



Genetic diversity, population structure, and clonal verification in South African avocado cultivars using single nucleotide polymorphism (SNP) markers

Raven Wienk^{1,2,3} · Marja Mostert-O'Neill^{2,3} · Nilwala Abeysekara⁴ · Patricia Manosalva⁵ · Barbie Freeman⁶ · Noëlani van den Berg^{1,2,3}

Received: 2 June 2022 / Revised: 1 November 2022 / Accepted: 2 November 2022
© The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2022

Abstract

Since the start of avocado cultivation in South Africa, superior rootstocks and fruit cultivars have been selected based on morphological traits, which is time consuming and expensive. Technological advances, such as the development of a single nucleotide polymorphism (SNP) genotyping platform for avocado, may reduce these limitations. Therefore, the aim of this study was to implement molecular marker technologies for the validation of clonal material, verification of horticultural varieties, and determining the genetic diversity and population structure of an avocado cultivar germplasm in South Africa. An avocado cultivar breeding population, containing 375 individuals, was genotyped using 384 SNP markers. Our affinity propagation analysis (APA) indicated a 10.74% mislabelling in the germplasm. The principal component analysis (PCA) and discriminate analysis of principal components (DAPC) suggested that the germplasm was admixed in relation to the three known avocado varieties, Guatemalan, Mexican, and West Indian. Additionally, the ancestral origins were determined for 27 individuals with unknown ancestry. Furthermore, the population diversity was assessed and revealed moderate levels of differentiation in the germplasm, suggesting a high level of gene flow between the different populations. This research highlights the value of clonal verification and horticultural variety identification—for the reliable propagation of material with desired traits. The accurate propagation of material and clonal identity could aid avocado growers to link morphological characters and stress tolerance to accurate genetic backgrounds, which could improve the selection of avocados for current and future environmental stressors, especially as Africa is set to be significantly impacted by climate change.

Keywords Avocado · SNPs · Population structure · Genetic admixture · Breeding

Communicated by C. Chen

✉ Raven Wienk
raven.wienk@up.ac.za

- ¹ Hans Merensky Chair in Avocado Research, University of Pretoria, Pretoria, South Africa
- ² Department of Biochemistry, Genetics and Microbiology, Faculty of Natural and Agricultural Sciences, University of Pretoria, Pretoria, South Africa
- ³ Faculty of Natural and Agricultural Sciences, Forestry and Agricultural Biotechnology Institute, University of Pretoria, Pretoria, South Africa
- ⁴ Dominican University of California, San Rafael, CA, USA
- ⁵ Department of Microbiology and Plant Pathology, University of California Riverside, Riverside, CA, USA
- ⁶ Subtropical Horticulture Research Station, USDA-ARS, Miami, FL, USA

Introduction

Avocado (*Persea americana*) comes mainly from three “varieties”, *P. americana* var. *americana* Mill. (“West Indian”), *P. americana* var. *guatemalensis* Williams. (“Guatemalan”), and *P. americana* var. *drymifolia* Blake. (“Mexican”) (Lahav and Lavi 2002; Wolstenholme 2003). Intraspecific reproduction between varieties has led to extensive hybridisation with varying agronomical traits (Davis et al. 1998; Ashworth and Clegg 2003). The South African avocado industry relies heavily on superior rootstocks and cultivars, with important morphological traits such as improved fruit yield, better fruit quality, and resistance/tolerance to abiotic and biotic factors, which are usually hybrids (Popenoe and Williams 1947). These hybrids typically show remarkable morphological similarities, making selection and verification of propagated material difficult (Popenoe and Williams

1947). These morphological trait assessments, selection, and development of new cultivars and rootstocks can extend over 20 years—which is a laborious, resource-intensive, and time-consuming process (Köhne 2005; Schaffer et al. 2013).

Advances in technology can now assess an individual on the genotypic level to genetically classify and verify horticultural varieties (Schnell et al. 2003) using molecular markers (Karp et al. 1997). Molecular markers are excellent for genomic and evolutionary studies, clonal verification, identifying cryptic relatedness among individuals, and identifying favourable genotypes linked to phenotypic performances in certain environmental conditions (Batley 2015). Consequently, these technologies have the potential to advance and improve genomic selection by reducing the time and costs involved in phenotyping large numbers of individuals, which is vital to the avocado industry (Clegg 2004).

The most current and popular molecular markers that have been used to study avocado are microsatellites and single nucleotide polymorphisms (SNPs). These markers have been used to investigate sequence nucleotide diversity (Chen et al. 2008), population structure (Chen et al. 2009; Ge et al. 2019a, 2019b; Juma et al. 2020), horticultural variety assignment (Chen et al. 2009), determine genetic diversity (Rubinstein et al. 2019; Ge et al. 2019b; Juma et al. 2020), clarify phylogenetic relationships (Ge et al. 2019a), provide clonal and cultivar verification (Kuhn et al. 2019c), and create linkage maps (Kuhn et al. 2019b).

No research has been conducted on the genetic diversity and population structure of avocados in South Africa. The aim of the present study was to address the lack of genetic diversity and population structure knowledge by sampling and SNP genotyping individuals from an avocado fruit cultivar population by using the 384 SNP platform developed by Kuhn et al. (2019c). An affinity propagation analysis (APA) (Frey and Dueck 2007) was used for clonal verification and identification of mislabelled individuals. Furthermore, the population structure and genetic diversity were investigated using principal component analysis (PCA) and discriminant analysis of principal components (DAPC) (Jombart et al. 2010). These results will be valuable in the establishment of molecular tools that can be used for the effective execution of conservation and breeding practices in the avocado industry.

Materials and methods

Biological material—germplasm accessions

This study used an avocado breeding population from Tzaneen, Limpopo (South Africa), which was selected and provided by Allesbeste™. It consisted of 375 fruiting cultivar individuals, of which 108 individuals were genetically

unique. As sample collection could be error prone, some trees were sampled in duplicate. Each accession had a unique “accession ID”, thus, individuals with identical “accession IDs” that were sampled from different trees were presumed to be genetically identical “clonal/clones”. These were sampled to confirm clonal identity and determine the technical error rate. Whereas, individuals with identical “accession IDs” that were sampled from the same tree were classified as “duplicates”; these were sampled to determine the machine error rate. Allesbeste™ provided proprietary material for this study, and as such, all accessions have been anonymised. All individuals that were genotyped are summarised in Supplementary File 1 (<https://doi.org/10.25403/UPresearchdata.19145087>).

SNP genotyping

Ten green, fleshy leaves at intermediate expansion with minimal damage were sampled from each tree. DNA was isolated from the leaf material by the USDA-ARS using the method described by Kuhn et al. (2017). Each avocado individual was genotyped with 384 bi-allelic SNP markers run on the Fluidigm EP1™ system with the 96.96 IFC (Fluidigm, San Francisco, CA, USA), with 91 DNA samples and five controls, as previously described by Kuhn et al. (2019c). Samples were SNP genotyped by the USDA-ARS. Additional published SNP genotypic data (Kuhn et al. 2019c) was incorporated for population structure analyses and horticultural variety verification; these individuals were from three horticultural varieties, believed to be ancestral to the South African germplasm. These individuals were labelled as “UCR”, that included two Guatemalan (G), six Mexican (M), and four West Indian (WI) individuals, provided by the University of California, Riverside, USA (Kuhn et al. 2019c).

Affinity propagation analysis and visualisation of genotypic data

The data was reformatted to proceed with downstream processes using a custom Perl (Version 5.28.1) script to extract and reformat the genotype information into four categories, as previously described by Kuhn et al. (2019c). Markers and individuals with greater than 5% missing data were removed in a recursive fashion, as previously described by Kuhn et al. (2019c). Consequently, 61 individuals and eight markers were removed and excluded from further analysis. Therefore, 326 individuals (including the 12 references and 107 unique accession IDs) and 376 markers remained from the cultivar data. This dataset was named “APA dataset”, as seen in Table 1.

Custom Python distance and similarity scripts were used to generate pairwise distances (Python—Version 3.8.6), as

Table 1 The number of individuals genotyped and the germplasm sources used in this study, including the published dataset—12 references

Germplasm	Population		Total
	Cultivar germplasm	Published horticultural references ^c	
Source	Allesbeste™	University of California, Riverside	
Location	Tzaneen, Limpopo, South Africa	Various Locations	
Number of individuals sampled (original dataset)	375	12	387
Number of individuals retained for APA analysis ^a (APA dataset)	314	12	326
Number of individuals retained for population analysis ^{a, b} (Population analysis dataset)	147	12	159

^aIndividuals with greater than 5% missing data were removed^bClonal or duplicate individuals were removed^cPublished horticultural reference individuals (Kuhn et al. 2019c)

described by Kuhn et al. (2019c). The similarity matrix was used to perform an affinity propagation analysis (APA) that generated clusters and aided in the identification of mislabelled individuals and confirmation of clonal material (Frey and Dueck 2007; Bodenhofer et al. 2011; Pedregosa et al. 2011; Kuhn et al. 2019a). Additionally, individuals were assigned silhouette scores (Rousseeuw 1987), as described by Kuhn et al. (2019a). Genotype statistics were obtained by the visualisation and sorting of the genotypic data by accession IDs, affinity groups, silhouette scores, and genotypic profiles in Microsoft Excel (2019), as described by Kuhn et al. (2019c).

