# CertiBase: a genetic database for cultivar certification and genetic breeding of pecan in Brazil

**Tales Poletto[1*], Igor Poletto[2], Cassiano Eric de Carvalho Marques[1], Ana Kelly de Sousa Silva[1], Keith Kubenka[3], Warren Chatwin[3], Oliver Gailing[4] and Valdir Marcos Stefenon[1]**

**Abstract:** *This study was aimed at developing a molecular genetic database that allows the certification of the pecan cultivars planted in Brazil, based on reference cultivars from the USA. We used a set of 10 nSSR and 10 cpSSR markers for cultivar identification and characterization. The database developed is an instrument for certification of cultivars, characterization of potential new cultivars, and selection of material for pecan breeding. The CertiBase database is composed of genotypes and haplotypes characterized for 40 reference pecan cultivars, while the CertiBase algorithm for plant identification is implemented in a user-friendly free software. Varietal certification of cultivars will bring more security to the entire production chain, in addition to enabling the characterization of the ancestry of potential new cultivars and establishment of breeding programs of pecan in Brazil.*

**\*Corresponding author:**
E-mail: talespoletto@gmail.com
*ORCID: 0000-0002-6162-4445*

[1] Universidade Federal de Santa Catarina, Centro de Ciências Agrárias, Rodovia Admar Gonzaga, 1346, Itacorubi, 88034-000, Florianópolis, SC, Brazil
[2] Universidade Federal do Pampa, Campus São Gabriel, Rua Aluízio Barros Macedo, s/n, BR 290, km 423, 97307-020, São Gabriel, RS, Brazil
[3] USDA - ARS Southern Plains Agricultural Research Center, 2881 F&B Road, College Station, TX 77845, United States
[4] University of Göttingen, Department of Forest Genetics and Forest Tree Breeding, Buesgenweg 2, 37077 Goettingen, Niedersachsen, Germany

## INTRODUCTION

Pecan [*Carya illinoinensis* (Wangenh.) K. Koch] is a fruit crop tree native to North America and cultivated in several countries around the world. Brazil is the fifth largest producer on the globe (INC 2023), planting cultivars imported from the USA. The commercial interest in pecan cultivation in Brazil started in the 1920s, with the introduction of the first cultivars in the country. New cultivars were imported from the USA in the 1940s and 1970s, expanding the cultivated area and cultivars' diversity. However, mainly during the early decades of the 20th century, the imported cultivars were cultivated and propagated without a strict guideline system. Consequently, the control of this material was lost, resulting in the commercialization of synonymous (different names for the same cultivar) and homonymous (different cultivars with the same name) cultivars. The same problem might occur also in established orchards, generating uncertainty for producers due to the lack of varietal certification of commercialized cultivars. Such hypothesis of synonymous and homonymous cultivars was recently supported through molecular analysis of plastid SSR markers (Nagel et al. 2023).

Moreover, since pecan's introduction in the country, open pollination among pecan cultivars occurred and farmers have carried out its propagation (part of the propagated plants were generated through seeds) and selection, observing

morpho-agronomic characteristics. This process resulted in several genotypes with unknown ancestry, but adapted to Brazilian ecological conditions, with great potential for registration as new cultivars (Poletto et al. 2020, Oliveira et al. 2023).

In this context, this study was aimed at developing a molecular genetic database that, based on reference cultivars imported from the USDA-ARS National Collection of Genetic Resources for Pecans and Hickories, allows the certification of the pecan cultivars planted in Brazil and the characterization of plants with potential for registration as new cultivars or importance for establishing germplasm collections for species breeding.

## MATERIAL AND METHODS

### Plant material

Samples of 40 cultivars were imported from the USDA-ARS National Collection of Genetic Resources for Pecans and Hickories (College Station and Brownwood, Texas), which is responsible for the breeding, stewardship, and maintenance of the pecan cultivars in the USA. These plants are the reference genotypes for our study. Leaf samples of adult trees representing 11 pecan cultivars commercially cultivated in Brazil (Barton, Desirable, Success, Jackson, Shawnee, Chickasaw, Farley, Stuart, Choctaw, Mahan, and Moneymaker) were collected in several orchards in the Rio Grande do Sul state. The pecan cultivars planted in Brazil were represented by two to six samples from each cultivar collected from different orchards, totalizing 35 samples. Since this cultivated material is commercialized by a few nurseries that vegetatively propagate the plantlets through grafting, this sampling is expected to represent the main genotypes of cultivars planted in Brazil. Additionally, 17 accessions without defined cultivar determination showing traits of interest for genetic breeding of pecan in Brazil (such as disease resistance, high yield, and fruit morphology) were sampled for molecular characterization and comparison with the reference cultivars.

