS50 (50 entries) and PC53 (53 entries),

1 all four were in a similar range. Nevertheless, PC53 did best in terms of standard deviation from the
2 random E-NE, whereas Cv-Dce30 did best in terms of approaching the potential maximum.

3

4 **Characteristics of the core collection**

5 The population structure and relatedness among SSR genotypes of 181 cultivars and landraces from
6 WHGC, indicated the existence of three levels of genetic structure (Fig 1, Fig. 2, and Fig. 3). In the first
7 level ($K = 2$) was observed a geographic pattern with one gene pool dominating the western and central
8 Europe (BI and CEU) and another gene pool dominating the southern Europe (IbeP, ItaP, and BBS). In
9 the second level ($K = 3$), there were still the gene pools frequent in the western and central Europe (BI
10 and CEU) and in the southern Europe (IbeP and ItaP), but there appeared a third gene pool (BBS) most
11 frequent in the Balkans and the Anatolian Peninsula. Finally, in the third level ($K = 5$) there was still the
12 BBS gene pool, while genotypes from northern Europe were further subdivided into the BI and CEU
13 gene pools and those from southern Europe into the IbeP and ItaP gene pools.

14 The high level of genetic similarity between cultivars grown in the Iberian and Italian Peninsula,
15 observed in $K = 2$ and $K = 3$ (Fig. 1), was already reported by other authors (Bocacci et al. 2006;
16 Gökirmak et al. 2009; Gürcan et al. 2010) and was a consequence of a high gene flow between western
17 and central Mediterranean basin (Bocacci and Botta 2010). Nevertheless, in $K = 5$ (Fig. 1) cultivars
18 from the Iberian Peninsula were separated from Italian ones (Bocacci et al. 2013) and a significant
19 genetic differentiation between the Spanish and Italian gene pools was reported in subsequent studies
20 (Bocacci and Botta 2010; Bocacci et al. 2013), indicating that northern Spain and southern Italy were
21 two independent hazelnut domestication areas (Bocacci et al. 2013). On the contrary, the genetic
22 similarity between cultivars from the British Islands and the Central Europe obtained with $K = 2$ and K
23 $= 3$ (Fig. 1) was not observed by Gökirmak et al. (2009). Indeed, the authors reported that these cultivars
24 clustered into separate groups as observed in our $K = 5$ stratification (Fig. 1). Thus, considering the data
25 reported in the literature, the most likely genetic structure of WHGC was composed by five Q groups
26 (CEU, Q1; BI, Q2; IbeP, Q3; ItaP, Q4; and BBS, Q5), and a more complex group of mosaics (M).
27 According to several authors (Gökirmak et al. 2009; Bocacci and Botta 2010; Bocacci et al. 2013)

1 these gene pools would be the result of an independent domestication of *C. avellana* that occurred in
2 different geographical areas: Central Europe, the British Islands, Spain, Italy, and Black Sea. In contrast,
3 the mosaic genotypes, which are found throughout our sampling range, represent a heterogeneous group
4 that may be indicative of recent admixture between distinct groups of cultivars.

5 The Bayesian clustering and admixture analysis can be considered a standard method to identify
6 the ancestral populations from which cultivars originated and quantify genetic relationships with
7 probabilities and proportions. Thus, it was helpful for suggesting the unknown origin of some cultivars.
8 ‘Mogulnuss’ (syn. ‘Riekchen’s Zellernuss’) and ‘Pallagrossa’ were placed into the CEU group, ‘Comen’
9 in the IbeP group and ‘White Filbert’ (syn. ‘Fructo rubro’) in the BI group, confirming the results
10 obtained by Gökirmak et al. (2009). On the contrary, ‘Jann’s/ Jean’s’ was placed in the group mainly
11 from the Italian Peninsula rather than from the Black Sea and ‘The Shah’ was placed in the admixed
12 group instead of the English group 2, as would have been expected from the findings of Gökirmak et al.
13 (2009). In our analysis ‘Sodlinger’ clustered into the BI group, although it was placed in a Spanish–
14 Italian group by Muehlbauer et al. (2014).

15 The 30 entries (Online Resource 1) included in our final core collection (Cv-Dce30) were from
16 different countries: Italy (17 entries, 56.7%), Spain (4 entries, 13.3%), Germany (4 entries, 13.3%),
17 Turkey (3 entries, 10%), and Slovenia (2 entries, 6.7%). They covered all six genetic groups obtained,
18 further confirming that Cv-Dce30 core collection was representative of the WHGC. Interestingly, half
19 of the entries were true-to-type cultivars from *ex situ* collections, while the other half were landraces
20 from *in situ* collections. The high number of landraces included in the Cv-Dce30 core collection
21 indicated that the hazelnut *on-farm* exploration conducted in southern Europe (Boccacci et al. 2013) has
22 genuinely contributed to expanding the existent hazelnut biodiversity in our collections. No reference
23 accessions were included as “kernel” in our core collections (van Hintum et al. 2000), but the most
24 popular hazelnut cultivar ‘Tonda Gentile delle Langhe’ (TGL), particularly appreciated by the industry
25 for the morphological, organoleptic, and nutritional characteristics of its nuts and kernels, was included
26 in all MS50, PC53, Cv30, and Cv-Dce30 subsets. Different reference cultivars could be added on a case-
27 by-case basis in different places where this core collection could be studied, such as: ‘Negret’ and

1 'Casina' in Spain, 'Barcelona' (syn. 'Fertile de Coutard') in France, 'Tonda Gentile Romana' and
2 'Tonda di Giffoni' in Italy, and 'Tombul' in Turkey.