The number of SNP differences, machine genotyping error, and technical error were calculated for each “clonal” and “duplicate” set of individuals. The machine genotyping error was calculated using the “duplicate” individuals—individuals sampled from the same tree multiple times. The technical error was calculated using “clonal” individuals—identical “accession ID” individuals sampled from different trees. Mislabelled individuals were identified in two ways; firstly, individuals with identical “accession IDs”, but had different genotypic SNP profiles beyond machine genotyping error, were classified as mislabelled type 1. Secondly, individuals with different “accession IDs”, but had similar genotypic SNP profiles within machine genotyping error rate, were classified as mislabelled type 2. Mislabelled individuals were highlighted in red in Supplementary File 1 (<https://doi.org/10.25403/UPresearchdata.19145087>).

Phylogenetic analysis

The APA dataset was used to perform a hierarchical cluster analysis to study the individuals with similar genetic characteristics and aid in the identification of mislabelled accessions. The dataset was aligned using MUSCLE (Edgar 2004)

and subsequently used to construct a condensed unweighted pair group method with arithmetic mean (UPGMA) tree (Sneath and Sokal 1973) using the maximum composite likelihood method, with the confidence examined using bootstrap values calculated for 1000 replicates in MEGA X (Kumar et al. 2018). The dendrogram was exported in the Newick format to be visualised and customised in Interactive tree of life (iTOL) v6 (Letunic and Bork 2019).

Population structure analysis

After the identification of the mislabelled individuals through the APA and phylogenetic analysis, individuals with the least missing data for each “clone” and “duplicate” within machine genotyping error rate were retained. Consequently, 167 individuals were removed and excluded from further analysis; these were highlighted in yellow in Supplementary File 1 (<https://doi.org/10.25403/UPresearchdata.19145087>). Thus, 159 cultivar individuals, of which 12 were published horticultural references, were retained and this second dataset was named “Population analysis dataset”, as seen in Table 1, and was used to perform the principal component analysis (PCA), discriminant analysis of principal components (DAPC), and structure and diversity analysis. Additionally, one non-polymorphic marker was detected during this analysis and removed—marker SHR-SPaS006061; thus, 375 markers remained for the structure analysis. The reduced dataset was reformatted in Microsoft Excel (2019) into a four-bit binary code with A as (1), C as (2), G as (3), and T as (4).

The PCA (Patterson et al. 2006; Reich et al. 2008), DAPC, and allele composition analysis were performed using the Adegenet package (Jombart 2008, Jombart et al. 2010), whereas the genetic diversity was determined using the MMOD 1.3.3 package (Winter 2012). These analyses were all performed in

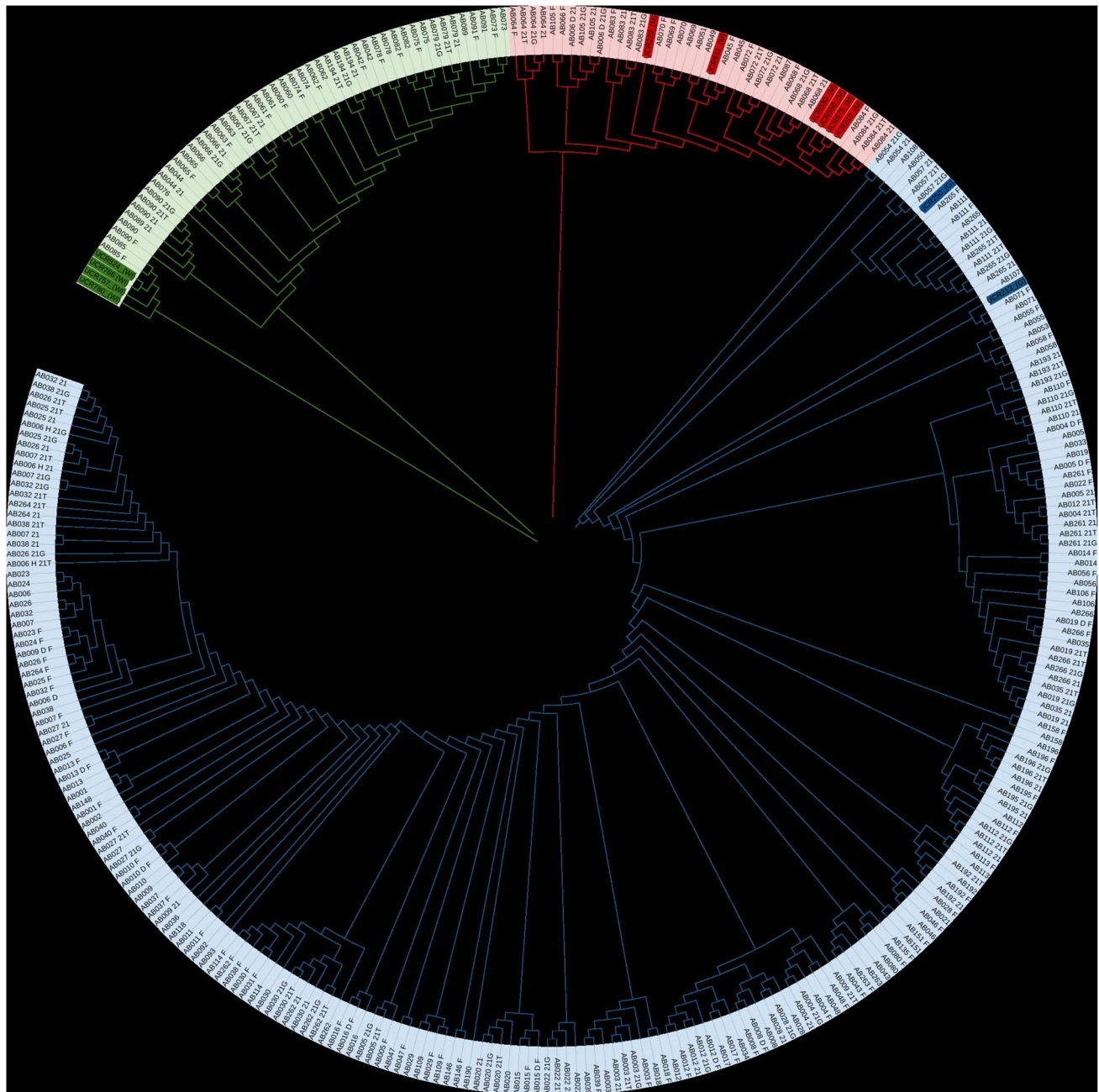


Fig. 1 Dendrogram of the 326 avocado trees constructed with UPGMA showing genetic relationships between the analysed samples. Leaves and branches were coloured according to the horticultural

variety based on the reference individuals. G, Guatemalan (blue); M, Mexican (red); WI, West Indian (green)

RStudio, version 1.3.1093 (RStudio Team 2016) using R version 4.0.3 (R Development Core Team 2020).

The PCA was performed to display the genetic relationships among individuals, genetically classify and verify the horticultural variety of individuals, and detect structure within the germplasm. The germplasm was analysed in relation to published SNP genotypic data (Kuhn et al. 2019c), which represented the three horticultural varieties

(Guatemalan, G; Mexican, M; West Indian, WI). The number of principal components (PCs) retained was based on preserving majority of the variance while retaining the fewest PCs (Jombart 2008). The variance explained by each PC was calculated as the ratio of each eigenvalue to the sum of all calculated eigenvalues.