### Molecular genetic analyses

Total DNA was isolated from all 92 samples using the CTAB method (Doyle and Doyle 1990). A first sample set included the 35 Brazilian plants and the 11 respective reference cultivars genotyped for 10 nuclear SSR (nSSR) (Grauke et al. 2003) and 10 chloroplast SSR (cpSSR) markers (Nagel et al. 2020, Nagel et al. 2023). A second genotyping involved all 40 reference cultivars and the 17 Brazilian plants without defined cultivar identification evaluated only for the nSSR markers. Both markers were amplified using the fluorescently labeled M13 tail method (Schuelke 2000), as described by Nagel et al. (2023). Amplified alleles were resolved through 50 cm capillary electrophoresis on an ABI 3500xL DNA Sequencer (ThermoFisher Scientific). The allele calling and the electropherograms' visualization were performed using the GeneMapper™ software (ThermoFisher Scientific). All automatically registered alleles were manually verified to identify and correct mismarked peaks.

The number of alleles (*A*) and the probability of identity (*PI*) of the nSSR and cpSSR individually and of the combined sets of markers were estimated with GenAlEx 6.4 (Peakall and Smouse 2012). The *PI* analysis indicates the statistical power of a set of markers for genetic tagging. Dendrograms were constructed in PAST 4.11 (Hammer et al. 2001) based on the Euclidean distance among samples using the concatenated genetic datasets (for the first sample set) or the Mahalanobis nSSR dataset (for the second sample set) and the neighbor-joining algorithm with 1000 bootstrap resampling.

### Building the CertiBase algorithm

Based on the generated genotypes of nSSR markers and haplotypes of the cpSSR markers, a software based on python script was built aiming at performing comparisons among the cultivars of the database (references) and particular pecan trees (query). The algorithm performs the pairwise comparison of the query to each reference and estimates the genetic similarity (*S*) based on the mean absolute percentage error of the alleles at each locus as $S = \frac{1}{N} \sum_{i=1}^{N} max \left( 0,1 - \frac{|q - r|}{r} \right)$, where *q* and *r* are the allele size, in base pairs, of the query and the reference, respectively. The '*max*' function is used to limit *S* estimations between 0.0 and 1.0, avoiding negative estimations of *S* when the query is compared with cultivars genetically distant showing relatively large differences in base pairs between alleles. For example, if the query is a sample of cultivar Barton, the comparison with the reference Barton will return values of *S* near 1.0 since small differences in allele sizes are accepted. However, the comparison of this query with another cultivar may return negative values for *S* for some loci, due to large differences in allele sizes. The CertiBase algorithm will perform the comparison among

genotypes/haplotypes of the query plants and the genotypes/haplotypes of the reference cultivars, returning a list of the three references with higher similarity to the query.

## RESULTS AND DISCUSSION

The use of SSR markers to aid in characterizing pecan accessions in bank germplasms and to support breeders in deciding on parental lines for hybrid crosses was established in the USDA Pecan Breeding Program in the USA (Grauke et al. 2015). Some years later, similar studies employed universal AFLP and S-SAP markers in Brazil (Poletto et al. 2020, Oliveira et al. 2023), characterizing selected genotypes of pecan cultivated in private orchards. Moreover, the sequencing of the plastid genome of cultivar Imperial (Nagel et al. 2020), one of the most planted cultivars in Brazil, enabled the development of species-specific cpSSR markers for pecan, which were shown to be useful for parental determination and characterization of misassigned cultivars (Nagel et al. 2023). In the present study, we intended to advance these studies using cpSSR and nSSR towards the analysis of misassignments, identification of genotypes without known origin, and characterization of genotypes with potential for registration as new cultivars.

