



OPEN

A divide-and-conquer approach for genomic prediction in rubber tree using machine learning

Alexandre Hild Aono^{1,2}, Felipe Roberto Francisco¹, Livia Moura Souza^{1,3}, Paulo de Souza Gonçalves⁴, Erivaldo J. Scaloppi Junior⁴, Vincent Le Guen^{5,6}, Roberto Fritsche-Neto⁷, Gregor Gorjanc², Marcos Gonçalves Quiles⁸ & Anete Pereira de Souza^{1,9}✉

Rubber tree (*Hevea brasiliensis*) is the main feedstock for commercial rubber; however, its long vegetative cycle has hindered the development of more productive varieties via breeding programs. With the availability of *H. brasiliensis* genomic data, several linkage maps with associated quantitative trait loci have been constructed and suggested as a tool for marker-assisted selection. Nonetheless, novel genomic strategies are still needed, and genomic selection (GS) may facilitate rubber tree breeding programs aimed at reducing the required cycles for performance assessment. Even though such a methodology has already been shown to be a promising tool for rubber tree breeding, increased model predictive capabilities and practical application are still needed. Here, we developed a novel machine learning-based approach for predicting rubber tree stem circumference based on molecular markers. Through a divide-and-conquer strategy, we propose a neural network prediction system with two stages: (1) subpopulation prediction and (2) phenotype estimation. This approach yielded higher accuracies than traditional statistical models in a single-environment scenario. By delivering large accuracy improvements, our methodology represents a powerful tool for use in *Hevea* GS strategies. Therefore, the incorporation of machine learning techniques into rubber tree GS represents an opportunity to build more robust models and optimize *Hevea* breeding programs.

Rubber tree (*Hevea brasiliensis*) has an elevated importance in the global economy, being almost the only feedstock for commercial rubber^{1,2}. Considering the long perennial vegetative cycle of *Hevea*, breeding programs aim to improve its yield production in order to reach the rapidly increasing rubber demand¹⁻³. Therefore, genomic approaches are needed in rubber tree breeding, especially considering its recent domestication history⁴. *H. brasiliensis* is a diploid species ($2n = 36$) with an elevated occurrence of duplicated regions in its genome ($\sim 70\%$)⁵⁻⁷, and this complex genomic organization has hindered the development of genomic strategies for breeding. However, with the improvement of next-generation sequencing (NGS) technologies and the consequent reduction in genotyping costs, data generation has become more efficient, providing more genomic resources in less time and with lower associated costs⁸. This greater availability of data improved precision in selection with higher genetic gains in various crops^{8,9} and, in rubber tree, could complement traditional approaches based on only phenotypic and pedigree information^{8,10}.

Various rubber tree genomic resources have become available in recent decades, such as a large set of different molecular markers¹¹⁻¹⁴, draft genomes^{5,6}, and, more recently, a chromosome-level assembled genome⁷. These data have already allowed the construction of saturated linkage maps with associated quantitative trait loci (QTLs), which were proposed as a tool for marker-assisted selection (MAS)¹⁵. Although QTLs for several traits have been identified in rubber tree^{4,15-20}, the amount of phenotypic variance explained by these identified QTLs is usually

¹Molecular Biology and Genetic Engineering Center (CBMEG), University of Campinas (UNICAMP), Campinas, SP, Brazil. ²The Roslin Institute and Royal (Dick) School of Veterinary Studies, The University of Edinburgh, Midlothian, UK. ³São Francisco University (USF), Itatiba, Brazil. ⁴Center of Rubber Tree and Agroforestry Systems, Agronomic Institute (IAC), Votuporanga, Brazil. ⁵Centre de Coopération Internationale en Recherche Agronomique pour le Développement (CIRAD), UMR AGAP, 34398 Montpellier, France. ⁶AGAP, CIRAD, INRAE, Institut Agro, Univ Montpellier, Montpellier, France. ⁷Genetics Department, Luiz de Queiroz College of Agriculture (ESALQ), University of São Paulo (USP), Piracicaba, SP, Brazil. ⁸Instituto de Ciência e Tecnologia (ICT), Universidade Federal de São Paulo (UNIFESP), São José dos Campos, SP, Brazil. ⁹Department of Plant Biology, Institute of Biology (IB), University of Campinas (UNICAMP), Campinas, SP, Brazil. ✉email: anete@unicamp.br

small¹⁹ because of the highly complex genetic architectures associated with growth and rubber production traits. The configuration of these phenotypes is controlled by many genes with small effects²¹, and weak QTLs may not be identified using existing methodologies^{2,22}, which prevents the identification of interindividual differences²³. Together with the environmental and genetic background restrictions of QTLs²⁴, these features limit the application of *Hevea* QTLs for MAS¹⁴. Consequently, novel genomic strategies that can assist in rubber tree breeding programs are needed, especially considering the time required to evaluate these phenotypes, the elevated costs, and the low female fertility in *H. brasiliensis*^{2,15,25}.

Aimed at solving such difficulties in many crops, genomic selection (GS) has arisen as a promising methodology for considerably reducing the required breeding cycle²⁶. GS has shown better performance than MAS^{27,28}, mainly because of its associated genetic gains²⁹ and reduced costs over a long time period³⁰. This strategy enables the selection of plants based on their estimated performance obtained with a large dataset of molecular markers^{8,31}, reducing breeding time by avoiding the need to evaluate a considerable number of phenotypes over different years²⁴. Using known phenotypic and genotypic information from a training population³², it is possible to create a predictive model that can be used to predict the breeding values of a testing population using only genotypic data⁸. This modeling is generally based on a mixed-effect regression method³³ and has already been demonstrated to be promising for several crops^{34–38}. In rubber tree²⁵, and² assessed the potential of GS for predicting stem circumference (SC) and rubber production (RP), respectively, simulating breeding schemes through cross-validation (CV) techniques.

There are several CV approaches for simulating a real application of GS in a plant breeding program. These methods take into account the population structure in the dataset and the appropriateness of applying the developed predictive model to a set of plants. There are basically three CV schemes in GS: (1) predicting traits in an untested environment using previously tested lines (CV0)⁸, (2) predicting new lines' traits that were not evaluated in any environment (CV1)³⁹, and (3) predicting traits that were evaluated in some environments but not in others (CV2)⁴⁰. These three scenarios were already evaluated in rubber tree² assessed the potential of GS in a within-family context using CV0 and CV1 methods, and²⁵ tested three different populations with CV1 and CV2. These initiatives represent the first attempts to use GS on rubber tree data, but with low associated predictive capabilities for some of the created CV schemes, mostly when prediction is performed with genotypes that have not already been tested.

Different approaches have been used in GS to create predictive models, including parametric and non-parametric methods^{24,26,41–45}. Significant differences in predictive capabilities have not been demonstrated when changing the predictive approach^{31,46,47}; thus, linking genotypes and phenotypes remains a great challenge^{23,48}, especially for plant species with high genomic complexity. In this context, more robust techniques for estimating these models with higher prediction capabilities are needed to expand the practical implementation of GS in rubber tree. Nonlinear techniques have already shown improved performance in representing complex traits with nonadditive effects^{9,49–51}, and, in this context, machine learning (ML) strategies have emerged as a promising set of tools for complementing these statistical nonlinear methods.

The objective of this work was to develop a genomic prediction approach for rubber tree data. Considering that ML methods have not been proven to have better performance than statistical methodologies for GS^{23,52}, we evaluated their efficiency in rubber tree, also suggesting a novel approach for constructing a predictive system with neural networks based on two-stage prediction: (1) subpopulation prediction and (2) phenotype estimation. Such a divisive approach was created considering a common paradigm in Computer Science: divide and conquer. For datasets with a clear subpopulation structure, such as rubber tree, the proposed approach represents a promising alternative for the development of predictive models.