DAPC was performed to determine the genetic differentiation between different clusters of individuals using

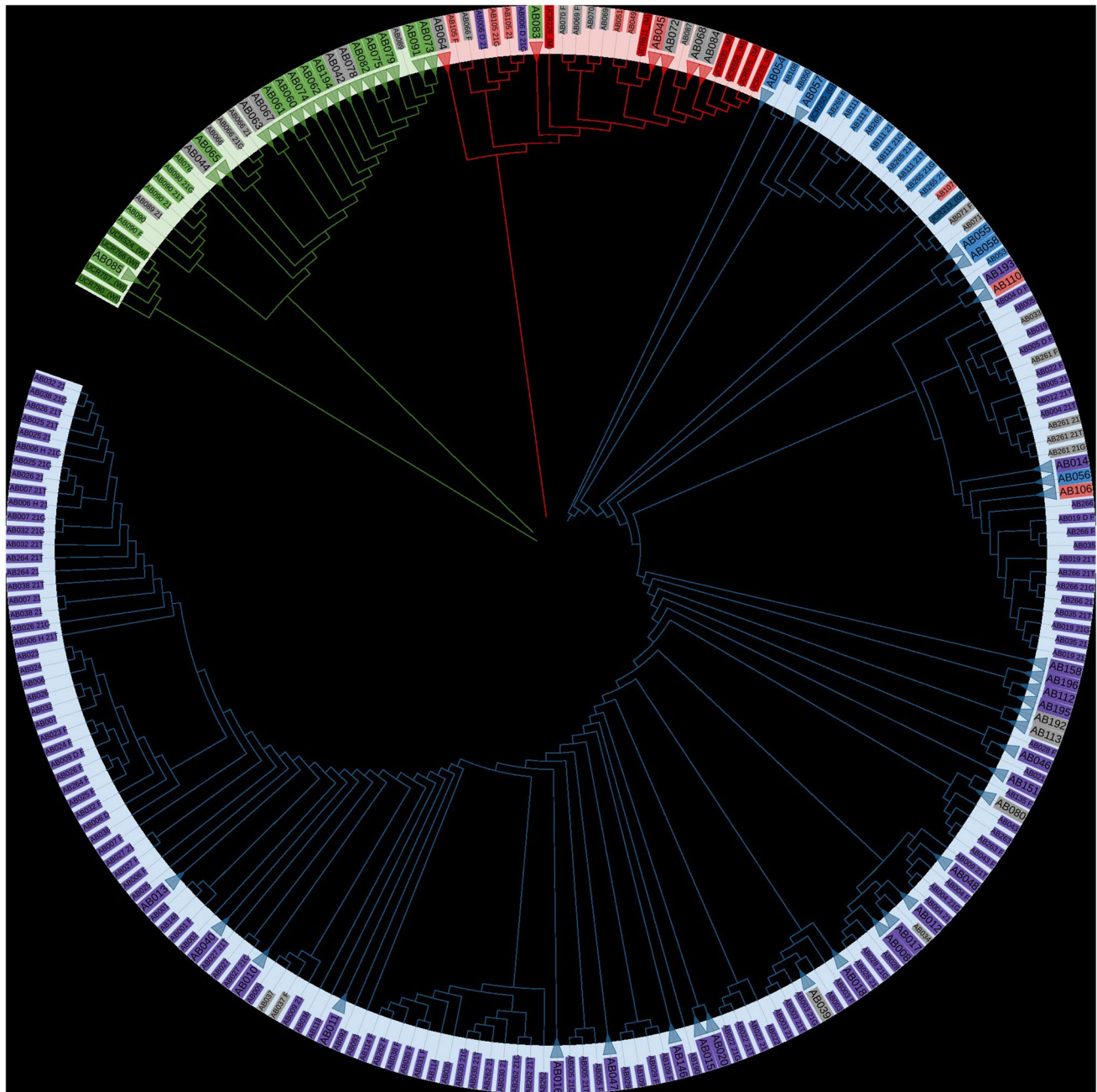


Fig. 2 Dendrogram of the 326 avocado trees constructed with UPGMA showing genetic relationships between the analysed samples. Leaves with identical accessions were collapsed into nodes. The reference individuals were coloured according to the horticultural variety (darker shades), G, Guatemalan (blue); M, Mexican (red); WI, West Indian (green); hybrids (purple); and unknown (grey)

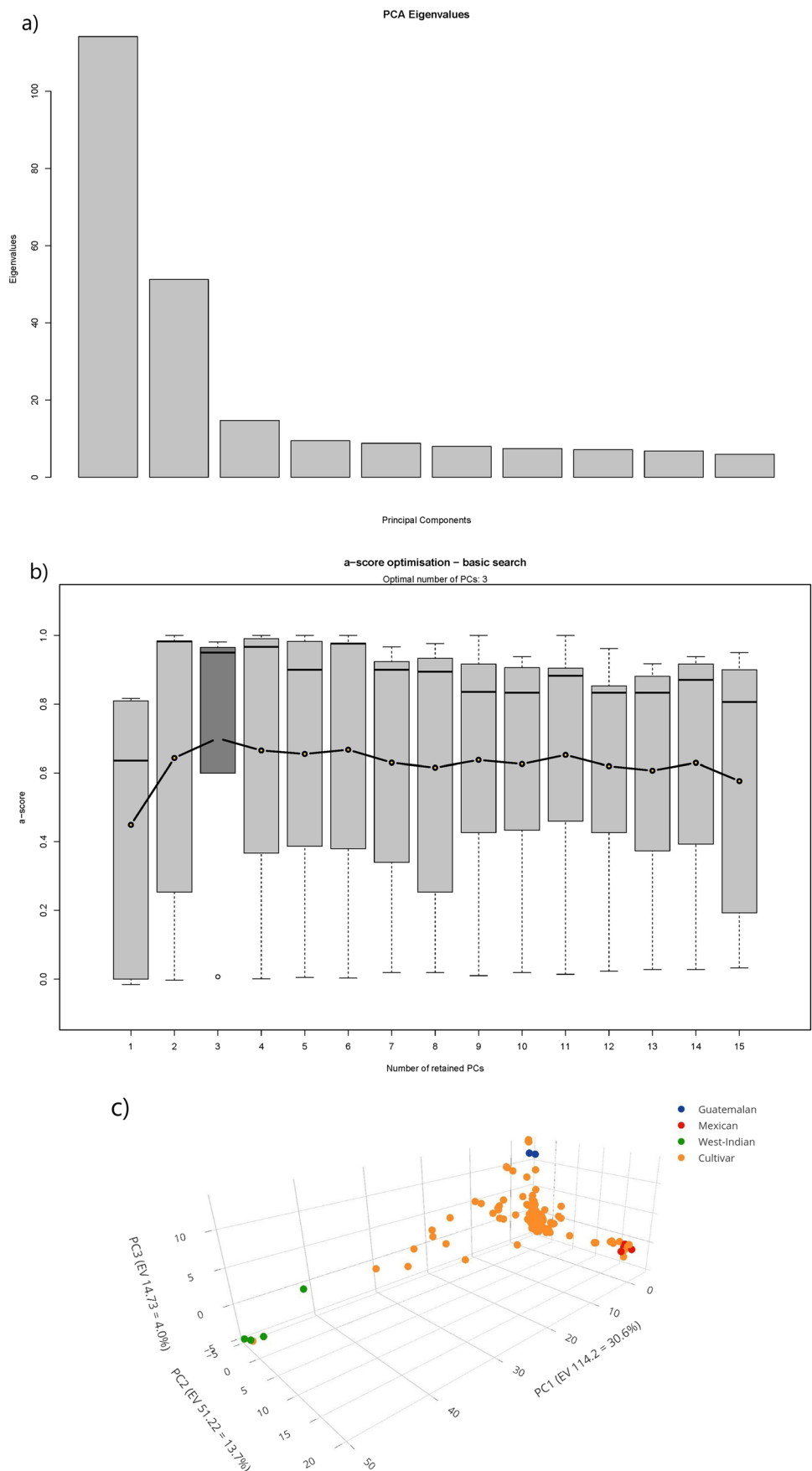
West Indian (green). The leaves/nodes representing the individuals were coloured according to horticultural variety information provided by Allesbeste™. G, Guatemalan (blue); M, Mexican (red); WI, West Indian (green); hybrids (purple); and unknown (grey)

the `find.clusters()` function to determine the number of groups (K) de novo, with the optimal K selected using the `diffNgroup` method (Jombart 2008). The number of PCs to retain was determined using the `optim.a.score()` function (Jombart 2008). The clusters are considered populations, as it may indicate the individual's horticultural variety. The allele composition analysis/membership probabilities were

displayed using the `compoplot()` function (Jombart 2008). PCA and DAPC data were imported and visualised using the Plotly R Chart Studio (Plotly Technologies Inc 2015).