3

4 **Conclusions**

5 The M-strategies employed in this work to build our core collections may be considered useful tools for
6 the conservation and characterization of hazelnut genetic resources. Among them, the CH multi-
7 strategy, optimizing Cv and Dce simultaneously with equal weight, was selected as the best strategy to
8 build our final core collection. The ability of each sampling strategy in capturing the diversity and
9 representativeness, and the effectiveness and quality of each sampling method, were performed using
10 various well-known approaches. Thus, our final core collection, representing most of the diversity
11 conserved in the European hazelnut germplasm repositories, could be used as a base for new research
12 into genotype x environment interactions focused on a minimum number of accessions. However,
13 reducing the number of selected accessions inevitably increases the probability of discarding genotypes
14 with agronomical traits of interest. Thus, it will be important to consider core collections combining
15 molecular markers, morphological and phenotypical traits, as well as resistance to biotic and abiotic
16 stresses. It is also pragmatic to include cultivars that are considered more influential in a determined
17 cultivation area. Finally, any approach toward core collections should remain dynamic and be revised
18 periodically, to include new accessions and information about new characterization methods (e.g.,
19 functional markers), as well as new methodologies aimed at increasing their representativeness.

20

21

22 **Declarations**

23 **Acknowledgments** In memory of my dad, Ugo Boccacci (October 6, 1946 - March 19, 2021)

24 **Funding** This work was funded by AGRI GEN RES Community Program (European Commission, Directorate-
25 General for Agriculture and Rural Development, under Council Regulation (EC) No. 870/2004) – SAFENUT
26 project ("Safeguard of almond and hazelnut genetic resources: from traditional uses to modern agro-industrial
27 opportunities"), AGRI GEN RES 068

1 **Conflicts of interest/Competing interests** Authors declare that they have no conflicts of interest/competing
2 interests

3 **Ethics approval** Not applicable

4 **Consent to participate** Not applicable

5 **Consent for publication** Not applicable

6 **Data archiving statement/Availability of data and material** The Online Material 1 (EMS1, “xlsx” format) and
7 Online Material 2 (EMS2, “pdf” format) are available in the Dryad Repository (<https://datadryad.org/stash>) as
8 "Hazelnut SSR database: genetic profiles of the accessions, list of synonyms, and true-to-type genotypes"
9 (doi:10.5061/dryad.cz8w9gj45). For private access during the review period, reviewers may share the unpublished
10 dataset using this temporary link:
11 https://datadryad.org/stash/share/fG-4pHenR4hQ_z5G7y18b43ytHnTLpy3w-F2lFJqljk

12 **Code availability** Not applicable

13 **Authors' contributions** Paolo Boccacci conceived the study and written the manuscript. Paolo Boccacci, Maria
14 Aramini, Daniela Torello Marinoni, Nadia Valentini, Mercè Rovira, Anita Solar, and Jean-Paul Sarraquigne
15 collected vegetal material from the collection fields. Paolo Boccacci, Maria Aramini, Matthew Ordidge, and
16 Daniela Torello Marinoni performed SSR analyses. Paolo Boccacci and Theo van Hintum performed data
17 elaborations. Loretta Bacchetta coordinated and Roberto Botta co-coordinated the SAFENUT project. All authors
18 reviewed the manuscript. All authors read and approved the final version of the manuscript

19

20 **References**

21 Bacchetta L, Rovira M, Tronci C, Aramini M, Drogoudi P, Silva AP, Solar A, Avanzato D, Botta R, Valentini N,
22 Boccacci P (2015) A multidisciplinary approach to enhance the conservation and use of hazelnut *Corylus*
23 *avellana* L. genetic resources. *Genet Resour Crop Evol* 62:649–663

24 Balas FC, Osuna MD, Domínguez G, Pérez-Gragera F, López-Corrales M (2014). *Ex situ* conservation of
25 underutilised fruit tree species: establishment of a core collection for *Ficus carica* L. using microsatellite
26 markers (SSRs). *Tree Genet Genomes* 10:703–710

27 Bassil NV, Botta R, Mehlenbacher SA (2005) Microsatellite markers in the hazelnut: isolation, characterization,
28 and cross-species amplification in *Corylus*. *J Am Soc Hort Sci* 130:543-549

29 Bassil NV, Postman J, Hummer K, Botu M, Sezer A (2009) SSR fingerprinting panel verifies identities of clones
30 in backup hazelnut collection of USDA genebank. *Acta Hort* 845: 95-102