Material and methods

Plant material and phenotypic characterization. The data used in this work were obtained with different experiments in two previous studies. The plant material and permissions for collecting rubber tree employed in the present study are in compliance with institutional, national, and international guidelines and legislation. Therefore, our analyses were conducted by separating the methodologies and considering two datasets: experimental group 1 (EG1) and experimental group 2 (EG2). EG1 includes 408 samples of three F_1 segregant populations obtained with crosses between (Pop1) GT1 and PB235 (30 genotypes)²⁵, (Pop2) GT1 and RRIM701 (127 genotypes)^{25,53}, and (Pop3) PR255 and PB217 (251 genotypes)^{4,19,25}. EG2 is based on an F_1 cross between RRIM600 and PB260 (330 samples)².

The parents of the crosses used are important clones for rubber tree breeding programs. PR255, PB235, PB260, and RRIM600 have high yield, and PB217 has considerable potential for long-term yield performance due to its slow growth process^{2,25}. PR255 and RRIM701 have good growth, and RRIM701 also presents an increased SC after initial tapping⁵⁴. The latex production is stable in PR255 and medium in RRIM600. Stable or medium latex production represents a good adaptation to several environments, as observed in GT1, a clone tolerant to wind and cold. Additionally, PB260 presents high female fertility⁵⁵, and PB235 is susceptible to tapping panel dryness⁵⁶.

In EG1 and EG2, we analyzed the SC trait. In EG1, Pop3 was planted in 2006 in a randomized block design in Itiquira, Mato Grosso State, Brazil, 17°24' 03" S and 54°44' 53" W^{4,19,25}. Each individual was represented by four grafted trees in each plot and four replications. Pop1 and Pop2 were planted in 2012 at the Center of Rubber Tree and Agroforestry Systems/Agronomic Institute (IAC - Brazil), 20°25' 00" S and 49°59' 00" W, following an augmented block design, with four blocks containing two clones per plot spaced 4 m apart for each trial, which was repeated four times^{25,53}.

Even though EG2 corresponds to only one cross, this population was planted following an almost complete block design at two different sites², which for convenience we named site 1 (S1) and site 2 (S2). In S1, 189 clones were planted in 2012 in Société des Caoutchoucs de Grand-Béréby (SOGB—Ivory Coast), 4° 40' 54" N and 7°

06' 05" W. In S2, 143 clones were planted in 2013 in Société Africaine de Plantations d'Hévéas (SAPH - Ivory Coast), 5° 19' 47.79" N and 4° 36' 39.74" W. This cross consisted of six blocks with randomized trees spaced 2.5 m apart and a mean number of ramets per clone of 11 for S1 (ranging between 7 and 17) and 13 for S2 (ranging between 5 and 20).

SC measurements of Pop3 in EG1 were obtained in four years (from 2007 to 2010) and those of Pop1 and Pop2 were obtained from 2013 to 2016, considering that growth traits are usually measured only during the first 6 years^{25,57}. According to the water distribution of the experiments installed, EG1 phenotypes were measured to supply information considering low-water (LW) and well-watered (WW) conditions; thus, Pop3 was evaluated in October 2007–2010 (LW) and in April 2008–2010 (WW), and Pop1 and Pop2 were evaluated in June 2013, December 2013, May 2014, November 2014, and June 2015–2016. SCs were measured for individual trees at 50 cm above ground level. For both phenotypes, the average per plot was calculated. SC in EG2 was measured at 1 m above ground level before tapping for 3 months every two days except on Sundays (with the beginning at 32 months after planting in S1 and 38 months after planting in S2).

Phenotypic data analysis. All phenotypic analyses were performed using R statistical software⁵⁸. EG1 and EG2 traits were analyzed with the following steps: (1) data distribution evaluation; (2) standardized normalization with the R package *bestNormalize*⁵⁹; (3) mixed-effect model creation and residual appropriateness verification through quantile-quantile (Q-Q) plots using the *breedR* package⁶⁰; (4) estimation of best linear unbiased predictions (BLUPs) based on the models created; (5) hierarchical clustering on BLUP values using a complete hierarchical clustering approach based on Euclidean distances and dendrogram visualization with the *ggtree* R package⁶¹; and (6) identification of phenotypic groups using the clustering approach of (5), with cluster numbers ranging between 2 and 5, and several clustering indexes implemented in the *NbClust* R package⁶².

In EG1, we employed the following statistical mixed-effect model:

$$Y_{ijk} = \mu + L_k + B_{jk} + W + G_{ik} + e_{ijk} \quad (1)$$

where Y_{ijk} corresponds to the phenotype of the i th genotype in the j th block and k th location. The phenotypic mean is represented by μ , and the fixed effects represent the contribution of the k th location (L_k), the j th block at the k th location (B_{jk}), and the watering condition of the measurement (W). The genotype G and the residual error e (nongenetic effects) represent the random effects.

EG2 SC phenotypes were modeled for each site (S1 and S2) according to the following statistical model:

$$Y_{ijk_r} = \mu + B_j + L_{kj} + R_{r_{kj}} + G_{ij} + e_{ijk_r} \quad (2)$$

where Y_{ijk_r} corresponds to the phenotype of the i th genotype positioned in the r th rank of the k th line in the j th block. The phenotypic mean is represented by μ , and the fixed effects represent the contribution of the j th block (B_j), the k th line of the j th block (L_{kj}), and the r th rank of the k th line in the j th block ($R_{r_{kj}}$). The genotype G and the residual error e (nongenetic effects) represent the random effects. Broad-sense heritability (H^2) was estimated as $H^2 = \sigma_g^2 / \sigma_p^2$, with σ_g^2 and σ_p^2 representing the genetic and phenotypic variances, respectively.

Genotyping process. DNA extraction from EG1 was described by^{19,53}, and the genotyping process was performed using a genotyping-by-sequencing (GBS) protocol⁶³ with *EcoT22I* restriction enzyme followed by Illumina sequencing using the HiSeq platform for Pop3 and the GAIIX platform for Pop1 and Pop2²⁵. EG1 genotype data analysis was performed as described by²⁵. In summary, raw sequencing reads were processed using the TASSEL 5.0 pipeline⁶⁴, with a minimum count of 6 reads for creating a tag. The tag mapping process was performed using Bowtie2 v.2.1⁶⁵ with the *very sensitive* algorithm and *H. brasiliensis* reference genome⁷. Single nucleotide polymorphisms (SNPs) were called with the TASSEL algorithm, and only biallelic SNPs were retained using VCFtools⁶⁶. These markers were filtered using the R package *snpReady*⁶⁷ with a maximum of 20% missing data for a SNP and 50% in an individual and a minimum allele frequency (MAF) of 5%. Missing data were imputed using the k-nearest neighbors⁶⁸ algorithm implemented in the *snpReady* package.

EG2 samples were genotyped with simple sequence repeat (SSR) markers, following the protocol for DNA extraction and genotyping described by⁶⁹. EG2 genotype data analysis was performed as described by². In summary, a total of 332 SSRs were used for S1²⁰ and 296 for S2². Missing data were imputed using BEAGLE 3.3.2⁷⁰ with 25 iterations of the phasing algorithm and 20 haplotype pairs to sample for each individual in an iteration. For evaluating the genotypic profile of individuals in EG1 and EG2, we performed principal component analyses (PCAs) in R statistical software⁵⁸ with the *ggplot2* package⁷¹. Additionally, for evaluating the overall correspondences between genotypic and phenotypic data, we colored the PCA scatter plots with the BLUPs estimated for SC trait, as performed by⁷².