For the comparison among the cultivars planted in Brazil and the correspondent reference cultivars from the USA, a total of 53 alleles (mean $A$ = 5.3, ranging from 2 to 12) were identified in the nSSR, and 20 alleles (mean $A$ = 2.0, ranging from 1 to 3) in the cpSSR data set (Table 1). According to the Probability of Identity analysis, the chance of misassigning a sample from a given cultivar is $PI$ = 0.018% using 10 cpSSRs, $PI$ = 0.000001% using 10 nSSRs, and $PI$ = 0.00000002% using both sets of markers combined (Figure 1). Given the higher polymorphism of the nSSR loci, these markers are more informative for the identification of the pecan cultivars. However, cpSSRs have been shown to be informative for the identification of pecan cultivars sharing the mother cultivar, given the exclusively maternal inheritance of the chloroplast and, consequently, of these markers (Nagel et al. 2023). Thus, the combination of both markers significantly increases the power and usefulness of the analysis.

Different sorts of molecular markers have been used for the characterization of cultivated and wild pecan germplasms (Oliveira et al. 2021), including RAPD (Elarabi and Elsoda 2017), AFLP (Poletto et al. 2020, Oliveira et al. 2023), S-SAP (Oliveira et al. 2023), ISSR (Elarabi and Elsoda 2017), nSSR (Grauke et al. 2011, Wang et al. 2022), cpSSR (Grauke et al. 2010, Nagel et al. 2023), SNPs (Bentley et al. 2019), and InDels (Mo et al. 2024). Each of such markers has advantages and drawbacks concerning genome coverage, polymorphism, dominance, informativeness, effort, and costs. While nSSR markers are highly polymorphic and mostly non-linked, the maternal heritability of cpSSRs makes these markers very informative for verifying the directionality of crosses and the parental structure of accessions of unknown origin with the potential to be registered as new cultivars (Grauke et al. 2015, Nagel et al. 2023).

### Genetic consistencies and misassignments of the Brazilian-grown cultivars

According to reports of the farmers from whom the Brazilian samples were collected, most cultivars grown in Brazil are morphologically consistent with the correspondent references, but some incongruencies have been described. In our genetic analysis with 10 nSSR and 10 cpSSR markers, all samples from cultivars Barton, Success, Jackson, Shawnee, Chickasaw, Stuart, Choctaw, Mahan, and Moneymaker clustered with the respective North American references in the dendrogram, with bootstrap supports higher than 61% (Figure 1), corroborating the morphological correspondence. On the other hand, two samples of cultivar Farley and one sample of cultivar Desirable failed to cluster with the respective references (Figure 1). This mismatch of samples of cultivars Desirable and Farley to the North American references in the genetic analysis supports the hypothesis of occurrence of cultivar misassignment in Brazil, at least in some orchards.

**Table 1.** Total number of alleles ($A$) for nuclear (nSSR) and plastid (cpSSR) SSR markers recorded from 35 accessions representing 11 pecan cultivars planted in Brazil