Measures of genetic diversity were evaluated with several “ F_{ST} analogues”, specifically, Nei's G_{ST} (Nei 1973; Nei and Chesser 1983), Hedrick's G_{ST} (Hedrick 2005), and Jost's D (Jost 2008) and estimators for H_s and H_t using the `diff_stats`

Fig. 3 Principal component analysis (PCA) of the 159 avocado cultivar germplasm using 375 single nucleotide polymorphisms (SNPs). **a** PCA eigenvalues of the analysis. **b** The optimal a-score indicating the number of principal components (PCs) to retain for analysis, indicating three PCs. **c** The eigenvalues and variance of each PC are found within parentheses on each axis. Individuals are represented as dots and the reference varieties are represented by G, Guatemalan (blue); M, Mexican (red); WI, West Indian (green); and the Allesbeste™ cultivar germplasm is represented in orange



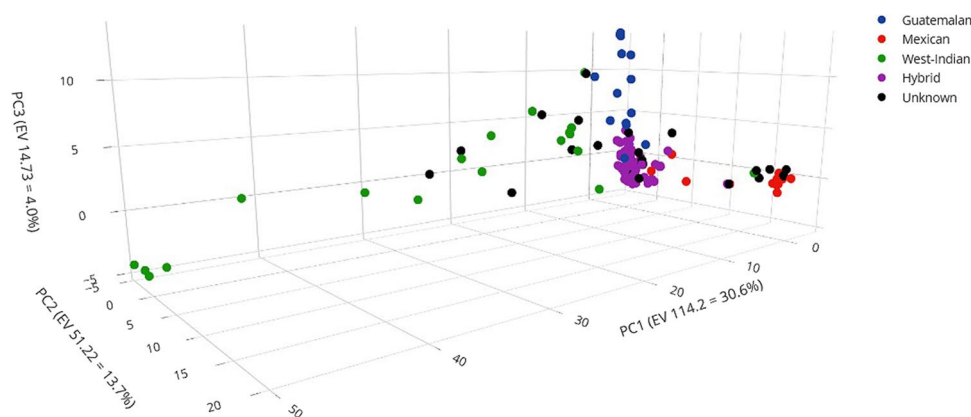


Fig. 4 Principal component analysis (PCA) of the 159 avocado cultivar germplasm—coloured according to the suspected horticultural varieties based on Allesbeste™ breeding records. The eigenvalues of each principal component are found within parentheses on each axis.

Individuals are represented as dots and the horticultural varieties are G, Guatemalan (blue); M, Mexican (red); WI, West Indian (green); G×M, Guatemalan×Mexican hybrid (purple); and unknown variety (black)

function (Meirmans and Hedrick 2011; Winter 2012). H_s and H_t are estimates of the heterozygosity expected for this population with and without sub-populations, respectively. Population divergence was estimated between all combinations of population clusters nested within varieties, using the pairwise_ G_{ST} -Nei, pairwise_ G_{ST} -Hendrick, and pairwise_ D functions; furthermore, the chao_bootstrap function was applied to the populations to determine the robustness of the analysis (Winter 2012). The 12 reference individuals were removed from the “Population analysis dataset” to prevent the reference samples from skewing the analysis.

Results

SNP genotyping statistics and affinity propagation analysis

After removing individuals and markers with more than 5% missing data, the cultivar population contained 326 individuals (including 12 reference individuals and 107 unique accession IDs) and 376 markers. Missing data per individual varied from 0 to 17 markers of the 376 markers; thus, the average missing data from individuals was 3.0 or 0.8%. Missing data for markers varied from 0 to 15 for the 326 individuals; thus, the average missing data of all markers was 2.6 or 0.7%. The heterozygous allele calls for individuals ranged from 2.4 (9/370, six missing data) for accession “UCR524 (WI)” to 75% (282/376) for accession “AB042”, and the heterozygous allele calls for markers ranged from 0 (0/324, two missing data) for SNP marker “SHRSPaS006061” to 79.3% (253/319, seven missing data) for SNP marker “SHRSPaS002697”. Average allele frequency over all markers for allele 1 was 33.8% and allele 2 was 33.7%.

The APA generated 64 cultivar groups for 326 individuals and groups varied from one to 43 individuals. The machine genotyping error ranged from 0 to 2.02% for accession “AB006” with 38 SNP differences. The cultivar germplasm technical genotyping error ranged from 0 to 1.46% for accessions “AB035 and AB266” with 11 SNP differences. The cultivar germplasm contained 35 individuals which were mislabelled, thus, indicating that approximately 10.74% mislabelling is present in the cultivar germplasm (21 individuals were type 1, four individuals were type 2, and 10 individuals were both types 1 and 2). Formatted data with affinity propagation groups, silhouette scores, and genotype data are recorded in Supplementary File 1 (<https://doi.org/10.25403/UPresearch.data.19145087>).

Phylogenetic analysis

The genetic distance matrices of the 326 avocado individuals were used to study the genetic relationships in the population through hierarchical clustering. The phylogenetic analysis indicated that the germplasm was divided into three main populations, based on the reference individuals which are highlighted in the darker shade of the respective colours, as seen in Fig. 1. Based on breeding records and suspected horticultural variety provided by industry, the individuals were coloured accordingly. The UPGMA-based dendrogram produced three major groups, some containing individuals from different horticultural varieties, pointing at genetic admixture between varieties, as seen in Fig. 2. Majority of the individuals from the phylogenetic analysis corresponded with the APA, with the exception of some of the mislabelled individuals.

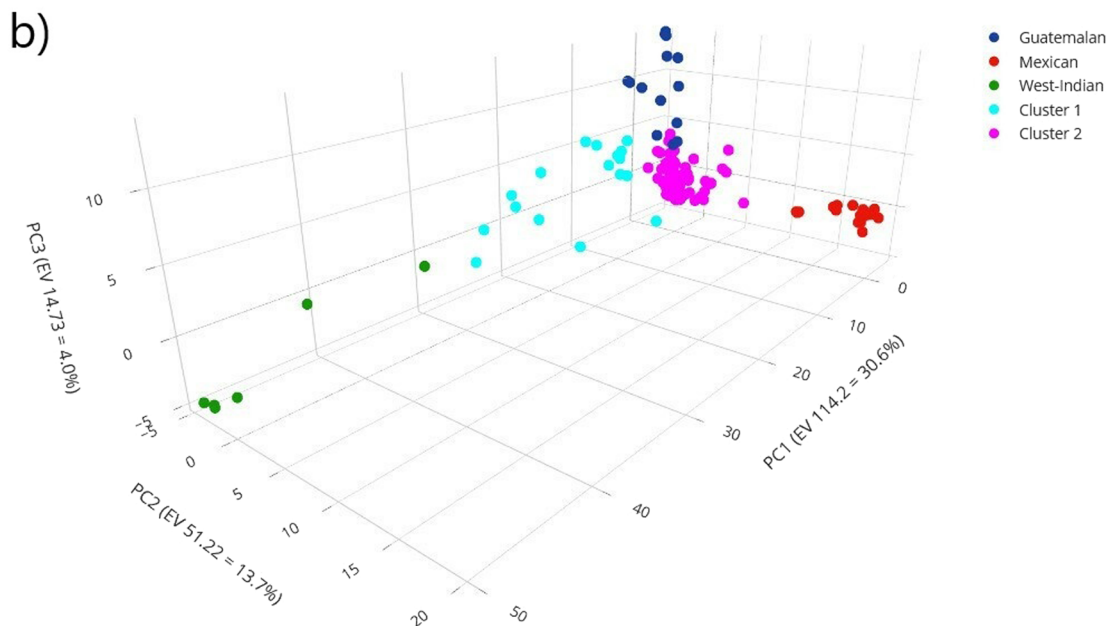
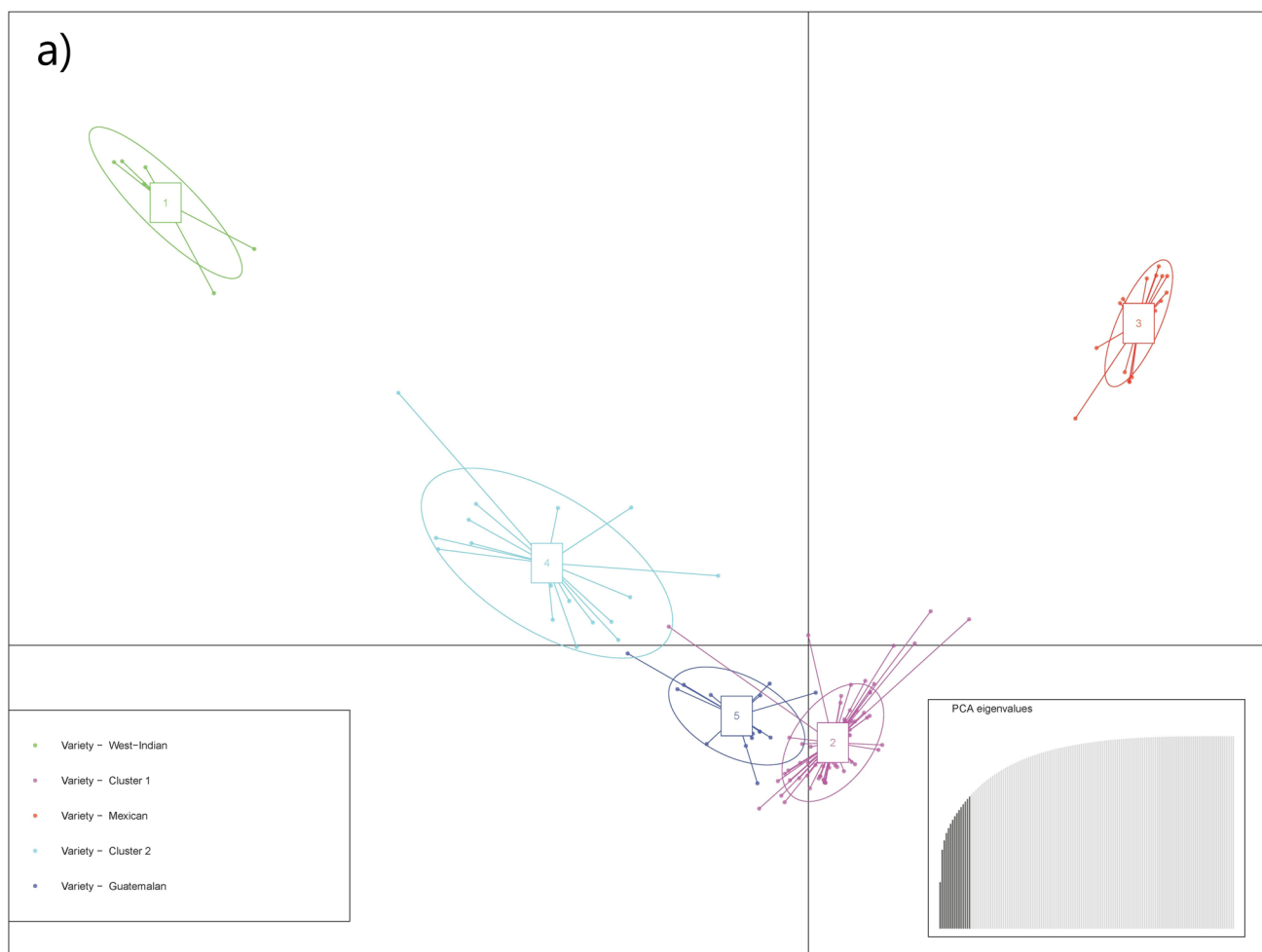