Statistical models for genomic prediction. We employed two different strategies for creating traditional genomic prediction models: Bayesian ridge regression (BRR)⁷³ and a single-environment, main genotypic effect model with a Gaussian kernel (SM-GK)⁷⁴. BRR and SM-GK models were implemented in the BGLR⁷⁵ and BGGE⁷⁶ R packages, respectively. Considering the genotype matrix with n individuals and p markers, BRR models were implemented considering the following:

$$y = 1\mu + Z\gamma + e \quad (3)$$

where y represents the BLUP values calculated based on the established mixed-effect models for phenotypic data analyses, μ the overall mean, Z the genotype matrix, e the residuals, and γ the vector of marker effects. In

SM-GK, Z is the incidence matrix of genetic effects, and γ is the vector of genetic effects with variance estimated through a Gaussian kernel calculated using the `snprReady` R package.

Genomic prediction via machine learning. For genomic prediction via ML, we selected the following algorithms considering a regression task: (a) AdaBoost⁷⁷, (b) multilayer perceptron (MLP) neural networks⁷⁸, (c) random forests⁷⁹, and (d) support vector machine (SVM)⁸⁰. To create these models, we used Python v.3 programming language together with the library `scikit-learn v.0.19.0`⁸¹. We also tested a combination of feature selection (FS) techniques for increasing the predictive accuracies⁸², using a combination of three different methods: (i) L1-based FS through an SVM model⁸⁰, (ii) univariate FS with Pearson correlations (and ANOVA for discrete variables) (p-value of 0.05), and (iii) gradient tree boosting⁸³. Such a strategy is based on marker subset selection, separating the markers identified by all of these methods together (intersection of the 3 approaches, named Inter3) or by at least two of them simultaneously (Inter2), and using such subsets for prediction.

To understand the subset selection, we performed functional annotation of the genomic regions underlying these markers selected through FS considering a 10,000 base-pair (bp) window for the up- and downstream regions. Using BLASTn software⁸⁴ (minimum e-value of 1e-6), these sequences were aligned against coding DNA sequences (CDSs) from the *Malpighiales* clade (*Linum usitatissimum* v1.0, *Manihot esculenta* v8.1, *Populus deltoides* WV94 v2.1, *Populus trichocarpa* v4.1, *Ricinus communis* v0.1, and *Salix purpurea* v5.1) of the Phytozome v.13 database⁸⁵. On the basis of significant correspondence, Gene Ontology (GO) terms⁸⁶ were retrieved.

Multilayer perceptron neural network. As the final approach for genomic prediction in EG1, we proposed the creation of neural networks with novel architectures for each of the biparental populations, using the Keras Python v.3 library for this task⁸⁷. We employed MLP networks, which have an architecture based on multiple layers and feedforward signal propagation⁸⁸.

For all the predictive tasks, we considered an MLP structure with two hidden layers (HLs) and used the mean absolute error (MAE) as the error function for training and defining the architecture of the networks. Additionally, 200 epochs were considered (batch size of 16). The training process of the networks was performed using the backpropagation strategy together with the Adam optimization algorithm⁸⁹, which aims to minimize the MAE by updating the synaptic weights using a gradient-based strategy that combines heuristics from a momentum term and RMSProp⁹⁰. The update process is based on a change of Δw_{ij} for each connection, considering the individual influence of a weight w_{ij} on the MAE value obtained with the gradient descent g_t in the iteration t calculated with $\partial MAE / \partial w_{ij}$ and used in the equation

$$\Delta w_{ij} = g_t \times \eta \frac{v_t}{\sqrt{s_t + \epsilon}} \quad (4)$$

where η is the learning rate representing the amount of change in the process of training, v_t is the exponential average of gradients along the weights w_i of layer i , and s_t is the exponential average of squares of gradients along w_i . The Adam optimizer employs two other hyperparameters for the optimization process (β_1 and β_2), which are used for the calculation of v_t ($v_t = \beta_1 \times v_{t-1} - (1 - \beta_1) \times g_t$) and s_t ($s_t = \beta_2 \times s_{t-1} - (1 - \beta_2) \times g_t^2$). We used $\beta_1 = 0.9$ and $\beta_2 = 0.999$ ⁸⁹. We tested the following configurations for the MLP hyperparameters: (a) number of neurons in the first HL, varying from 1 to $\sqrt{(q+2)m} + 2\sqrt{m/(q+2)}$ (m individuals and q output neurons in the output layer); (b) number of neurons in the second HL, varying from 1 to $q\sqrt{m/(q+2)}$; (c) rectified linear activation (ReLU), sigmoid and hyperbolic tangent activation functions; and (d) learning rates of 0.005, 0.001, and 0.0001. The performed tests for the network definition were based on the upper bounds established by^{91,92}.

Proposed approach and validation strategies. Each of the sets of hyperparameters estimated for the MLP networks was used to create a joint and single system for prediction in EG1, which we indicate as part of a divide-and-conquer approach created for genomic prediction (Fig. 1). Considering an individual as part of a dataset subpopulation that has a specific phenotypic distribution, we propose the use of a two-stage prediction process based on the following steps: (1) creating four different neural networks according to different hyperparameter searches and the training data (division step), (2) predicting which subpopulation an unlabeled observation belongs to according to the network induced for this task (prediction 1 and conquer step), and (3) predicting its phenotypic performance based on the network trained specifically for the subpopulation predicted (prediction 2 and final conquer step).

CV1 was the strategy employed for the selection of data for evaluating the models' performance due to its reduced bias when splitting the dataset and the low prediction accuracies described²⁵. We first separated a test dataset using 10% of the genotypes with a stratified holdout strategy implemented in the `scikit-learn Python v.3 module`⁸¹. The stratification was performed only in EG1 and was based on the subpopulation structure present in the dataset. For all the models evaluated in this work (statistical and ML based), the same dataset split was considered in every round of CV.

The remaining 90% of the genotypes were used as the development set for defining the networks' architecture and for evaluating the overall models' performance through a stratified k-fold approach ($k = 4$) with 50 repetitions (subpopulation stratification). The predictive accuracy in every CV split was evaluated by comparing the predicted and real BLUPs by measuring (1) the Pearson correlation coefficient (R) and (2) the mean absolute percentage error (MAPE). For the subpopulation prediction task, we evaluated the classification accuracy (ratio between the number of correctly predicted data and the total number of predictions). For each trait, we compared the predictive accuracy differences using ANOVA and multiple comparisons by Tukey's test with the `agricolae R package`⁹³.

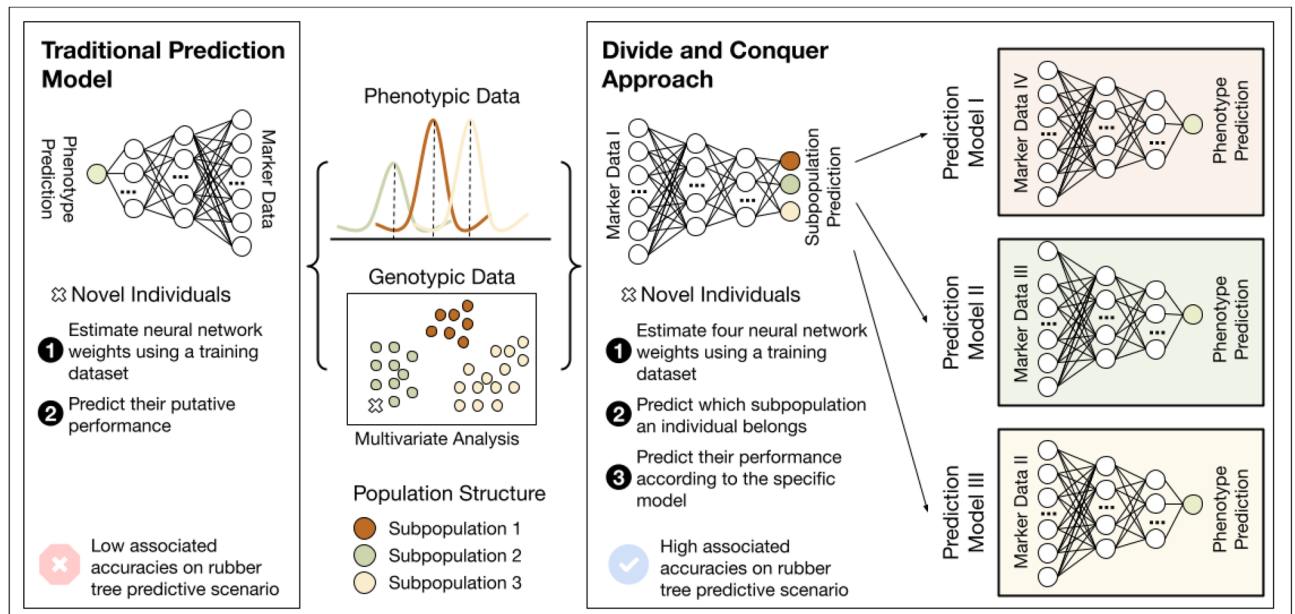


Figure 1. Overview of the approach proposed. Based on a divide-and-conquer strategy with different neural networks combined into a single model (part 1), individuals with unknown phenotypic performance (a) are classified into a subpopulation using a specific neural network (part 2) and (b) have their phenotypic values estimated through an induced network specific to the subpopulation they belong to (part 3).