- 1 Belaj A, del Carmen Dominguez-García M, Atienza SG, Martín Urdiroz N, De la Rosa R, Satovic Z, Martín A,
2 Kilian A, Trujillo I, Valpuesta V, Del Río C (2012). Developing a core collection of olive (*Olea europaea* L.)
3 based on molecular markers (DArTs, SSRs, SNPs) and agronomic traits. *Tree Genet Genomes* 16:76
- 4 Bernard A, Barreneche T, Donkpegan A, Lheureux F, Dirlewanger E (2020). Comparison of structure analyses
5 and core collections for the management of walnut genetic resources. *Tree Genet Genomes* 8:365–378
- 6 Boccacci P, Botta R (2010) Microsatellite variability and genetic structure in hazelnut (*Corylus avellana* L.)
7 cultivars from different growing regions. *Sci Hortic* 124:128-133
- 8 Boccacci P, Akkak A, Bassil NV, Mehlenbacher SA, Botta R (2005) Characterization and evaluation of
9 microsatellite loci in European hazelnut (*Corylus avellana* L.) and their transferability to other *Corylus* species.
10 *Mol Ecol Notes* 5:934-937
- 11 Boccacci P, Akkak A, Botta R (2006) DNA-typing and genetic relationships among European hazelnut (*Corylus*
12 *avellana* L.) cultivars using microsatellite markers. *Genome* 49:598-611
- 13 Boccacci P, Rovira M, Botta R (2008) Genetic diversity of hazelnut (*Corylus avellana* L.) germplasm in
14 northeastern Spain. *HortScience* 43:667-672
- 15 Boccacci P, Aramini M, Valentini N, Bacchetta L, Rovira M, Drogoudi P, Silva AP, Solar A, Calizzano F,
16 Erdorğan V, Cristofori V, Ciarmiello LF, Contessa C, Ferreira JJ, Marra FP, Botta R (2013) Molecular and
17 morphological diversity of *on-farm* hazelnut (*Corylus avellana* L.) landraces from southern Europe and their
18 role in the origin and diffusion of cultivated germplasm. *Tree Genet Genomes* 9:1465–1480
- 19 Botta R, Molnar TJ, Erdorğan V, Valentini N, Torello Marinoni D, Mehlenbacher S (2019). Hazelnut (*Corylus*
20 spp.) Breeding. In: Al-Khayri JM, Jain SM, Johnson DV (eds.) *Advances in Plant Breeding Strategies: Nut*
21 *and Beverage crops*. Springer Nature, Switzerland, Volume 4, pp 157-219
- 22 De Beukelaer H, Smýkal P, Davenport GF, Fack V (2012) Core Hunter II: fast core subset selection based on
23 multiple genetic diversity measures using Mixed Replica search. *BMC Bioinformatics* 13:312
- 24 Di Guardo M, Scollo F, Ninot A, Rovira M, Hermoso JF, Distefano G, La Malfa S, Batlle I (2019) Genetic structure
25 analysis and selection of a core collection for carob tree germplasm conservation and management. *Tree Genet*
26 *Genomes* 15: 41
- 27 Díez CM, Imperato A, Rallo L, Barranco D, Trujillo I (2012). Worldwide core collection of olive cultivars based
28 on simple sequence repeat and morphological markers. *Crop Sci* 52:211-221

- 1 El Bakkali A, Haouane H, Moukhli A, Costes E, Van Damme P, Khadari B (2013) Construction of core collections
2 suitable for association mapping to optimize use of Mediterranean olive (*Olea europaea* L.) genetic resources.
3 PLoS ONE 8:e61265
- 4 Escribano P, Viruel MA, Hormaza JI (2008) Comparison of different methods to construct a core germplasm
5 collection in woody perennial species with simple sequence repeat markers. A case study in cherimoya
6 (*Annona cherimola*, Annonaceae), an underutilised subtropical fruit tree species. Ann Appl Biol 153:25–32
- 7 Evanno G, Regnaut S, Goudet J (2005) Detecting the number of clusters of individuals using the software
8 STRUCTURE: a simulation study. Mol Ecol 14:2611-2620
- 9 FAO (1996) Global plan of action for the conservation and sustainable utilization of plant genetic resources for
10 food and agriculture. Food and Agriculture Organization, Rome
- 11 FAOSTAT (2021) <http://www.fao.org/faostat/en/?#data>. Accessed 06 May 2021
- 12 Franco J, Crossa J, Taba S, Shands H (2005) A sampling strategy for conserving genetic diversity when forming
13 core subsets. Crop Sci 45:1035–1044
- 14 Franco J, Crossa J, Warburton ML, Taba S (2006) Sampling strategies for conserving maize diversity when
15 forming core subsets using genetic markers. Crop Sci 46:854–864
- 16 Ferreira JJ, Garcia-González C, Tous J, Rovira M (2010) Genetic diversity revealed by morphological traits and
17 ISSR markers in hazelnut germplasm from northern Spain. Plant Breed 129:435–441
- 18 Gökirmak T, Mehlenbacher SA, Bassil NV (2009) Characterization of European hazelnut (*Corylus avellana* L.)
19 cultivars using SSR markers. Genet Resour Crop Evol 56:147-172
- 20 Gouesnard B, Bataillon TM, Decoux G, Rozale C, Schoen DJ, David JL (2001) MSTRAT: an algorithm for
21 building germplasm core collections by maximizing allelic or phenotypic richness. J Hered 92:93–94
- 22 Guichoux E, Lagache L, Wagner S et al (2011) Current trends in microsatellite genotyping. Mol Ecol Resour
23 11:591–611
- 24 Gürcan K, Mehlenbacher SA, Erdoğan V (2010) Genetic diversity in hazelnut (*Corylus avellana* L.) cultivars from
25 Black Sea countries assessed using SSR markers. Plant Breed 129:422–434
- 26 Hummer KE (2001) Hazelnut genetic resources at the Corvallis repository. Acta Hort 556:21–24
- 27 Jakobsson M, Rosenberg NA (2007) CLUMPP: a cluster matching and permutation program for dealing with label
28 switching and multimodality in analysis of population structure. Bioinformatics 23:1801–1806