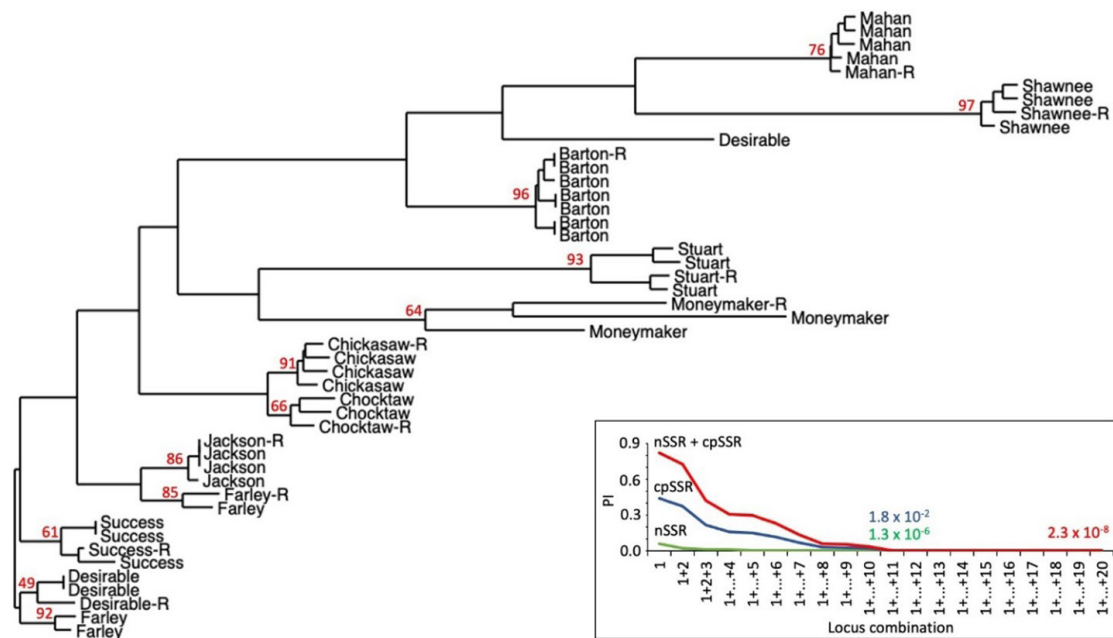| nSSR | | cpSSR | |
|---|---|---|---|
| **Marker** | **N. alleles** | **Marker** | **N. alleles** |
| pm-cin4 | 12 | cpCil1 | 2 |
| pm-ca10 | 2 | cpCil2 | 2 |
| pm-cin13 | 5 | cpCil4 | 2 |
| pm-cin20 | 2 | cpCil5 | 2 |
| pm-cin22 | 5 | cpCil6 | 1 |
| pm-ga23 | 5 | cpCil7 | 2 |
| pm-cin27 | 6 | cpCil8 | 2 |
| pm-ga38 | 5 | cpCil9 | 2 |
| pm-ga39 | 4 | cpCil10 | 2 |
| pm-ga41 | 7 | cpCil12 | 3 |
| Overall mean nSSR | 5.3 | Overall mean cpSSR | 2.0 |
| Total nSSR | 53 | Total cpSSR | 20 |

**Figure 1.** Neighbor-joining dendrogram of the pecan cultivars planted in Brazil and the reference cultivars (names followed by -R), based on the Euclidean distance of the combined genetic datasets (nSSR and cpSSR markers). Numbers at nodes are the bootstrap support after 1000 replicates. The inset shows the Probability of Identity (*PI*) of the SSR markers.

Based on the same 10 cpSSR markers used in this study, the analysis of 24 pecan genotypes including eight commercial cultivars and three accessions without defined identification, putative misassignment of pecan cultivars was reported by Nagel et al. (2023) and further studies including nSSR markers were suggested. Similarly to our study with cpSSR and nSSR markers, samples of cultivar Barton clustered together based on the cpSSR haplotypes, showing point divergence in a single locus, while samples of cultivar Shawnee had 100% similarity to each other (Nagel et al. 2023).

In addition to synonymy and homonymy due to a mix of plants or confusion concerning the interpretation of morphological traits, the observed mismatches may also originate from cross-pollination between cultivars, with the predominance of morphological characteristics of one of the parental trees. The electropherograms of the nSSR marker pm-cin4 demonstrated the matching of one sample of cultivar Barton from Brazil to the reference pattern (Figure 2a) and the mismatch of one sample of cultivar Desirable to the North American reference (Figure 2b). This mismatched sample has alleles 124 and 158 for the pm-cin4 locus and is likely the result of a cross between cultivars Desirable and Shawnee (Figure 2b) since these cultivars hold the alleles with 124 bp (reference cultivar Desirable) and 158 bp (reference cultivar Shawnee). Note that the allele size in the electropherogram (Figure 2) is 18 base pairs longer due to the M13 tail and is discounted to determine the cultivar genotype if using a PCR protocol without such tail in the forward primer.

**Potential of plants without defined cultivar for genetic breeding and new cultivar registration**

The comparison between the 17 Brazilian accessions without cultivar definition and the 40 reference cultivars suggests that some Brazilian plants can be assigned to known cultivars. Considering the clustering analysis, which results from the genetic similarity based on the nuclear markers (Figure 3), samples BG71, BG80, and BG85, for instance, might be Jackson cultivar, BG6, BG49, and BG31 might be Chetopa (Figure 3). While the mean overall distance was $d = 1.36$, the mean distance for these plants to the reference was $d = 0.58$ (BGs to Jackson) and $d = 0.61$ (BGs to Chetopa). The distance over the 40 reference cultivars was $d = 1.39$, while the distance of the cluster BGs/Jackson to the other plants was $d = 1.14$, the distance of the cluster BGs/Chetopa to the others was $d = 1.20$ (Figure 3). On the other hand, several samples do not cluster with the reference cultivars and might result from open pollination. For instance, samples BG39, BG57, and BG79 form an individual cluster (mean

$d$ = 1.37), as occurs with samples BG29 and BG66 ($d$ = 1.72), and BG97 ($d$ = 1.98) (Figure 3). Thus, the genetic differentiation among these pecan trees and known pecan cultivars is equivalent to the genetic differentiation among different cultivars.