For EG1, four different MLP architectures were estimated: (a) subpopulation prediction, (b) BLUP prediction for Pop1, (c) BLUP prediction for Pop2, and (d) BLUP prediction for Pop3. After defining the network hyperparameters with the development set, all of these structures were joined into a single predictive system that was used for the final prediction. In addition to evaluating the predictive performance through the CV scenarios created, we also checked the performance of the model for a leave-one-out (LOO) CV configuration.

Results

Phenotypic and genotypic data analyses. The raw phenotypic data were evaluated considering the experimental groups proposed. EG1 (Supplementary Fig. S1) had reduced values compared to those of EG2 (Supplementary Fig. S2) due to the different heights and years of stem measurements. However, for the normalized SC values (Supplementary Figs. S3–S5), such an evident discrepancy was not observed. By modeling the phenotypic measures with the mixed-effect models established and contrasting the raw values with the normalized ones through Q-Q plots, we observed that the residuals obtained with the normalized measurements in EG1 (Supplementary Fig. S6) and EG2 (Supplementary Figs. S7, S8) were more appropriate. Heritabilities (H^2) were estimated as 0.55 for EG1, 0.83 for EG2-S1 and 0.93 for EG2-S2, which is in accordance with the findings of^{2,25}.

Interestingly, BLUPs from EG1 (Supplementary Fig. S9) and EG2-S1 (Supplementary Fig. S10) presented reduced variability when compared to that of BLUPs estimated for EG2-S2 (Supplementary Fig. S10). This observation is corroborated by the hierarchical clustering analyses performed for these experimental groups. EG1 (Supplementary Fig. S10) and EG2-S1 (Supplementary Fig. S12) could be divided into three phenotypic groups according to the best data partitioning scheme established through NbClust clustering indexes⁶², and EG2-S2 could be arranged into 5 such groups (Supplementary Fig. S13). Therefore, it was expected that for the genomic prediction step, EG2-S2 would represent a more difficult task due to its higher data variability.

SNP calling in EG1 was performed according to the TASSEL pipeline. Of the 363,641 tags produced, approximately 84.78% could be aligned against the *H. brasiliensis* reference genome, which generated 107,466 SNPs. These markers were filtered separately for each population using the parameters established, and then these separated datasets were combined through intersection comparisons, yielding a final dataset of 7414 high-quality SNP markers. For EG2 predictions, 332 and 296 SSR markers were used for EG2-S1 and EG2-S2, respectively.

Using these datasets, we performed PCAs for EG1 (Supplementary Fig. S14) and EG2 (Supplementary Fig. S15). In the figures, the colors of the genotypes correspond to their BLUP values, and their shapes correspond to population structure in EG1 and site in EG2. As expected, for the SC trait, there were no clear associations between markers and BLUPs, underlining the challenge of creating genomic prediction models. Additionally, the subpopulation structure in EG1 was evident.

Genomic prediction. From the BLUP and marker datasets, we fit genomic prediction models using the traditional statistical approaches (BRR and SM-GK) and the ML algorithms (AdaBoost, MLP, RF, and SVM) selected. For EG1 (Supplementary Fig. S16), EG2-S1 (Supplementary Fig. S17) and EG2-S2 (Supplementary Fig. S18), no substantial changes were observed when changing the prediction approach. After applying Tukey's

Prediction scenario	Inter2	Inter3
Subpopulation prediction	224	17
GT1 x PB235	345	20
GT1 x RRIM701	454	62
PR255 x PB217	591	119

Table 1. Feature selection strategies performed on the marker dataset considering the intersection among the three methods established (Inter3) and the intersection among at least two out of the three methods established (Inter2).

multiple comparisons test, we found equivalent performance values for SVM, SM-GK and BRR for all the experimental groups. The worst performance was observed for MLP, however, considering the default architectures employed in scikit-learn⁸¹.

Additionally, we also tested the inclusion of FS techniques for increasing model performance in ML algorithms. Using the Inter2 approach, we selected 539 (~7.27%), 69 (~20.78%) and 82 (~27.70%) markers for EG1, EG2-S1 and EG2-S2, respectively. For Inter3, 113 (~1.52%), 8 (~2.41%) and 15 (~5.07%) markers were identified. This SNP subsetting approach was beneficial for EG1 (Supplementary Fig. S19A), EG2-S1 (Supplementary Fig. S20) and EG2-S2 (Supplementary Fig. S21); however, there were less pronounced improvements for data from EG2 sites, which was expected because of the limited SSR marker dataset. We considered that, even with increased predictive accuracies, to achieve better results, a wider set of markers would be required. Then, we considered the best strategy for EG2-S1 to be the combination of the Inter2 FS approach with SVM and that for EG2-S2 to be the combination of Inter3 FS with the AdaBoost ML algorithm.

Even though FS approaches boosted prediction accuracies for EG1, when analyzing model performance by calculating the Pearson correlation between the real and predicted BLUPs for each family separately, we observed that this better performance was caused by the predictions coming from the family with the largest number of individuals, which showed a clear inefficiency of the model for the other families. However, when analyzing predictive power within families (Supplementary Fig. S19B), such an approach was not sufficient for obtaining a reliable prediction with this evident data stratification. In this context, different from EG2, we developed an approach specific to datasets similar to EG1, i.e., a methodology with high capabilities to supply accurate predictions, even considering the subpopulation structure present in a dataset.

Considering a genomic prediction problem based on the creation of a regression model for a dataset containing genotypes that belong to different groups of genetically similar individuals, we modeled such a task by dividing the prediction into different stages (Fig. 1) and creating a divide-and-conquer approach for prediction. The basis of such an approach is that closely related genotypes will share QTLs that might not be the same in another group of genotypes. Therefore, we created a different neural network for each biparental population (divide part), coupled with an intrapopulation system of FS and with a different form of hyperparameter estimation. Following this division part, the separated systems were combined using an additional step (the conquer part). To do so, another neural network was created to infer which subpart of the system should be used for prediction.

Feature selection at the subpopulation level. The selection of subsets of markers was performed according to each EG1 network using the four different tasks: (i) subpopulation prediction, (ii) EG1-Pop1 BLUP prediction, (iii) EG1-Pop2 BLUP prediction, and (iv) EG1-Pop3 BLUP prediction. As expected, each FS strategy returned a different quantity of markers (Table 1). For each subset of markers selected considering Inter2 and Inter3, we evaluated their performance using the ML algorithms selected. Some of the models created for task (i) did not present any mistakes (Supplementary Fig. S22), which was expected due to the subpopulation structure present in the dataset and their evident linear separability. For this task, we considered the most suitable FS strategy to be the Inter2 approach.