1 Kim KW, Chung HK, Cho GT, Ma KH, Gwag CD, Kim TS, Cho EG, Park YJ (2007) PowerCore: a program
2 applying the advanced M strategy with a heuristic search for establishing core sets. *Bioinformatics* 23:2155–
3 2162

4 Köksal AI (2000) Inventory of hazelnut research, germplasm and references. REU technical series. FAO-CIHEAM

5 Kopelman NM, Mayzel J, Jakobsson M, Rosenberg NA, Mayrose I (2015) Clumpak: a program for identifying
6 clustering modes and packaging population structure inferences across K. *Mol Ecol Resour* 15: 1179-1191

7 Lassois L, Denancé C, Ravon E, Guyader A, Guisnel R, Hibrand-Saint-Oyant L, Poncet C, Lasserre-Zuber P,
8 Feugey L, Durel CE (2016). Genetic diversity, population structure, parentage analysis, and construction of
9 core collections in the French apple germplasm based on SSR markers. *Plant Mol Biol Rep* 34:827–844

10 Le Cunff L, Fournier-Level A, Laucou V, Vezzulli S, Lacombe T, Adam-Blondon AF, Boursiquot JM, This P
11 (2008). Construction of nested genetic core collections to optimize the exploitation of natural diversity in *Vitis*
12 *vinifera* L. subsp. *sativa*. *BMC Plant Biol* 8:31

13 Liang W, Dondini L, De Franceschi P, Paris R, Sansavini S, Tartarini S (2015) Genetic diversity, population
14 structure and construction of a core collection of apple cultivars from Italian germplasm. *Plant Mol Biol Rep*
15 33:458–473

16 Liu KJ, Muse SV (2005). PowerMarker: An integrated analysis environment for genetic marker analysis.
17 *Bioinformatics* 21:2128–2129

18 Liu Q, Song Y, Liu L, Zhang M, Sun J, Zhang S, Wu J (2015) Genetic diversity and population structure of pear
19 (*Pyrus* spp.) collections revealed by a set of core genome-wide SSR markers. *Tree Genet Genomes* 11:128

20 Marita JM, Rodriguez JM, Nienhuis J (2000) Development of an algorithm identifying maximally diverse core
21 collections. *Genet Resour Crop Evol* 47:515–526

22 Mehlenbacher SA (2018) Advances in genetic improvement of hazelnut. *Acta Hort* 1226:1-12

23 Miranda C, Urrestarazu J, Santesteban LG, Royo JB, Urbina V (2010) Genetic diversity and structure in a
24 collection of ancient Spanish pear cultivars assessed by microsatellite markers. *J Am Soc Hortic Sci* 135:428-
25 437

26 Muehlbauer MF, Honig JA, Capik JM, Vaiciunas JN, Molnar TJ (2014) Characterization of eastern filbert blight-
27 resistant hazelnut germplasm using microsatellite markers. *J Am Soc Hortic Sci* 139:399–432

28 Nei M (1987) *Molecular evolutionary genetics*. Columbia Univ. Press, New York, NY

29 Odong TL, van Heerwaarden J, Jansen J, van Hintum TJJ, van Eeuwijk FA (2011) Statistical techniques for
30 defining reference sets of accessions and microsatellite markers. *Crop Sci* 51(6):2401–2411