A more complete comparison including morphological traits and plastid SSR markers is necessary to characterize the putative parental cultivars of such pecan accessions, enabling the registration of new cultivars. Moreover, some cultivated genotypes lacking identified cultivars and ancestry have a high potential for future breeding programs of pecan in Brazil (Oliveira et al. 2023). Thus, since these plants hold traits of agronomic interest such as high yield and resistance to pecan scab, they should be included in active germplasm collections of pecan, towards establishing a reference collection for genetic breeding of this crop tree in Brazil. Considering that the expansion of planted areas of pecan in Brazil is limited by the poor performance of commercial cultivars in some regions due to climatic conditions, the development of cultivars adapted to the Brazilian environment is imperative (Nagel et al. 2023).

Moreover, varietal certification of commercial cultivars will bring more security to the entire production chain, from the production of seedlings to the identification of dubious cultivars in orchards, helping to avoid fruit mixing during harvesting. In addition, the genetic and agronomic characterization of selected plants is essential for the establishment of active germplasm banks and genetic breeding programs for this species in Brazil.
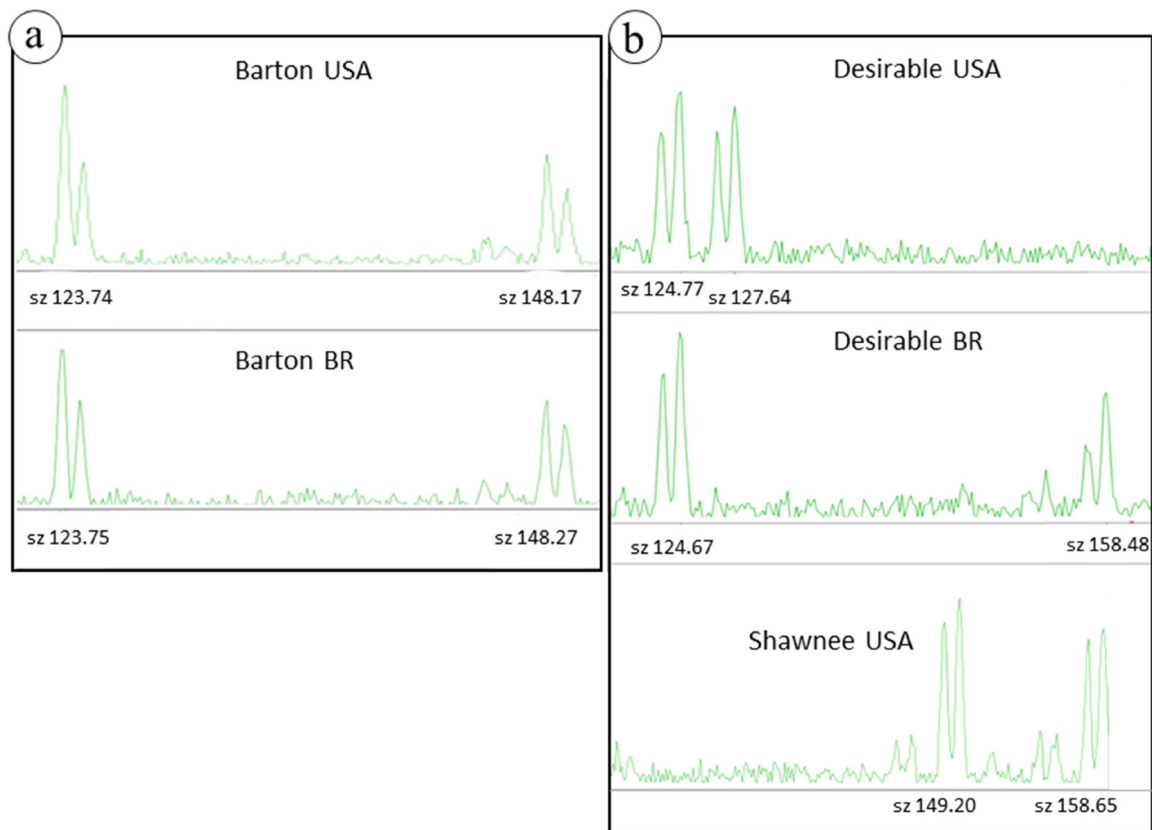