Fig. 5 Discriminant analysis of principal components (DAPC) of the 159 avocado cultivar germplasm using 375 single nucleotide polymorphisms (SNPs). **a** Scatterplot shows the first two PCs of the DAPC for $K=5$, with clusters shown by different colours and inertia ellipses, while dots represent individuals. The PCA eigenvalue plot is inset on the bottom right. **b** Scatterplot shows the first three PCs of the DAPC for $K=5$; the eigenvalues and variance of each principal component are found within parentheses on each axis. Individuals are represented as dots and the varieties are represented by G, Guatemalan (blue); M, Mexican (red); WI, West Indian (green); cluster 1 (cyan); and cluster 2 (magenta)

Principal component analysis and population structure analysis

The PCA was used to study the genetic relationships in the cultivar germplasm. The first three eigenvalues were 114.15, 51.22, and 14.73, respectively. The variance explained by the first three PCs were 30.6%, 13.7%, and 4.0%, respectively; hence, the overall variation was 48.3%. The eigenvalues of the analysis showed that majority of the genetic variance was captured by the first three PCs, as seen from the PCA eigenvalue bar graph in Fig. 3 a and the a-score optimisation in Fig. 3 b. The scatterplot of the first three PCs for the cultivar germplasm indicated that the reference individuals for Guatemalan (G, blue), Mexican (M, red), and West Indian (WI, green) had well-defined clusters. The Allesbeste™ cultivar germplasm (orange) appeared to show a cline between the reference individuals, indicating possible genetic admixture, as seen in Fig. 3 c. The WI cluster separated from the G and M clusters along the first PC, whereas the M cluster separated from the G cluster along the second PC. The Allesbeste™ cultivar germplasm clustered mainly between the G and M clusters, with most individuals grouping closer to the G cluster, as seen in Fig. 3 c.

Further analysis of the cultivar germplasm based on the above PCA, with the individuals now coloured according to the suspected horticultural variety provided by Allesbeste™ based on breeding records, revealed that the majority of the individuals were G and or G×M hybrids. The cultivar germplasm had 27 individuals with unknown horticultural variety, as seen in Fig. 4, which were resolved with 16 individuals assigned as G, eight as M, and three as WI. Additionally, there were 17 misclassified individuals, which were reassigned.

The DAPC was used to investigate the population differentiation between groups of individuals and identify clusters of genetically related individuals. Based on information from literature and industry, hybrids are common between the horticultural varieties; therefore, a DAPC analysis was performed from $K=2$ until $K=7$ to identify potential hybrids in the germplasm. The eigenvalues of the analysis showed that majority of the genetic variance was again captured by the first three PCs. According to the diffNgroup method, the optimum number of genetic clusters were $K=5$,

which was best supported and appeared to be the most biologically relevant scatterplot. This scatterplot shows the first two PCs of the DAPC for $K=5$, as seen in Fig. 5 a. Clusters are shown by different colours and inertia ellipses, while dots represent individuals, indicating the Guatemalan (blue), Mexican (red), and West Indian (green), cluster 1 (cyan, possible G×WI hybrids), and cluster 2 (magenta, possible G×M hybrids). Three groups of genetically closer clusters can be identified, Guatemalan, cluster 1, and cluster 2, as seen in Fig. 5 a. This scatterplot also indicates that majority of the West Indian accessions are hybrids. Additionally, the scatterplot was shown using the first three PCs of the DAPC of the cultivar germplasm for $K=5$, indicating the Guatemalan (blue), Mexican (red), and West Indian (green), cluster 1 (cyan), and cluster 2 (magenta), as seen in Fig. 5 b. The scatterplot showed a cline, indicating genetic admixture between the genetic clusters. The cultivar germplasm consisted of 9.4% G, 13.2% M, 3.8% WI, 10.7% cluster 1 (possible G×WI hybrids), and 62.9% cluster 2 (possible G×M hybrids). Majority of the results from the DAPC matched the suspected horticultural variety provided by Allesbeste™.

The allele composition analysis of the cultivar germplasm indicated the inferred structure and membership probabilities, where each individual is represented by a coloured bar with length proportional to the estimated membership to each cluster (Pritchard et al. 2000), as seen in Fig. 6. Majority of the germplasm individuals were composed of the G cluster and cluster 2 (G×M hybrid), which corresponds to the DAPC results. The reference individuals are located in the enclosed area, from individual 147 to 159 in the genomic composition plot, as seen in Fig. 6. All genomic composition plots from $K=2$ until $K=7$ is recorded in Supplementary Fig. 1.

Diversity analysis

The genetic differentiation and diversity present in the germplasm population were evaluated with three “ F_{ST} analogues”, Nei’s G_{ST} , Hedrick’s G_{ST} , and Jost’s D , additionally H_s and H_t are estimates of the heterozygosity expected for this population with and without the sub-populations defined in the data, respectively. This analysis indicated that the WI vs. cluster 1 (Nei’s $G_{ST}=0.050$, Hedrick’s $G_{ST}=0.148$, Jost’s $D=0.058$) had the least genetic differentiation, whereas M vs. WI (Nei’s $G_{ST}=0.525$, Hedrick’s $G_{ST}=0.847$, Jost’s $D=0.509$) had the highest genetic differentiation, as seen in Table 2.

Discussion

The aim of this study was to use a set of previously developed SNP markers for the validation of clonal material, verification of horticultural variety, and determination of