- 1 Odong TL, Jansen J, van Eeuwijk FA, van Hintum TJJ (2013) Quality of core collections for effective utilisation
2 of genetic resources review, discussion and interpretation. *Theor Appl Genet* 126:289–305
- 3 Öztürk SC, Balık Hİ, Balık SK, Kızılcı G, Duyar Ö, Doğanlar S, Frary A (2017) Molecular genetic diversity of
4 the Turkish national hazelnut collection and selection of a core set. *Tree Genet Genomes* 13:113
- 5 Paetkau D, Calvert W, Stirling I, Strobeck C (1995) Microsatellite analysis of population structure in Canadian
6 polar bears. *Mol Ecol* 4:347–354
- 7 Pereira-Lorenzo S, Ramos-Cabrer AM, Barreneche T, Mattioni C, Villani F, Díaz-Hernández MB, Martín LM,
8 Martín A (2017) Database of European chestnut cultivars and definition of a core collection using simple
9 sequence repeats. *Tree Genet Genomes* 13:114
- 10 Perrier X, Jacquemoud-Collet J (2006) DARwin software. Available from: <http://darwin.cirad.fr/>. Accessed 5 May
11 2021
- 12 Pritchard JK, Stephens M, Donnelly P (2000) Inference of population structure using multilocus genotype data.
13 *Genetics* 155:945-959
- 14 Ramasamy RK, Sumathy Ramasamy S, Bindroo BB, Naik VG (2014) STRUCTURE PLOT: a program for
15 drawing elegant STRUCTURE bar plots in user friendly interface. *SpringerPlus* 3:431
- 16 Rovira M, Hermoso JF, Romero AJ (2017) Performance of hazelnut cultivars from Oregon, Italy, and Spain, in
17 Northeastern Spain. *HortTechnology* 27(5):631-638
- 18 Schoen DJ, Brown AHD (1993) Conservation of allelic richness in wild crop relatives is aided by assessment of
19 genetic markers. *Proc Natl Acad Sci USA* 90:10623–10627
- 20 Štajner N, Tomić L, Ivanišević D, Korać N, Cvetković-Jovanović T, Beleski K, Angelova E, Maraš V, Javornik
21 B (2014). Microsatellite inferred genetic diversity and structure of Western Balkan grapevines (*Vitis vinifera*
22 L.). *Tree Genet Genomes* 10:127–140
- 23 Thachuk C, Crossa J, Franco J, Dreisigacker S, Warburton M, Davenport GF (2009) Core Hunter: an algorithm
24 for sampling genetic resources based on multiple genetic measures. *BMC Bioinformatics* 10:243
- 25 Thomas MR, Matsumoto S, Cain P, Scott NS (1993) Repetitive DNA of grapevine: classes present and sequences
26 suitable for cultivar identification. *Theor Appl Genet* 86:173-180
- 27 Thompson MM, Lagerstedt HB, Mehlenbacher SA (1996) Hazelnuts. In: Janick J, Moore JN (eds) *Fruit breeding:*
28 *nuts*, vol 3. Wiley, New York, pp 125–184
- 29 Valentini N, Calizzano F, Boccacci P, Botta R (2014) Investigation on clonal variants within the hazelnut (*Corylus*
30 *avellana* L.) cultivar ‘Tonda Gentile delle Langhe’. *Sci Hortic* 165:303-310

- 1 van Hintum, T.J.L., Brown A.H.D., Spillane C., Hodgkin T. (2000) Core collections of plant genetic resources. IPGRI
- 2 Technical Bulletin No. 3. International Plant Genetic Resources Institute, Rome
- 3 Wagner H.W., Sefc K.M. (1999) IDENTITY 4.0. Centre for Applied Genetics, University Agricultural Sciences,
- 4 Vienna
- 5 Wang Y., Zhang J., Sun H., Ning N., Yang L. (2011) Construction and evaluation of a primary core collection of
- 6 apricot germplasm in China. *Sci Hort* 128:311–319
- 7

1 **Table 1** Variability parameters for different core subsets compared with the whole collection. In bold are indicated
 2 the best core subset obtained from each sampling strategy

Sampling strategy	Subset code	Subset size	<i>A</i>	<i>GD</i>	<i>Ho</i>	<i>PIC</i>
Whole collection	WHGC	181	118	0,79	0,80	0,76
MSTRAT	MS10	10	86 ^a	0,82	0,76	0,80
	MS20	20	100	0,83	0,81	0,81
	MS30	30	108	0,82	0,80	0,80
	MS40	40	109	0,82	0,82	0,80
	MS50	50	118	0,81	0,81	0,79
	MS19	19	99	0,83	0,81	0,81
Power Core	PC53	53	118	0,81	0,81	0,79
Core Hunter single - Dce	Dce10	10	86 ^a	0,85 ^a	0,77	0,83 ^a
	Dce20	20	93 ^a	0,85 ^a	0,78	0,83 ^a
	Dce30	30	101	0,84 ^a	0,77	0,83 ^a
	Dce40	40	103	0,84 ^a	0,78	0,82 ^a
	Dce50	50	103	0,84 ^a	0,77	0,82 ^a
Core Hunter single - Mr	Mr10	10	77 ^a	0,82	0,63 ^a	0,79
	Mr20	20	87 ^a	0,82	0,65 ^a	0,79
	Mr30	30	94 ^a	0,82	0,66 ^a	0,79
	Mr40	40	98	0,82	0,69 ^a	0,80
	Mr50	50	103	0,82	0,70 ^a	0,80
Core Hunter single - Cv	Cv10	10	94 ^a	0,83	0,88	0,81
	Cv20	20	109	0,80	0,83	0,77
	Cv30	30	118	0,81	0,85	0,79
	Cv40	40	118	0,78	0,82	0,76
	Cv50	50	118	0,79	0,81	0,76
Core Hunter single - He	He10	10	87 ^a	0,85 ^a	0,87	0,83 ^a
	He20	20	97	0,85 ^a	0,84	0,84 ^a
	He30	30	102	0,85 ^a	0,88	0,83 ^a
	He40	40	104	0,85 ^a	0,86	0,83 ^a
	He50	50	107	0,84 ^a	0,84	0,83 ^a
Core Hunter single - Sh	Sh10	10	89 ^a	0,85 ^a	0,86	0,83 ^a
	Sh20	20	103	0,85 ^a	0,85	0,83 ^a
	Sh30	30	107	0,85 ^a	0,86	0,83 ^a
	Sh40	40	109	0,85 ^a	0,85	0,83 ^a
	Sh50	50	112	0,84 ^a	0,84	0,82 ^a
Core Hunter multi - Cv-Dce	Cv-Dce10	10	94 ^a	0,83 ^a	0,87	0,82 ^a
	Cv-Dce20	20	108	0,84 ^a	0,83	0,82 ^a
	Cv-Dce30	30	118	0,82	0,83	0,80
	Cv-Dce40	40	118	0,83 ^a	0,82	0,81 ^a
	Cv-Dce50	50	118	0,83 ^a	0,79	0,81 ^a