For EG1-Pop1 (Supplementary Fig. S23), EG1-Pop2 (Supplementary Fig. S24) and EG1-Pop3 (Supplementary Fig. S25), the best accuracies were observed for the combination Inter2-SVM. However, considering the overall performance with the other algorithms, the best approach for SNP subsetting was Inter3. For this reason, we selected this strategy for the BLUP prediction task. Interestingly, there was no intersection between these three Inter3 datasets in the populations; the only case of overlap was a single SNP marker in Pop2 and Pop3.

From the genomic regions flanking these markers selected for BLUP prediction, we could retrieve several instances of correspondence between rubber tree sequences and CDSs from the *Malpighiales* clade in the Phytosome database. From the 20 markers used in Pop1 for prediction, 62 in Pop2, and 119 in Pop3, we found CDS correspondence for the genomic regions related to 8 (40%), 27 (~43.55%) and 48 (~40.32%) SNPs, respectively. Even though there was no obvious complementarity among these markers due to the absence of intersections, we found GO terms with similar biological processes (Supplementary Tables S1–S3), indicating common molecular processes related to these genomic regions.

Neural network creation. With the marker dataset established through FS for EG1 subtasks, we estimated the best hyperparameter configuration for creating the networks proposed: (i) subpopulation prediction in EG1 (Supplementary Fig. S26), (ii) BLUP prediction in EG1-Pop1 (Supplementary Fig. S27), (iii) BLUP prediction in EG1-Pop2 (Supplementary Fig. S28), and (iv) BLUP prediction in EG1-Pop3 (Supplementary Fig. S29). With

Neural network	N-1HL	N-2HL	LR	AF
EG1 (Subpopulation Prediction)	45	25	0.005	Rectified linear activation
EG1 (BLUP Prediction in GT1 x PB235)	10	3	0.005	Rectified linear activation
EG1 (BLUP Prediction in GT1 x RRIM701)	30	7	0.005	Rectified linear activation
EG1 (BLUP Prediction in PR255 x PB217)	42	4	0.005	Rectified linear activation

Table 2. Hyperparameter definition for each one of the created neural networks in experimental groups 1 (EG1) and 2 (EG2) considering (i) the number of neurons selected for the first hidden layer (N-1HL), (ii) the number of neurons selected for the second hidden layer (N-2HL), (iii) the learning rate (LR), and (iv) the activation function (AF).

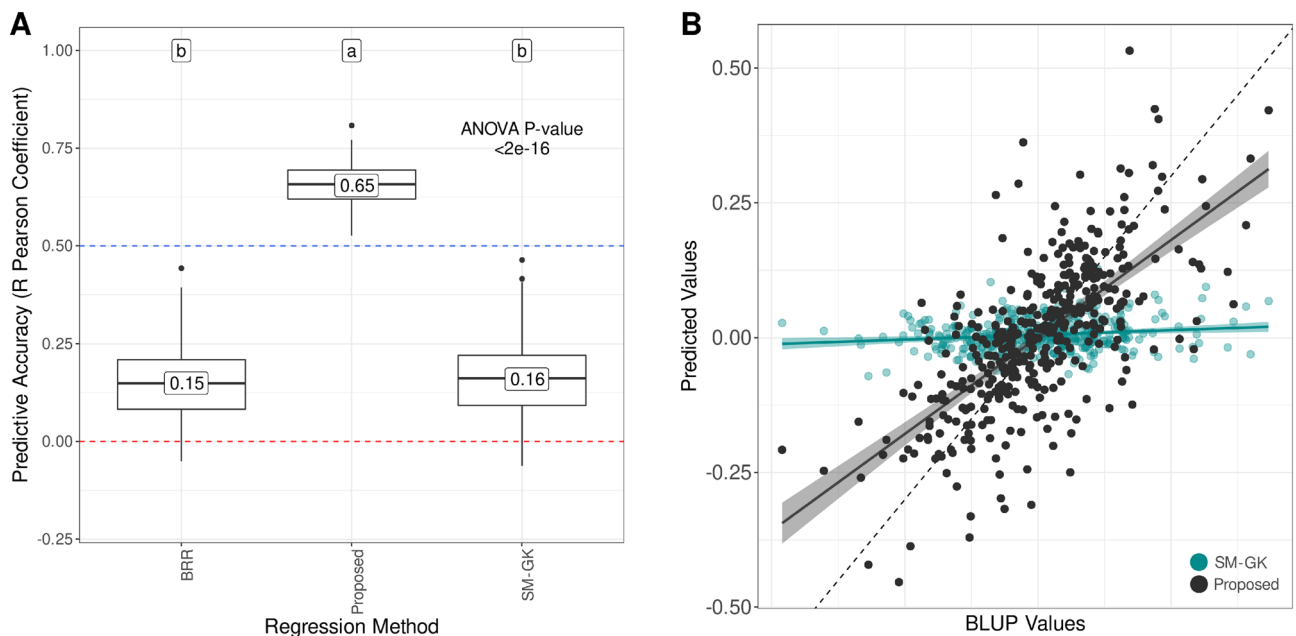


Figure 2. Predictive accuracies for stem circumference BLUP prediction in experimental group 1 (EG1) considering (A) a fourfold cross validation (CV) scheme (50 times repeated) and (B) a leave-one-out CV strategy. The models used for prediction were a single-environment model with a nonlinear Gaussian kernel (SM-GK), Bayesian ridge regression (BRR), and the proposed strategy using the divide-and-conquer approach. The labels indicate the results from Tukey's multiple comparison test.

the exception of network (i), which is a classification task, for each hyperparameter combination, we evaluated the MAPE and R Pearson coefficient values using the development set to select the best configuration for prediction. For network (i), several hyperparameter combinations returned prediction capabilities without mistakes (Supplementary Fig. S26), which led us to select the configuration with the minimum value for the loss function (Table 2).

For networks (ii), (iii) and (iv), we selected the best hyperparameter combination by evaluating the plot profiles. We selected the combinations closest to the right corner of the plots (Supplementary Figs. S27–S29), ideally representing the best MAPE and R Pearson coefficient simultaneously. Interestingly, for the four networks, the best activation function was ReLU, and the learning rate was 0.005, only changing the quantity of neurons in the established HLs. An evaluation of the predictive performance of these networks compared to the traditional genomic prediction approaches with k-fold CV built in the development set revealed significant improvement and effective performance in each population, different from the FS performed using these datasets combined (Supplementary Fig. S19).

The network modeled for EG1-Pop1 showed the largest increases (Supplementary Fig. S30), with a mean improvement of 9 times the initial obtained accuracies. EG1-Pop2 (Supplementary Fig. S31) and EG1-Pop3 (Supplementary Fig. S32) showed increases of 7 and 3 times, respectively. In addition to such significant improvements, the models' performance was also more stable, with the predictive accuracies having a narrow distribution, as observed in the boxplots' conformations.

Divide-and-conquer approach. All of the individual networks were combined to create the proposed approach in EG1. Compared to the traditional approaches, this approach showed a mean improvement of 4 times the initial accuracies (Fig. 2A) in the k-fold evaluations. Moreover, BRR and SM-GK presented equivalent

performance values. Additionally, when analyzing the performance of the development set for predicting the BLUP values of genotypes from the test set, we found Pearson R coefficients of 0.39, 0.42, and 0.81 for BRR, SM-GK, and the proposed approach, respectively, showing the methodology's efficiency even for data not in the development set.

As the final step in model evaluation, we performed a LOO CV split to check whether an increase in the training data improves prediction accuracy. By contrasting the real BLUP values with the predicted values, we found R Pearson coefficients of 0.14, 0.16 and 0.68 for BRR, SM-GK, and the proposed approach, respectively. The regression curve clearly indicates the proposed approach's appropriateness for rubber tree data (Fig. 2B).