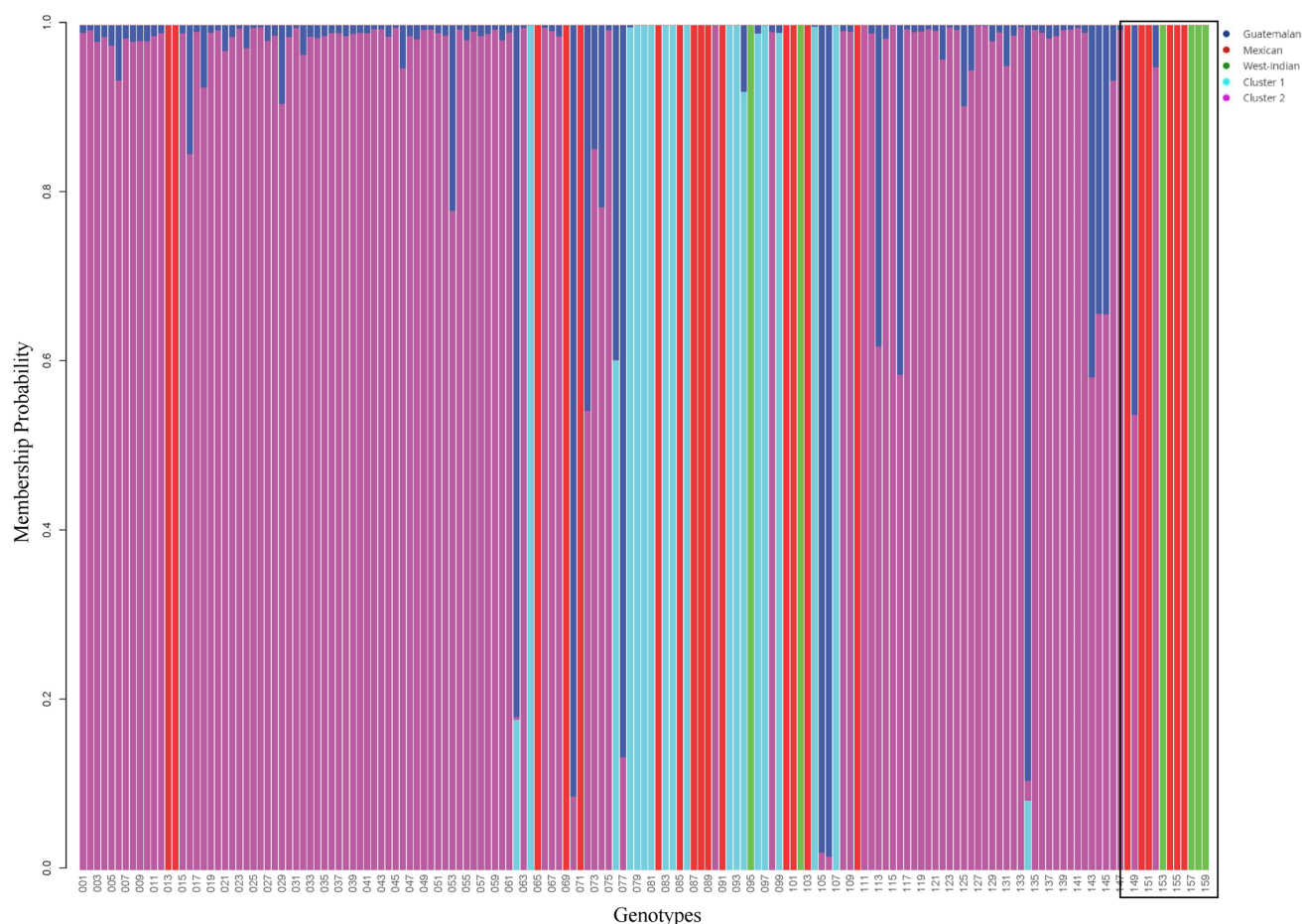


Fig. 6 Genomic composition plot of the discriminant analysis of principal components (DAPC) indicating the cluster's composition for the cultivar germplasm for 159 genotypes. Each thin vertical line in the bar plot represents one individual and each colour represents one inferred ancestral population. The length of each colour in a vertical bar represents the proportion of that individual's ancestry that is

derived from the inferred ancestral population corresponding to that colour. The same colour in different individuals indicates that they belong to the same cluster, indicating they share the same ancestral population. Clusters: Guatemalan (blue), Mexican (red), West-Indian (green), cluster 1 (cyan), and cluster 2 (magenta). Reference individuals are located in the enclosed area

Table 2 Global population pairwise F_{ST} comparison among the populations at $K=5$ identified by the discriminant analysis of principal components

	Hs	Ht	Pairwise Nei's G_{ST}	Pairwise Hedrick's G_{ST}	Pairwise Jost's D
G vs. M	0.226	0.313	0.278	0.562	0.225
G vs. WI	0.287	0.438	0.345	0.719	0.423
G vs. cluster 1	0.375	0.429	0.126	0.359	0.173
G vs. cluster 2	0.285	0.306	0.069	0.181	0.059
M vs. WI	0.187	0.394	0.525	0.847	0.509
M vs. cluster 1	0.287	0.399	0.28	0.613	0.313
M vs. cluster 2	0.199	0.266	0.252	0.502	0.167
WI vs. cluster 1	0.356	0.374	0.05	0.148	0.058
WI vs. cluster 2	0.254	0.44	0.422	0.796	0.499
Cluster 1 vs. cluster 2	0.345	0.424	0.187	0.481	0.242

Hs and Ht are estimates of the heterozygosity expected for this population with and without the subpopulations defined in the data, respectively

Clusters: Guatemalan (G), Mexican (M), West-Indian (WI), cluster 1 (possible G × WI hybrids), cluster 2 (possible G × M hybrids)

the genetic diversity and population structure of an avocado cultivar breeding population in South Africa.

An APA was used to identify mislabelled individuals and confirm clonal material. An APA uses all points simultaneously with no genetic assumptions to determine which individuals would best serve as epitomes and the clustering occurs naturally, thus, decreasing erroneous results (Frey and Dueck 2007). Previously, Kuhn et al. (2019c) used an APA to identify 38 mislabelled individuals in the USDA-ARS Subtropical Horticulture Research Station (SHRS) germplasm collection, thus, indicating 13% mislabelling. Similarly, in this study, the APA identified 35 mislabelled individuals, thus, indicating approximately 10.74% mismatch ratio in the cultivar germplasm. Mislabelling in breeding populations can occur in every phase of avocado production, including incorrect identification in the field, propagation, as well as during procurement of samples and during genotyping (Kuhn et al. 2019c). It is important to identify mislabelled individuals in germplasms to prevent the propagation of incorrect material, which could be used for budwood purposes. Therefore, identifying these mislabelled individuals may improve breeding efficiency and deployment, while reducing loss of resources and time.

The PCA was used to verify the horticultural variety of individuals in the cultivar germplasm. PCA identifies genetic structures among individuals in the absence of any assumption about the underlying population genetic model (Patterson et al. 2006; Reich et al., 2008), as well as summarises the overall variability in a population. Based on the PCA, the majority of the South African cultivar germplasm grouped between the Guatemalan and Mexican varieties, and the population appeared to show a cline, rather than well-defined clusters, indicating evidence of genetic admixture. Thus, the cultivar germplasm appeared to consist mainly of $G \times M$ hybrids. Furthermore, the cultivar germplasm had 27 individuals of unknown horticultural variety, and 17 individuals with misclassified horticultural variety, which were resolved. Genetic admixture among avocado populations is attributed to the extensive hybridisation between varieties; and this is common as avocado varieties do not have sterility barriers (Davis et al. 1998; Ashworth and Clegg 2003). Hybrids allow for a desirable blend of important traits in one individual, such as disease resistance and improved yield. Unfortunately, PCA summarises the overall variability in a population and requires an aforementioned definition of clusters to study population structures; thus, these drawbacks warranted further investigation through DAPC.

The DAPC was used to determine the population structure of the cultivar germplasm, as it is a multivariate model which assesses the genetic differentiation between different clusters of individuals into groups, while maximising between-group variability and minimising within-group variation (Fisher 1936; Lachenbruch and Goldstein 1979;

Jombart 2008; Jombart and Ahmed 2011). DAPC has a few advantages, such as the probabilistic assignment of individuals to groups (like Bayesian approaches) and the visual assessment of structures for different population genetic models (Jombart et al. 2010). In this study, the DAPC allowed for the verification of the horticultural variety of individuals in the breeding population. Based on the DAPC, the cultivar germplasm consisted of 9.4% Guatemalan, 13.2% Mexican, 3.8% West Indian, 10.7% cluster 1 (possible $G \times WI$ hybrids), and 62.9% cluster 2 (possible $G \times M$ hybrids). The high percentage of Guatemalan, Mexican, and possible $G \times M$ hybrids in the germplasm is coherent, as the most popular cultivar grown worldwide is hass (Crane et al. 2013), which is a $G \times M$ hybrid (Rendón-Anaya et al. 2019). Furthermore, the Guatemalan variety has high fruit averages and horticultural quality, whereas the Mexican variety has a desirable fruit size (Bergh and Ellstrand 1986) and has shown some tolerance and resistance to *Phytophthora cinnamomi* (Sánchez-González et al. 2019), which are valuable traits in the industry. Furthermore, Guatemalan and Mexican varieties are typically grown in less tropical areas (Williams 1977), such as avocado-growing regions in South Africa.

Interestingly, majority of the West Indian accessions in the cultivar germplasm appeared to be $G \times WI$ hybrids, even though the industry records indicated these are West Indian accessions. These $G \times WI$ hybrids have been known to have an early harvest period and bridges harvesting gaps (Bergh and Ellstrand 1986), which could explain the presence of cluster 1 (possible $G \times WI$ hybrids) in the germplasm. However, there does not appear to be any $M \times WI$ hybrids within the population. This may be due to lack of sampling or due to the lack of breeding of $M \times WI$ hybrids in South Africa. Some West Indian individuals are more tolerant to salinity and calcareous soils (Ben-Ya'acov and Michelson 1995), which is not favoured by most avocado cultivars grown in South Africa. Most commercial avocado rootstocks and cultivars are hybrids (Popenoe and Williams 1947); hence, it is important to correctly identify the horticultural variety of individuals, as this affects the ability of breeding programmes to select accurate and superior individuals. A concern involved in this study is the precise DAPC assignment of individuals, as it may be skewed by the lack of reference samples utilised during analysis (Ottewell et al. 2016). Furthermore, it is important to curate more avocado germplasms in South Africa to include potential $M \times WI$ hybrids and improve our understanding of the population. An informative addition to this study would involve linking the genotypic data with phenotypic data to provide a more rounded description of the germplasm at hand.