**Figure 2.** Allelic patterns of the nuclear marker pm-cin4 for cultivar certification or parental analysis. **a)** Electropherogram of the North American reference (Barton USA) and one genotype of the Brazilian cultivar Barton, showing a true-to-type pattern. **b)** Electropherogram of the North American reference and one off-type genotype of the Brazilian cultivar Desirable. The Shawnee cultivar seems to be a parent of the off-type Brazilian genotype since the reference cultivar Shawnee USA holds the allele 158.65, equivalent to allele 158.48 in Desirable BR. The size of the alleles in the electropherogram is larger than registered in the database, due to the 18 bp of the M13-tail used in the PCR amplification (see Material and Methods).

## The CertiBase software

The CertiBase software was validated using a query sample previously identified morphologically and molecularly as Barton in this study. This query sample grouped in the neighbor-joining dendrogram with the reference cultivar with a bootstrap support of 96% (Figure 1). Using our database as the entry of the software (Figures 4a-b), the CertiBase algorithm returned 92.5% similarity of the query sample to cultivar Barton, 89.2% to Success, and 87.9% to Jackson (Figure 4c), supporting the expected identification of the sample.

Having verified cultivar collections and established protocols for such verification is required to provide the scientific community with a confident database with methods and tools necessary to characterize and use germplasm collections (Grauke et al. 2015). Literman et al. (2022) proposed a method based on SNPs for identifying pecan cultivars and hybrids among *Carya* species. Although very informative, this method relies on genome sequencing, comprising more complex laboratory skills, equipment, and funds. On the other hand, the CertiBase method relies on tools and skills commonly used in molecular laboratories.

The genetic database developed here is intended to support varietal certification, identification of potential new cultivars, and selection of the most relevant plants to create a base population for genetic breeding of pecan in Brazil. The CertiBase database is composed of genotypes and haplotypes characterized for 40 reference pecan cultivars obtained from the USDA-ARS National Collection of Genetic Resources for Pecans and Hickories, while the CertiBase algorithm is implemented in a user-friendly software in which the handler will choose the database (an excel sheet with the genotypes/haplotypes of each reference cultivar) and the query (an excel sheet with the genotypes/haplotypes of the plants intended to be evaluated) files (Figure 4a-b). These files are simple Excel sheets with a previously determined sequence of the nSSR and cpSSR markers assessed. The user must list the genotypes/haplotypes of the query plants in the sheet, in the same sequence that it is in the database sheet (described in the users' manual of the CertiBase), enabling the correct comparison. Finally, the software will run the algorithm and compute the similarities, returning the result with the proportion of similarity of the query to the three most similar reference cultivars of the database (Figure 4c). This program can also be used for other species given that the user creates its own database. The CertiBase software and database sheet are freely accessible for download at https://github.com/Yugaren/LFDGV/tree/main/CertiBase/CertiBase and can be run under Windows®, Linux®, or iOS® operational systems.
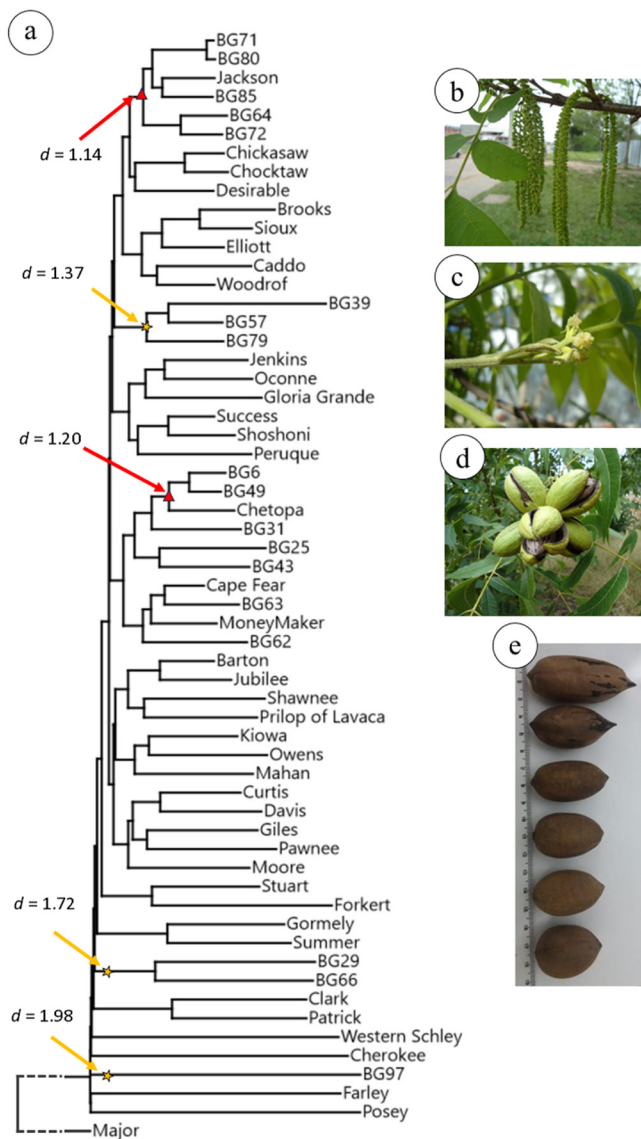