Discussion

GS has emerged as a potential tool for application in plant breeding programs^{34–38,94,95}. In rubber tree, previously obtained results^{2,25} have demonstrated the potential of such a technique for reducing breeding cycles. Because of the strong commercial rubber demand, there have been many economic incentives for rubber tree production in more environments beyond its natural range^{1,3}. Considering the difficulty of achieving ideal conditions for cultivating *H. brasiliensis* and the rubber demand, the development of more efficient varieties is needed. However, *Hevea*'s long life cycle considerably reduces breeding efficiency¹⁵. Therefore, the application of GS in rubber tree represents an alternative for achieving the desired rubber production in less time by replacing clone trials and reducing the long period of phenotypic evaluation².

The main objective of rubber tree breeding programs is to increase latex production with rapid growth⁴. Increased SC development can be associated with several rubber tree characteristics, such as growth⁹⁶, latex production²⁵, and drought resistance⁹⁷. Due to the high versatility of SC in evaluating rubber trees^{98–101}, we proposed to develop more effective models for predicting this trait, providing a method to be incorporated into the estimation of tree performance. The lack of high genotype variability in the datasets used represents a real scenario for rubber tree breeding programs²⁵, which face the difficulty of generating a population². In addition to the within-family approach suggested for GS with full-sib families by², the use of interconnected families is a common strategy for perennial species^{22,102,103}.

Using these dataset configurations, we evaluated ML algorithms as a more accurate methodology for predicting SC, a complex trait² obtained a mean accuracy for rubber production in a CV0 scenario of 0.53, which increased to 0.56 when selecting a set of markers based on heterozygosity values. In a CV1 scheme, the mean values ranged between 0.33 and 0.60. In the proposed work, we observed even lower accuracies when using SC instead of rubber production, which is in accordance with the findings of²⁵. In²⁵, the authors achieved mean accuracies ranging between 0.19 and 0.28 in a CV1 scenario, contrasted with a CV2 scheme with values ranging between 0.84 and 0.86. For unknown tested genotypes, the predictive accuracies in rubber tree are low, and the inclusion of GS in *Hevea* breeding programs is therefore still not feasible.

Using the traditional approaches for prediction, we achieved LOO configurations of 0.14 and 0.16 for the BRR and SM-GK approaches, respectively, which is similar to what²⁵ observed. The BRR and SM-GK methodologies were selected to represent a parametric and a semiparametric approach¹⁰⁴. Different from BRR, which estimates marker effects, SM-GK estimates genotype effects through a relationship matrix obtained with a reproducing kernel⁷⁶. Even though²⁵ found similar results when using a linear and a nonlinear kernel for the estimation of the genomic relationship matrix¹⁰⁵, considered GK to have a more flexible structure and a higher associated performance. Therefore, considering these findings together with the fact that no significant differences have been found among statistical models for GS^{31,46,47}, we selected only these two statistical models for predictive evaluation.

Even though some previous attempts did not reveal significant differences in employing ML in GS compared with traditional linear regression methodologies^{32,33,39,52,106}, this is not what we observed in our study, which corroborates the findings of^{23,31,107,108}. This discrepancy may be explained by the different strategies used in the ML algorithms, especially distinct neural network architectures, training methodologies, and CV scenarios. The design of neural network architectures is an important step in using deep learning for prediction because differences in the definition of topologies can lead to decreased accuracies³¹.

Several factors are known to influence prediction accuracy in GS, such as the relationship between the individuals used to train models and those that will be predicted²¹, the size and structure of the populations used²⁴, the trait heritability¹⁰⁹, the marker density¹¹⁰, and the linkage disequilibrium (LD) between the set of markers used and the associated QTLs¹¹¹. This last aspect is especially critical in the datasets employed because of the limited set of markers obtained through GBS and SSR genotyping. Considering the reduced accuracies obtained with the CV1 technique already described in^{2,25}, it was expected that when using a K-fold strategy, the same observations would be found for the traditional regression models.

One of the main challenges in GS is the high dimensionality of the features in the datasets because the number of SNPs is much larger than the number of phenotypic observations¹¹² ('large p , small n ' problem). Although a greater saturation of markers enables an increase in the probability of finding LD, a larger number of markers in the same LD block does not contribute to better prediction performance¹¹⁰. In this context, FS techniques may be an alternative strategy for building a predictive model, considering that not all markers are related to a specific phenotype¹¹³ and that the quantity required for this task directly depends on the complexity and genetic architecture of the traits used¹¹⁰. Therefore, like^{23,82,114–118}, and¹¹⁹, we decided to test the prediction improvements by using an FS technique to enhance network performances.

Subset selection showed improvements for EG2 (Supplementary Figs. S20, S21); however, there were no sizeable improvements because of the genetic complexity of SC¹²⁰ and the low density of SSR markers¹²¹. In EG1, although an overall improvement in prediction accuracy was observed (Supplementary Fig. S19), when evaluating the intrapopulation predictive accuracy, we observed clear inefficiency of the approach, probably caused by the different allele substitution effects between the three subpopulations employed¹¹¹. In such a scenario with

unbalanced interconnected families, novel approaches are needed, and in this work, we have proposed the use of a divide-and-conquer strategy.

In computer science, the divide-and-conquer paradigm is based on the principle that if a problem is not simple enough to be solved directly, it can be divided into subproblems, and their results can be combined¹²². In our prediction task, the BLUPs of the populations could not be properly predicted together; thus, we separated the problem into different networks for prediction, combining the strategy into a single network structure. Such an approach has already been applied to the development of neural network architectures^{123–126}; however, such a formulation has not been explored in genomic prediction. In addition to increasing prediction accuracies, such an approach can reduce the time required for network training and hyperparameter estimation¹²⁴, supply superior model interpretability without loss of performance¹²⁷, and be used in combination with other models¹²⁸, including traditional genomic prediction methods. Considering that in genomic prediction, most of the scenarios include different population structures, such a paradigm can benefit the application and development of GS strategies.

In our dataset, most of the observed variance within SNP markers was caused by population structure, which is clearly shown by the PCA results (Supplementary Fig. S14). As this strong variability can be associated with several genomic regions and influence various traits differently and simultaneously in the populations¹²⁹, we hypothesize that traditional genomic prediction models are not capable of capturing these interpopulation differences related to SC QTLs. This is the main reason why performing FS on these unbalanced datasets together was not a promising strategy in our study. As intrapopulation QTLs are not transferable to other populations, the main effects on phenotypic variation are specific to the within-population genetic structure¹³⁰. In this sense, the prediction task in single populations can be seen as simpler than that in multiple populations¹³¹, which was the basis for developing the divide-and-conquer strategy. Considering the specific effects of causal genetic variants within populations^{132,133}, we tried to incorporate such factors into separate networks with their specific hyperparameter optimization processes.

Interestingly, FS steps performed in the three different populations of EG1 returned different markers, but these markers were putatively associated with genes acting in similar biological processes. GO mRNA splicing was found in the intersection set of markers selected for the three populations. The occurrence of genetic variation related to such a regulatory process may influence the transcription of diverse mRNAs from the same gene in different ways. Such diversity of molecules may be related to differences in phenotypic performance, leading to increased plant capabilities^{134–136}. Additionally, base-excision repair was found in both Pop1 and Pop3, which represents a very important defense pathway for maintaining genomic integrity¹³⁷ and is clearly essential for rubber tree growth and development¹³⁸. Due to the increased quantity of individuals in Pop2 and Pop3, more GO categories were found, including important processes for plant growth, such as response to different types of stress and several metabolic processes¹²⁰.

Different studies have reported the use of deep learning for genomic prediction with various datasets, including for humans^{23,113}, sows¹⁰⁷, and plant species such as soybean¹⁰⁸, wheat^{31–33,39,52}, maize³³, and strawberry and blueberry¹⁰⁶. Even though all of these studies used deep learning, the neural network creation approaches were not the same; some of them included architectures of convolutional neural networks (CNNs)^{106,107,113}, while others included MLPs^{32,33,39,52} or both approaches^{23,31,108}. There is no consensus on the efficiency of neural networks for genomic prediction; however, we decided to use such an architecture for combining multiple training processes into a single predictive structure.