Genetic diversity was determined with “ F_{ST} analogues” that assessed the within and among population variation. MMOD is a package that allows three different “ F_{ST} analogues” to be evaluated, Nei G_{ST} , Hendrick’s G_{ST} , and Jost’s

D , which is comparable between studies (Winter 2012). These “ F_{ST} analogues” and their combined use will allow more robust analyses of population structure than what is achievable with only F_{ST} (Meirmans and Hedrick 2011). Some previously reported F_{ST} values for avocado germplasms among the three varieties were 0.19, 0.22, and 0.25, reported by Boza et al. (2018), Guzmán et al. (2017), and Gross-German and Viruel (2013), respectively, whereas lower F_{ST} values of 0.061 and 0.05 were reported by Juma et al. (2020) and Cañas-Gutiérrez et al. (2019) respectively. In this study, the “ F_{ST} analogues” indicated that the West Indian vs. cluster 1 (possible $G \times WI$ hybrids; Nei $G_{ST} = 0.050$) had the least genetic differentiation, whereas Mexican vs. West Indian had the highest genetic differentiation (Nei $G_{ST} = 0.525$).

These studies show the varying levels of diversity in numerous avocado germplasms worldwide. These diversity levels can be affected by the type and number of markers used, the number of individuals and populations assessed, comparable reference samples, and different parameters used for the analysis. Genetic diversity allows for a species to adapt to various environmental conditions and stressors (Schleif 1993), such as climate change and resistance to new emerging pathogens and pests. The cultivar germplasm analysed in this study contained moderate differentiation between varieties and hybrid clusters. The “ F_{ST} analogue” values in this study were similar to other studies, such as Guzmán et al. (2017) and Gross-German and Viruel (2013). Moderate levels of differentiation in the germplasm suggest interbreeding between the three varieties, which is seen with cluster 1 (possible $G \times WI$ hybrids) and cluster 2 (possible $M \times G$ hybrids) in this study. Majority of the cultivar germplasm (62.9%) grouped into cluster 2 (possible $G \times M$ hybrids); this would correlate with industry breeding records.

To our knowledge, this study presents the first molecular genetic assessment of an avocado cultivar germplasm in South Africa. In the present study, molecular marker technology was used to identify mislabelled individuals, validate clonal material, verify horticultural variety, and determine population structure and genetic diversity. The results from the study may prevent the future propagation of incorrect material, establish proper management and conservation strategies, and lastly, improve cultivar breeding efficiency by aiding in the selection of avocado with the ability to cope with changing environments and emerging pests and pathogens. Molecular markers are a powerful and important tool for avocado breeding programmes.

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1007/s11295-022-01573-8>.

Acknowledgements The authors would like to thank the Forestry and Agricultural Biotechnology Institute (FABI) and the University of Pretoria for the use of their facilities and equipment. Furthermore,

the authors would like to thank Dr David Kuhn for the custom affinity propagation scripts. Lastly, I would like to thank Allesbeste™ for providing the plant material.

Author contribution RW contributed to the study design, experimental design, sample curation, formal analysis, investigation, visualisation, and drafting/writing/editing of the manuscript. NVDB contributed to the study conceptualisation and design, experimental design, project administration, resources, supervision, and funding. MMON and NA were responsible for methodology and technical assistance. PM provided the horticultural reference. BF extracted, processed, and performed the SNP genotyping. All co-authors contributed to writing/editing of the manuscript. All authors contributed to and approved the final manuscript.

Funding The authors would like to thank the Hans Merensky Foundation© and Allesbeste™ for funding.

Data availability The cultivar germplasm analysed during this study is available in the University of Pretoria Research Repository [Supplementary File 1—<https://doi.org/10.25403/UPresearchdata.19145087>].

Declarations

Ethics approval This study received specific approval by the appropriate ethics committee for research involving plants.

Conflict of interest The authors declare no competing interests.

References

- Ashworth V, Clegg M (2003) Microsatellite markers in avocado (*Persea americana* Mill.): genealogical relationships among cultivated avocado genotypes. *J Hered* 94:407–415
- Batley J (2015) Plant genotyping. Humana Press, Dordrecht, London
- Ben-Ya'acov A, Michelson E (1995) Avocado rootstocks. Vol. 17 (Janick J, ed.) John Wiley and Sons, Inc, New York, NY, pp 381–429
- Bergh B, Ellstrand N (1986) Taxonomy of the avocado. *Calif Avocado Soc Yearb* 70:135–146
- Bodenhofer U, Kothmeier A, Hochreiter S (2011) APCluster: an R package for affinity propagation clustering. *Bioinformatics* 27:2463–2464
- Boza EJ, Tondo CL, Ledesma N, Campbell RJ, Bost J, Schnell RJ et al (2018) Genetic differentiation, races and interracial admixture in avocado (*Persea americana* Mill.), and *Persea* spp. evaluated using SSR markers. *Genet Resour Crop Evol* 65:1195–1215
- Cañas-Gutiérrez GP, Arango-Isaza RE, Saldamando-Benjumea CI (2019) Microsatellites revealed genetic diversity and population structure in Colombian avocado (*Persea americana* Mill.) germplasm collection and its natural populations. *J Plant Breed Crop Sci* 11:106–119
- Chen H, Morrell P, De La Cruz M, Clegg M (2008) Nucleotide diversity and linkage disequilibrium in wild avocado (*Persea americana* Mill.). *J Hered* 99:382–389
- Chen H, Morell P, Ashworth V, De La Cruz M, Clegg M (2009) Tracing the geographic origins of major avocado cultivars. *J Hered* 100:56–65
- Clegg M (2004) Application of molecular markers to avocado improvement. 24–28. California Avocado Commission, Proceedings of the California Avocado Research Symposium, University of California, Riverside. http://www.avocadosource.com/arac/symposium_2004/arac2004_pg_24.pdf. Accessed 21 May 2020