**Figure 3.** Molecular characterization of pecan accessions without cultivar definition. **a)** Neighbor-joining dendrogram based on the Mahalanobis distance of the 40 reference cultivars and the 17 Brazilian accessions without cultivar definition (samples with BG prefix). The mean pairwise distance among all samples (reference cultivars and accessions without cultivar definition) is *d* = 1.36. **b)** Staminate flowers of cultivar Barton. **c)** Pistillate flowers of cultivar Barton. **d)** Fruits of cultivar Barton. **e)** Morphological diversity of pecan nuts cultivated in Brazil (from the top: Mahan, Sumner, Choctaw, Pawnee, Desirable, Shoshoni).
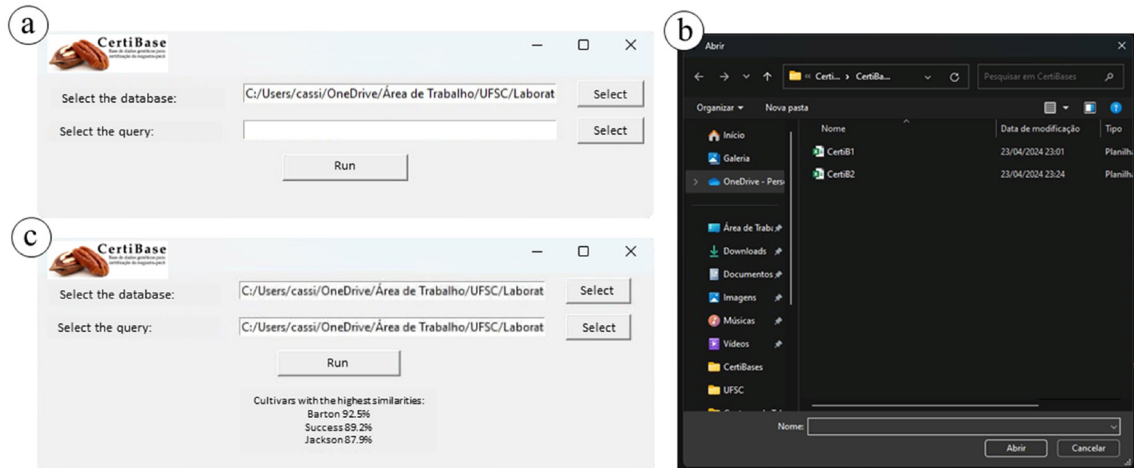
**Figure 4.** Overall view of the CertiBase software. **a)** The user must choose the database and the query files (excel sheets) using the "Select" button. **b)** A search window will open for file choice. **c)** After uploading both files (database and query sheets), the similarity is computed and presented in the lower part of the window for the three cultivars with the highest similarities to the query plant.

## ACKNOWLEDMENTS

## DATA AVAILABILITY

The datasets generated and/or analyzed during the current research are available for download at https://github.com/Yugaren/LFDGV/tree/main/CertiBase/CertiBase.