For each of the neural network architectures, we employed an MLP structure. We did not include convolutional operations because of the reduced quantity of markers obtained through FS. Additionally, CNNs were developed for extracting unknown patterns from the dataset, and as we hypothesized that FS operations might work as indicators of QTL regions, such operations would not be necessary. To define the most promising network architecture, we used a grid search, testing different combinations of hyperparameters as already performed in relation to GS strategies^{32,33,39,52}. Although other researchers have used the ‘trial and error’ approach to define the network topology¹³⁹, we preferred to develop a strategy that could be replicated in other predictive scenarios, especially with other traits and crops.

The approximation of functions through neural networks was supported first based on¹⁴⁰ and later on¹⁴¹, which extended the theorem of¹⁴⁰, proving that any continuous function can be represented by a neural network with one HL containing $2n + 1$ nodes (n features) and a more complex activation function than that usually employed by current researchers⁹². It has already been proven that one HL is capable of universal approximation by using a complex activation function^{91,142–145}; however, when using regular functions, such as sigmoid and ReLU functions, there is reduced efficiency of such networks. In this context¹⁴⁶, suggested that two HLs could be a solution for this reduced efficiency. In addition, the usage of an additional HL can substantially reduce the total number of required nodes for a satisfactory predictive capability⁹², and it has already been shown that some problems can be solved only by the use of two HLs^{143,147,148}. In practical situations, a neural network architecture with two HLs generalizes better than that with one and has been considered a superior approach^{143,149}. Therefore, in our study, we decided to include two HLs in our proposed architecture, representing a network with more complex training complexity¹⁵⁰.

Concerning the quantity of hidden neurons in a neural network, many researchers have developed different strategies, aiming at increasing accuracy and prediction while decreasing errors¹³⁹.⁹¹ has already proven that in a network architecture with two HLs, the number of nodes required to achieve a reasonable predictive accuracy with m samples and q output neurons is $\sqrt{(q + 2)m} + 2\sqrt{m/(q + 2)}$ in the first HL and $q\sqrt{m/(q + 2)}$ in the second HL. However, the quantity of suggested nodes tends to lead to overfitting of the training data with any arbitrary small error¹³⁹, and considering the capability of predicting unknown data, these values can be considered the maximum number of nodes in an artificial neural network structure⁹². The lower bound for hidden neurons was already proposed by¹⁵¹, which can be useful for accelerating the learning speed, but there was no evidence on separating this quantity across HLs, and the study was based on an MLP with 3 HLs¹³⁹. Thus, in

our architecture definition, we decided to test a large quantity of neurons, considering the findings of⁹¹, as our upper bound.

The created network coupling the population-specific architectures could increase the initial prediction capabilities by more than four times. Such an improvement represents the first attempt to develop a ML strategy for genomic prediction in rubber tree, with a high potential to be adapted to other species with the same data configuration. Considering a broader scenario with distantly related genotypes belonging to a population with undefined structure, this same approach could be applied. Instead of relying on the predefined stratification, clustering analyses could be performed and used for the divide part. Such a practice is already common in breeding, i.e., taking advantage of population structure for model prediction through multivariate techniques^{152–155}. Taking into account the importance of such group configuration in the differentiation of multiple traits^{156–158}, the strategy developed represents a promising approach for several plant species with a difficult prediction scenario.

The use of GS in rubber tree can optimize breeding programs, and the incorporation of ML techniques can be seen as a new possibility for building more robust models with higher associated prediction capabilities. By using data from rubber tree breeding programs, we were able to generate promising predictive results for a highly complex trait and a novel strategy for prediction, which has significant potential to enhance selection efficiency, and reduce the length of the selection cycle. Although our results confirmed the efficiency of the methodology proposed for rubber tree data, to properly evaluate the full potential of the method in other species and broader scenarios, our approach should be investigated in further studies with more genetically diverse populations in contrasting environments.

Data availability

All the genotypic data from this study are available in the Supplementary Material and under NCBI accessions PRJNA540286 (ID: 5440286) (GT1 × PB235 and GT1 × RRIM701) and PRJNA541308 (ID: 541308) (PR255 × PB217).

Received: 19 April 2022; Accepted: 13 September 2022

Published online: 26 October 2022

References

- Warren-Thomas, E., Dolman, P. M. & Edwards, D. P. Increasing demand for natural rubber necessitates a robust sustainability initiative to mitigate impacts on tropical biodiversity. *Conserv. Lett.* **8**, 230–241 (2015).
- Cros, D. *et al.* Within-family genomic selection in rubber tree (*hevea brasiliensis*) increases genetic gain for rubber production. *Ind. Crops Prod.* **138**, 111464 (2019).
- Ahrends, A. *et al.* Current trends of rubber plantation expansion may threaten biodiversity and livelihoods. *Glob. Environ. Change* **34**, 48–58 (2015).
- Rosa, J. R. B. F. *et al.* Qtl detection for growth and latex production in a full-sib rubber tree population cultivated under suboptimal climate conditions. *BMC Plant Biol.* **18**, 223 (2018).
- Lau, N.-S. *et al.* The rubber tree genome shows expansion of gene family associated with rubber biosynthesis. *Sci. Rep.* **6**, 28594 (2016).
- Tang, C. *et al.* The rubber tree genome reveals new insights into rubber production and species adaptation. *Nat. Plants* **2**, 1–10 (2016).
- Liu, J. *et al.* The chromosome-based rubber tree genome provides new insights into spurge genome evolution and rubber biosynthesis. *Mol. Plant* **13**, 336–350 (2020).
- Roorkiwal, M. *et al.* Genomic-enabled prediction models using multi-environment trials to estimate the effect of genotype × environment interaction on prediction accuracy in chickpea. *Sci. Rep.* **8**, 1–11 (2018).
- González-Camacho, J. M. *et al.* Applications of machine learning methods to genomic selection in breeding wheat for rust resistance. *Plant Genome* **11**, 170104 (2018).
- Hayes, B. J., Lewin, H. A. & Goddard, M. E. The future of livestock breeding: Genomic selection for efficiency, reduced emissions intensity, and adaptation. *Trends Genet.* **29**, 206–214 (2013).
- Lepinasse, D. *et al.* A saturated genetic linkage map of rubber tree (*hevea* spp.) based on rflp, aflp, microsatellite, and isozyme markers. *Theor. Appl. Genet.* **100**, 127–138 (2000).
- Venkatachalam, P., Priya, P., Gireesh, T., Amma, C. S. & Thulaseedharan, A. Molecular cloning and sequencing of a polymorphic band from rubber tree [*hevea brasiliensis* (muell.) arg.]: The nucleotide sequence revealed partial homology with proline-specific permease gene sequence. *Curr. Sci.* **90**, 1510–1515 (2006).
- Nakkanong, K., Nualsri, C. & Sdoodee, S. Analysis of genetic diversity in early introduced clones of rubber tree (*hevea brasiliensis*) using rapid and microsatellite markers. *Songklanakaraj J. Sci. Technol.* **30**, 553–560 (2008).
- de Souza, L. M. *et al.* Development of single nucleotide polymorphism markers in the large and complex rubber tree genome using next-generation sequence data. *Mol. Breed.* **36**, 115 (2016).
- An, Z. *et al.* A high-density genetic map and qtl mapping on growth and latex yield-related traits in *hevea brasiliensis* müll. arg. *Ind. Crops Prod.* **132**, 440–448 (2019).
- Lepinasse, D. *et al.* Identification of qtls involved in the resistance to south american leaf blight (*microcyclus ulei*) in the rubber tree. *Theor. Appl. Genet.* **100**, 975–984 (2000).
- Le Guen, V. *et al.* Bypassing of a polygenic *microcyclus ulei* resistance in rubber tree, analyzed by qtl detection. *New Phytol.* **173**, 335–345 (2007).
- Le Guen, V. *et al.* A rubber tree's durable resistance to *microcyclus ulei* is conferred by a qualitative gene and a major quantitative resistance factor. *Tree Genet. Genomes* **7**, 877–889 (2011).
- Souza, L. M. *et al.* Qtl mapping of growth-related traits in a full-sib family of rubber tree (*hevea brasiliensis*) evaluated in a sub-tropical climate. *PLoS One* **8**, e61238 (2013).
- Tran, D. M. *et al.* Genetic determinism of sensitivity to *corynespora cassiicola* exudates in rubber tree (*hevea brasiliensis*). *PLoS one* **11**, e0162807 (2016).
- Washburn, J. D., Burch, M. B., Franco, V. & José, A. Predictive breeding for maize: Making use of molecular phenotypes, machine learning, and physiological crop models. *Crop Sci.* **60**, 622–38 (2019).
- Muranty, H. *et al.* Accuracy and responses of genomic selection on key traits in apple breeding. *Horticult. Res.* **2**, 1–12 (2015).
- Bellot, P., de Campos, G. & Pérez-Enciso, M. Can deep learning improve genomic prediction of complex human traits?. *Genetics* **210**, 809–819 (2018).