- Crane J, Douhan G, Faber B, Arpaia M, Bender G, Balerdi C et al (2013) The avocado botany, production and uses: cultivars and rootstocks. In: Schaffer B, Wolstenholme B, and Whaley A (eds) CABI, pp 200–233
- Davis J, Henderson D, Kobayashi M, Clegg M (1998) Genealogical relationships among cultivated avocado as revealed through RFLP analyses. *J Hered* 89:319–323
- Edgar RC (2004) MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res* 32:1792–1797
- Fisher RA (1936) The use of multiple measurements in taxonomic problems. *Ann Eugen* 7:179–188
- Frey BJ, Dueck D (2007) Clustering by passing messages between data points. *Science* 315:972–976
- Ge Y, Tan L, Wu B, Wang T, Zhang T, Chen H et al (2019a) Transcriptome sequencing of different avocado ecotypes: de novo transcriptome assembly, annotation, identification and validation of EST-SSR markers. *Forests* 10:411
- Ge Y, Zhang T, Wu B, Tan L, Ma F, Zou M et al (2019b) Genome-wide assessment of avocado germplasm determined from specific length amplified fragment sequencing and transcriptomes: population structure, genetic diversity, identification, and application of race-specific markers. *Genes* 10:215
- Gross-German E, Viruel M (2013) Molecular characterization of avocado germplasm with a new set of SSR and EST-SSR markers: genetic diversity, population structure, and identification of race-specific markers in a group of cultivated genotypes. *Tree Genet Genomes* 9:539–555
- Guzmán LF, Machida-Hirano R, Borrayo E, Cortés-Cruz M, Espíndola-Barquera MdC, Heredia García E (2017) Genetic structure and selection of a core collection for long term conservation of avocado in Mexico. *Front Plant Sci* 8:243
- Hedrick PW (2005) A standardized genetic differentiation measure. *Evolution* 59:1633–1638
- Jombart T (2008) adegenet: a R package for the multivariate analysis of genetic markers. *Bioinformatics* 24:1403–1405
- Jombart T, Ahmed I (2011) *adegenet* 1.3-1: new tools for the analysis of genome-wide SNP data. *Bioinformatics* 27:3070–3071
- Jombart T, Devillard S, Balloux F (2010) Discriminant analysis of principal components: a new method for the analysis of genetically structured populations. *BMC Genet* 11:1–15
- Jost L (2008) G_{ST} and its relatives do not measure differentiation. *Mol Ecol* 17:4015–4026
- Juma I, Geleta M, Nyomora A, Saripella GV, Hovmalm HP, Carlsson AS et al (2020) Genetic diversity of avocado from the southern highlands of Tanzania as revealed by microsatellite markers. *Hereditas* 157:1–12
- Karp A, Edwards KJ, Bruford M, Funk S, Vosman B, Morgante M et al (1997) Molecular technologies for biodiversity evaluation: opportunities and challenges. *Nat Biotechnol* 15:625–628
- Köhne S (2005) Selection of avocado scions and breeding of rootstocks in South Africa. New Zealand and Australia Avocado Grower's Conference, Tauranga. https://www.avocadosource.com/Journals/AUSNZ/AUSNZ_2005/KohneS2005.pdf. Accessed 30 Apr 2020
- Kuhn D, Bally I, Dillon N, Innes D, Groh A, Rahaman J et al (2017) Genetic map of mango: a tool for mango breeding. *Front Plant Sci* 8:577
- Kuhn D, Livingstone D III, Richards J, Manosalva P, van den Berg N, Chambers A (2019a) Application of genomic tools to avocado (*Persea americana*) breeding: SNP discovery for genotyping and germplasm characterization. *Sci Hortic* 246:1–11
- Kuhn D, Dillon N, Bally I, Groh A, Rahaman J, Warschefsky M et al (2019b) Estimation of genetic diversity and relatedness in a mango germplasm collection using SNP markers and a simplified visual analysis method. *Sci Hortic* 252:156–168
- Kuhn D, Groh A, Rahaman J, Freeman B, Arpaia M, van den Berg N et al (2019c) Creation of an avocado unambiguous genotype SNP database for germplasm curation and as an aid to breeders. *Tree Genet Genomes* 15:71
- Kumar S, Stecher G, Li M, Knyaz C, Tamura K (2018) MEGA X: Molecular evolutionary genetics analysis across computing platforms. *Mol Biol Evol* 35:1547–1549
- Lachenbruch P, Goldstein M (1979) Discriminant analysis. *Biometrics* 35:69–85
- Lahav E, Lavi U (2002) Genetics and classical breeding. In: Whaley A, Schaffer B, Wolstenholme B (eds.) The avocado: Botany, production and uses. CAB International, Wallingford, pp 39–69
- Letunic I, Bork P (2019) Interactive tree of life (iTOL) v4: recent updates and new developments. *Nucleic Acids Res* 47:W256–259
- Meirmans PG, Hedrick PW (2011) Assessing population structure: F_{ST} and related measures. *Mol Ecol Resour* 11:5–18
- Nei M (1973) Analysis of gene diversity in subdivided populations. *Proc Natl Acad Sci* 70:3321–3323
- Nei M, Chesser RK (1983) Estimation of fixation indices and gene diversities. *Ann Hum Genet* 47:253–259
- Ottewill KM, Bickerton DC, Byrne M, Lowe AJ (2016) Bridging the gap: A genetic assessment framework for population-level threatened plant conservation prioritization and decision-making. *Divers Distrib* 22:174–188
- Patterson N, Price AL, Reich D (2006) Population structure and eigenanalysis. *PLoS Genet* 2:e190
- Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O et al (2011) Scikit-learn: machine learning in Python. *J Mach Learn Res* 12:2825–2830
- Plotly Technologies Inc (2015) Collaborative data science. (Plotly Technologies Inc, ed.) Montréal, QC. <https://chart-studio.plotly.com>
- Popenoe W, Williams L (1947) The expedition to Mexico of October 1947. *Calif Avocado Soc Yearb* 1947:22–28
- Pritchard J, Stephens M, Donnelly P (2000) Inference of population structure using multilocus genotype data. *Genetics* 155
- R Development Core Team (2020) R: a language and environment for statistical computing. R Foundation for Statistical Computing
- Reich D, Price AL, Patterson N (2008) Principal component analysis of genetic data. *Nat Genet* 40:491–492
- Rendón-Anaya M, Ibarra-Laclette E, Méndez-Bravo A, Lan T, Zheng C, Carretero-Paulet L et al (2019) The avocado genome informs deep angiosperm phylogeny, highlights introgressive hybridization, and reveals pathogen-influenced gene space adaptation. *Proc Natl Acad Sci* 116:17081–17089
- Rousseeuw P (1987) Silhouettes - a graphical aid to the interpretation and validation of cluster-analysis. *J Comput Appl Math* 20:53–65
- RStudio Team (2016) RStudio: Integrated development for R. RStudio, PBC, Boston, MA. <http://www.rstudio.com>
- Rubinstein M, Eshed R, Rozen A, Zviran T, Kuhn D, Irihimovitch V et al (2019) Genetic diversity of avocado (*Persea americana* Mill.) germplasm using pooled sequencing. *BMC Genomics* 20:379
- Sánchez-González EI, Gutiérrez-Soto JG, Olivares-Sáenz E, Gutiérrez-Díez A, Barrientos-Priego AF, Ochoa-Ascencio S (2019) Screening progenies of Mexican race avocado genotypes for resistance to *Phytophthora cinnamomi* Rands. *HortScience* 54:809–813
- Schaffer B, Wolstenholme B, Whaley A (2013) The avocado: botany, production and uses. CABI, Oxfordshire
- Schleif R (1993) Genetics and molecular biology. Johns Hopkins University Press, Baltimore, MD
- Schnell R, Brown J, Olano C, Power E, Krol C, Kuhn D et al (2003) Evaluation of avocado germplasm using microsatellite markers. *J Am Soc Hortic Sci* 128:881–889

- Sneath PH, Sokal RR (1973) Numerical taxonomy. The principles and practice of numerical classification 1973. WH Freeman and Company, San Francisco
- Williams L (1977) The avocado, a synopsis of the genus *Persea*, subg. *Persea*. Econ Bot 31:315–320
- Winter DJ (2012) MMod: an R library for the calculation of population differentiation statistics. Mol Ecol Resour 12:1158–1160
- Wolstenholme B (2003) Avocado rootstocks: what do we know; are we doing enough research? South Afr Avocado Growers' Assoc Yearb 26:106–112

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.

Terms and Conditions

Springer Nature journal content, brought to you courtesy of Springer Nature Customer Service Center GmbH (“Springer Nature”).

Springer Nature supports a reasonable amount of sharing of research papers by authors, subscribers and authorised users (“Users”), for small-scale personal, non-commercial use provided that all copyright, trade and service marks and other proprietary notices are maintained. By accessing, sharing, receiving or otherwise using the Springer Nature journal content you agree to these terms of use (“Terms”). For these purposes, Springer Nature considers academic use (by researchers and students) to be non-commercial.

These Terms are supplementary and will apply in addition to any applicable website terms and conditions, a relevant site licence or a personal subscription. These Terms will prevail over any conflict or ambiguity with regards to the relevant terms, a site licence or a personal subscription (to the extent of the conflict or ambiguity only). For Creative Commons-licensed articles, the terms of the Creative Commons license used will apply.

We collect and use personal data to provide access to the Springer Nature journal content. We may also use these personal data internally within ResearchGate and Springer Nature and as agreed share it, in an anonymised way, for purposes of tracking, analysis and reporting. We will not otherwise disclose your personal data outside the ResearchGate or the Springer Nature group of companies unless we have your permission as detailed in the Privacy Policy.

While Users may use the Springer Nature journal content for small scale, personal non-commercial use, it is important to note that Users may not:

1. use such content for the purpose of providing other users with access on a regular or large scale basis or as a means to circumvent access control;
2. use such content where to do so would be considered a criminal or statutory offence in any jurisdiction, or gives rise to civil liability, or is otherwise unlawful;
3. falsely or misleadingly imply or suggest endorsement, approval, sponsorship, or association unless explicitly agreed to by Springer Nature in writing;
4. use bots or other automated methods to access the content or redirect messages
5. override any security feature or exclusionary protocol; or
6. share the content in order to create substitute for Springer Nature products or services or a systematic database of Springer Nature journal content.

In line with the restriction against commercial use, Springer Nature does not permit the creation of a product or service that creates revenue, royalties, rent or income from our content or its inclusion as part of a paid for service or for other commercial gain. Springer Nature journal content cannot be used for inter-library loans and librarians may not upload Springer Nature journal content on a large scale into their, or any other, institutional repository.

These terms of use are reviewed regularly and may be amended at any time. Springer Nature is not obligated to publish any information or content on this website and may remove it or features or functionality at our sole discretion, at any time with or without notice. Springer Nature may revoke this licence to you at any time and remove access to any copies of the Springer Nature journal content which have been saved.

To the fullest extent permitted by law, Springer Nature makes no warranties, representations or guarantees to Users, either express or implied with respect to the Springer nature journal content and all parties disclaim and waive any implied warranties or warranties imposed by law, including merchantability or fitness for any particular purpose.

Please note that these rights do not automatically extend to content, data or other material published by Springer Nature that may be licensed from third parties.

If you would like to use or distribute our Springer Nature journal content to a wider audience or on a regular basis or in any other manner not expressly permitted by these Terms, please contact Springer Nature at

onlineservice@springernature.com