3 *A*, number of alleles; *GD*, genetic diversity; *Ho*, observed heterozygosity; *PIC*, polymorphism information content
 4 ^aStatistically significant difference, Dunnett's test ($P < 0.05$)

1 **Table 2** Comparison of the best core subsets selected by each sampling method

Sampling strategy	Subset code	Subset size	<i>MR</i>	<i>CE</i>	<i>SH</i>	<i>HE</i>	<i>NE</i>	<i>PN</i>	Cv (%)
Whole collection	WHGC	181	0.62	0.80	4.15	0.79	4.94	0.00	118 (100)
MSTRAT	MS50	50	0.64	0.82	4.24	0.81	5.48	0.00	118 (100)
Power Core	PC53	53	0.64	0.83	4.26	0.81	5.58	0.00	118 (100)
Core Hunter single	Cv30	30	0.63	0.82	4.26	0.81	5.50	0.00	118 (100)
Core Hunter multi	Cv-Dce30	30	0.64	0.84	4.29	0.82	5.73	0.00	118 (100)

2 *MR*, Modified Rogers distance; *CE*, Cavalli-Sforza and Edwards distance; *SH*, Shannon's diversity index;
3 *HE*, expected proportion of heterozygous loci; *NE*, number of effective alleles; *PN*, proportion of non-
4 informative alleles; *CV*, allele coverage

5
6 **Table 3** Quality evaluation of each sampling method based on the average distance between each accession and the
7 nearest entry (*A-NE*) and average distance between each entry and the nearest neighbouring entry (*E-NE*)

Sampling strategy	Subset code	Subset size	<i>A-NE</i>	<i>pA-NE</i>	<i>rA-NE</i>	<i>std dev</i>	<i>E-NE</i>	<i>pE-NE</i>	<i>rE-NE</i>	<i>std dev</i>
Whole collection	WHGC	181	0.000	0.000	0.000	0.000	0.306	0.306	0.306	0.000
MSTRAT	MS50	50	0.273	0.232	0.279	0.006	0.433	0.563	0.387	0.018
Power Core	PC53	53	0.268	0.225	0.270	0.006	0.444	0.555	0.383	0.017
Core Hunter single	Cv30	30	0.356	0.290	0.348	0.008	0.480	0.627	0.420	0.024
Core Hunter multi	Cv-Dce30	30	0.354	0.290	0.348	0.008	0.499	0.627	0.420	0.024

8 *A-NE* and *E-NE*, realized values; *pA-NE* and *pE-NE*, potential optimal values; *rA-NE* and *rE-NE*, average values from
9 1,000 random sets and the corresponding standard deviation (*std dev*)