## REFERENCES

Bentley N, Grauke LJ and Klein P (2019) Genotyping by sequencing (GBS) and SNP marker analysis of diverse accessions of pecan (*Carya illinoinensis*). **Tree Genetics & Genomes 15**: 8.

Doyle JJ and Doyle JL (1990) Isolation of plant DNA from fresh tissue. **Focus 12**: 39-40.

Elarabi NI and Elsoda AS (2017) Molecular evaluation of genetic diversity among seven genotypes of pecan (*Carya Illinoinensis*). **Journal of Agricultural Chemistry and Biotechnology 8**: 27-33.

Grauke LJ, Iqbal MJ, Reddy AS and Thompson TE (2003) Developing microsatellite DNA markers in pecan. **Journal of the American Society for Horticultural Science 128**: 374-380.

Grauke LJ, Klein R, Grusak MA and Klein P (2015) The Forest and the trees: Applications for molecular markers in the repository and pecan breeding program. **Acta Horticulturae 1070**: 109-126.

Grauke LJ, Mendoza-Herrera MA and Binzel ML (2010) Plastid microsatellite markers in *Carya*. **Acta Horticulturae 859**: 237-246.

Grauke LJ, Mendoza-Herrera MA, Miller AJ and Wood BW (2011)

Geographic patterns of genetic variation in native pecans. **Tree Genetics & Genomes 7**: 917-932.

Hammer Ø, Harper DAT and Ryan PD (2001) Past: Paleontological statistics software package for education and data analysis. **Palaeontologia Electronica 4**: 4.

INC - International Nut and Dried Fruit (2023) **Nuts & dried fruits: statistical yearbook 2022/2023.** International Nut and Dried Fruit, Reus, 79p.

Literman RA, Ott BM, Wen J, Grauke LJ, Schwartz RS and Handy SM (2022) Reference-free discovery of nuclear SNPs permits accurate, sensitive identification of Carya (hickory) species and hybrids. **Applications in Plant Sciences 10**: e11455.

Mo Z, Lou W, Zhang Y, Hu L, Zhai M and Xuan J (2024) Insertion/deletion variant characterization and marker development in pecan [*Carya illinoinensis* (Wangenh.) K. Koch]. **Scientia Horticulturae 325**: 112660.

Nagel J, Machado LO, Lemos RPM, Matielo CBD, Poleto T, Poletto I and Stefenon VM (2020) Structural, evolutionary and phylogenomic features of the plastid genome of *Carya illinoinensis* cv. Imperial. **Annals of Forest Research 63**: 3-18.

Nagel J, Poleto T, Muniz MFB, Poletto I, Oliveira JNM and Stefenon VM (2023) Species-specific plastid SSR markers reveal evidence of cultivar misassignments in Brazilian pecan [*Carya illinoinensis* (Wangenh.) K. Koch] orchards. **Genetic Resources and Crop Evolution 70**: 971-980.

Oliveira LO, Beise DC, Santos DD, Nagel JC, Poletto T, Poletto I and Stefenon VM (2021) Molecular markers in *Carya illinoinensis* (Juglandaceae): from genetic characterization to molecular breeding, **The Journal of Horticultural Science and Biotechnology 96**: 560-569.

Oliveira LO, Santos DD, Beise DC, Poleto T, Poetto I, Muniz MFB, Oliveira JNM and Stefenon VM (2023) Old but still good: genetic diversity of ancient pecan genotypes from southern Brazil. **Anais da Academia Brasileira de Ciências 95**: e20220885.

Peakall R and Smouse PE (2012) GenAlEx 6.5: genetic analysis in Excel. Population genetic software for teaching and research - an update. **Bioinformatics 28**: 2537-2539.

Poletto T, Poletto I, Silva LMM, Muniz MFB, Reiniger LRS, Richards N and Stefenon VM (2020) Morphological, chemical and genetic analysis of southern Brazilian pecan (*Carya illinoinensis*) accessions. **Scientia Horticulturae 262**: 108863.

Schuelke M (2000) An economic method for the fluorescent labelling of PCR fragments. **Nature Biotechnology 18**: 233-234.

Wang X, Chatwin W, Hilton A and Kubenka K (2022) Genetic diversity revealed by microsatellites in genus Carya. **Forests 13**: 188.