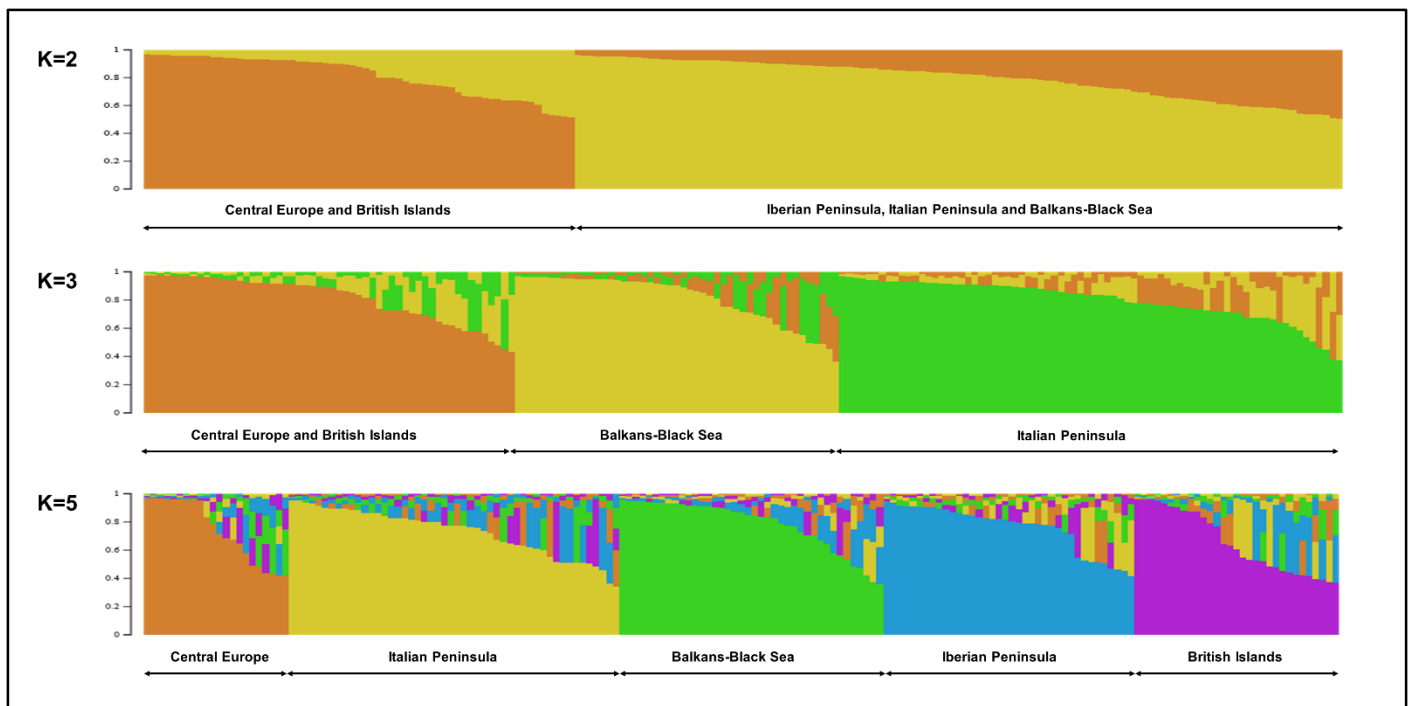


Fig. 1 Population structure and hierarchical organization of genetic relatedness of 181 genotypes from the whole hazelnut germplasm collection (WHGC) at $K = 2$, $K = 3$, and $K = 5$ as inferred by STRUCTURE software

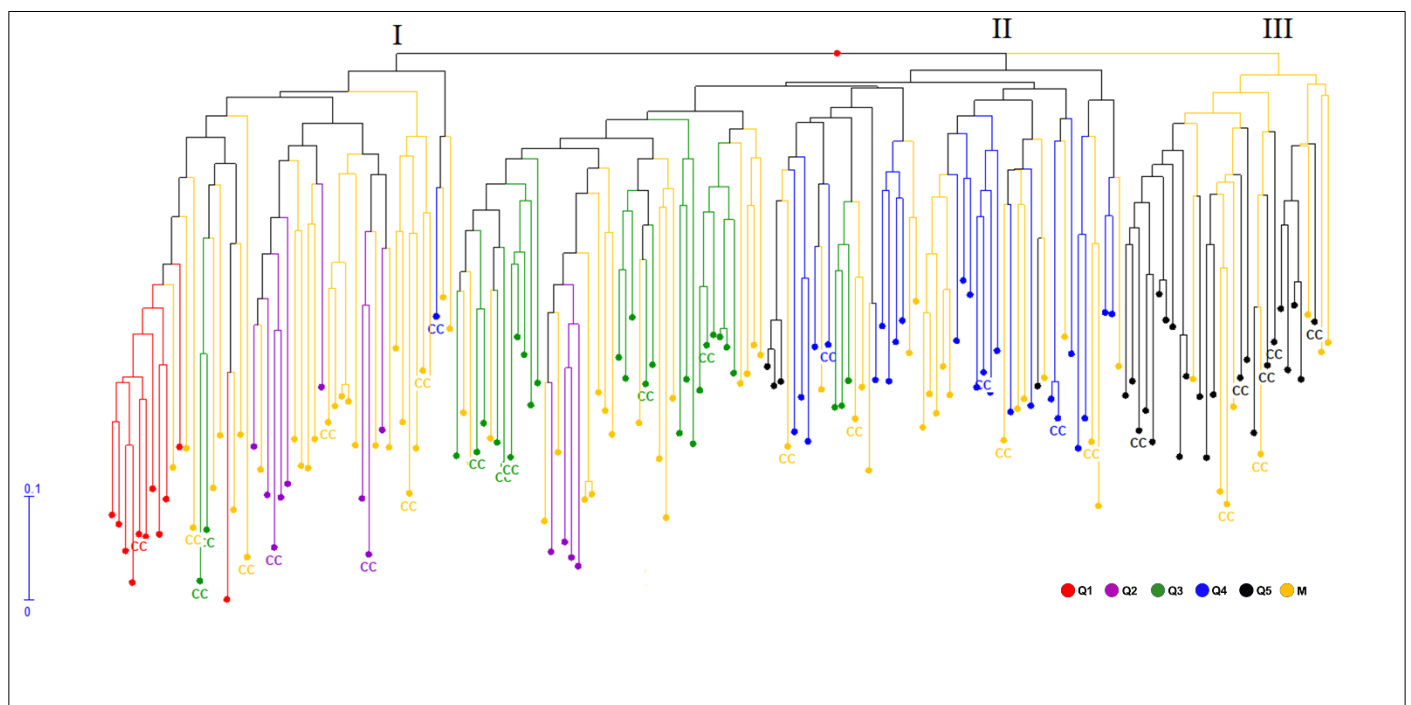


Fig. 2 Neighbor-joining dendrogram based on the Dice similarity index showing the relationships among 181 hazelnut genotypes from WHGC. Genotypes are colored according to their assignment to the different gene pools, as inferred by STRUCTURE software at $K = 5$: Central Europe (Q1), British Islands (Q2), Iberian Peninsula (Q3), Italian Peninsula (Q4), Balkans-Black Sea (Q5), and mosaic group (M). Entries of the final core collection (Cv-Dce30) are reported as CC

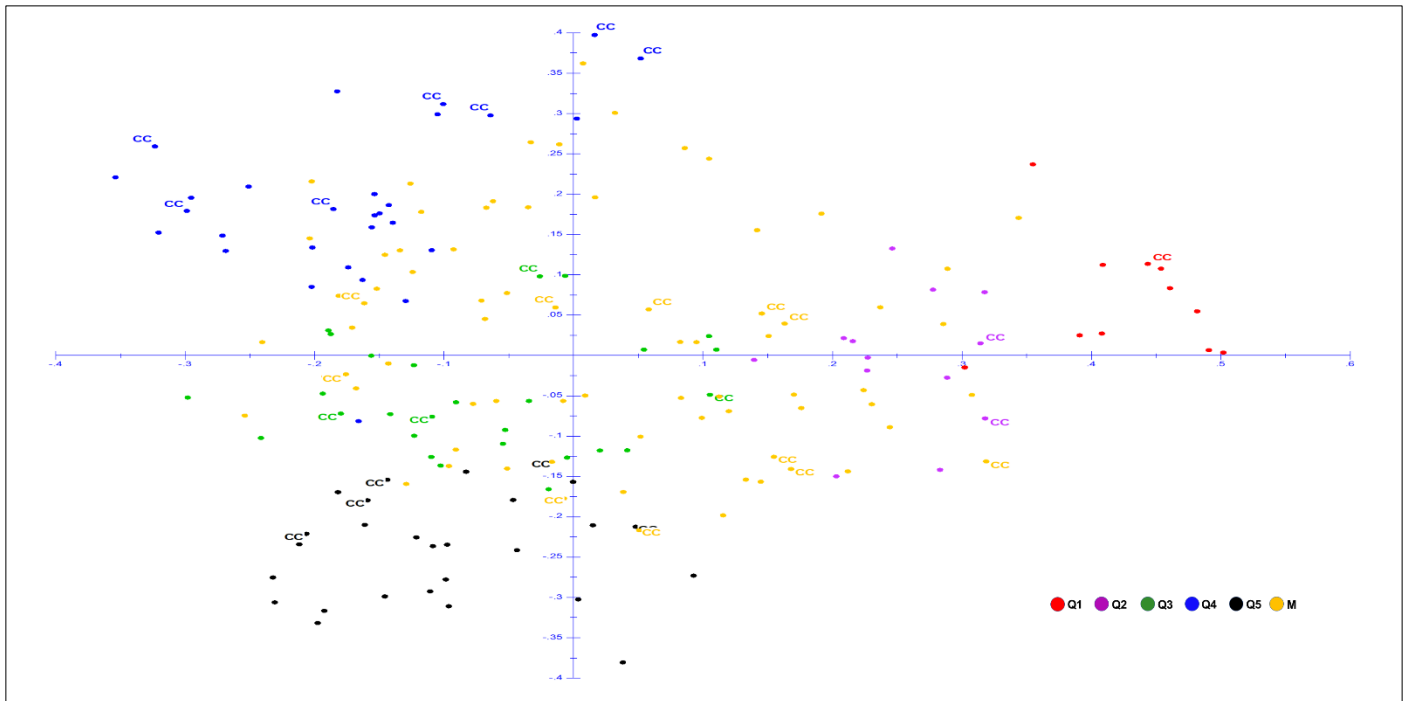


Fig. 3 Two-dimensional PCoA scatterplot of 181 hazelnut genotypes from WHGC based on Dice's distance. Genotypes are colored according to their assignment to the different gene pools, as inferred by STRUCTURE software at $K = 5$: Central Europe (Q1), British Islands (Q2), Iberian Peninsula (Q3), Italian Peninsula (Q4), Balkans-Black Sea (Q5), and mosaic group (M). Entries of the final core collection (Cv-Dce30) are reported as CC