24. Crossa, J. *et al.* Genomic selection in plant breeding: Methods, models, and perspectives. *Trends Plant Sci.* **22**, 961–975 (2017).
25. Souza, L. M. D. *et al.* Genomic selection in rubber tree breeding: A comparison of models and methods for managing $g \times e$ interactions. *Front. Plant Sci.* **10**, 1353 (2019).
26. Hayes, B. *et al.* Prediction of total genetic value using genome-wide dense marker maps. *Genetics* **157**, 1819–1829 (2001).
27. Bernardo, R. & Yu, J. Prospects for genomewide selection for quantitative traits in maize. *Crop Sci.* **47**, 1082–1090 (2007).
28. Heffner, E. L., Lorenz, A. J., Jannink, J.-L. & Sorrells, M. E. Plant breeding with genomic selection: Gain per unit time and cost. *Crop Sci.* **50**, 1681–1690 (2010).
29. Albrecht, T. *et al.* Genome-based prediction of testcross values in maize. *Theor. Appl. Genet.* **123**, 339 (2011).
30. Wang, X., Xu, Y., Hu, Z. & Xu, C. Genomic selection methods for crop improvement: Current status and prospects. *Crop J.* **6**, 330–340 (2018).
31. Ma, W. *et al.* A deep convolutional neural network approach for predicting phenotypes from genotypes. *Planta* **248**, 1307–1318 (2018).
32. Crossa, J. *et al.* Deep kernel and deep learning for genome-based prediction of single traits in multi-environment breeding trials. *Front. Genet.* **10**, 1168 (2019).
33. Montesinos-López, O. A. *et al.* Multi-trait, multi-environment deep learning modeling for genomic-enabled prediction of plant traits. *G3 Genes Genomes Genet.* **8**, 3829–3840 (2018).
34. Zhao, Y. *et al.* Accuracy of genomic selection in European maize elite breeding populations. *Theor. Appl. Genet.* **124**, 769–776 (2012).
35. Spindel, J. *et al.* Genomic selection and association mapping in rice (*Oryza sativa*): Effect of trait genetic architecture, training population composition, marker number and statistical model on accuracy of rice genomic selection in elite, tropical rice breeding lines. *PLoS Genet.* **11**, e1004982 (2015).
36. Crossa, J. *et al.* Genomic prediction of gene bank wheat landraces. *G3 Genes Genomes Genet.* **6**, 1819–1834 (2016).
37. Xavier, A., Muir, W. M. & Rainey, K. M. Assessing predictive properties of genome-wide selection in soybeans. *G3 Genes Genomes Genet.* **6**, 2611–2616 (2016).
38. Wolfe, M. D. *et al.* Prospects for genomic selection in cassava breeding. *Plant Genome* **10**, plantgenome2017-03 (2017).
39. Montesinos-López, O. A. *et al.* Multi-trait, multi-environment genomic prediction of durum wheat with genomic best linear unbiased predictor and deep learning methods. *Front. Plant Sci.* **10**, 1311 (2019).
40. Jarquín, D. *et al.* Increasing genomic-enabled prediction accuracy by modeling genotype \times environment interactions in Kansas wheat. *Plant Genome* **10**, plantgenome2016-12 (2017).
41. VanRaden, P. Genomic measures of relationship and inbreeding. *INTERBULL Bull.* **37**, 33 (2007).
42. VanRaden, P. M. Efficient methods to compute genomic predictions. *J. Dairy Sci.* **91**, 4414–4423 (2008).
43. De Los Campos, G. *et al.* Predicting quantitative traits with regression models for dense molecular markers and pedigree. *Genetics* **182**, 375–385 (2009).
44. Jannink, J.-L., Lorenz, A. J. & Iwata, H. Genomic selection in plant breeding: From theory to practice. *Brief. Funct. Genomics* **9**, 166–177 (2010).
45. Endelman, J. B. Ridge regression and other kernels for genomic selection with R package rrrblup. *Plant Genome* **4**, 250–255 (2011).
46. Roorikwal, M. *et al.* Genome-enabled prediction models for yield related traits in chickpea. *Front. Plant Sci.* **7**, 1666 (2016).
47. Varshney, R. K. Exciting journey of 10 years from genomes to fields and markets: Some success stories of genomics-assisted breeding in chickpea, pigeonpea and groundnut. *Plant Sci.* **242**, 98–107 (2016).
48. Harfouche, A. L. *et al.* Accelerating climate resilient plant breeding by applying next-generation artificial intelligence. *Trends Biotechnol.* **37**, 1217–35 (2019).
49. González-Camacho, J. *et al.* Genome-enabled prediction of genetic values using radial basis function neural networks. *Theor. Appl. Genet.* **125**, 759–771 (2012).
50. Pérez-Rodríguez, P. *et al.* Comparison between linear and non-parametric regression models for genome-enabled prediction in wheat. *G3 Genes Genomes Genet.* **2**, 1595–1605 (2012).
51. Crossa, J. *et al.* Genomic prediction in CIMMYT maize and wheat breeding programs. *Heredity* **112**, 48–60 (2014).
52. Montesinos-López, O. A. *et al.* A benchmarking between deep learning, support vector machine and Bayesian threshold best linear unbiased prediction for predicting ordinal traits in plant breeding. *G3 Genes Genomes Genet.* **9**, 601–618 (2019).
53. Conson, A. R. *et al.* High-resolution genetic map and QTL analysis of growth-related traits of *Hevea brasiliensis* cultivated under suboptimal temperature and humidity conditions. *Front. Plant Sci.* **9**, 1255 (2018).
54. Romain, B. & Thierry, C. Rubberclones (*Hevea* clonal descriptions) (2011).
55. Baudouin, L., Baril, C., Clément-Demange, A., Leroy, T. & Paulin, D. Recurrent selection of tropical tree crops. *Euphytica* **96**, 101–114 (1997).
56. Sivakumaran, S., Haridas, G. & Abraham, P. Problem of tree dryness with high yielding precocious clones and methods to exploit such clones. *Proc. Coll. Hevea* **88**, 253–267 (1988).
57. Rao, G. P. & Kole, P. Evaluation of Brazilian wild *Hevea* germplasm for cold tolerance: Genetic variability in the early mature growth. *J. For. Res.* **27**, 755–765 (2016).
58. Team, R. C. *et al.* R: A language and environment for statistical computing. (2013).
59. Peterson, R. Estimating normalization transformations with bestNormalize. URL <https://github.com/CompPersonR/bestNormalize> (2017).
60. Muñoz, F. & Sanchez, L. breedR: Statistical Methods for Forest Genetic Resources Analysts (2019). R package version 0.12-4.
61. Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. T.-Y. ggTreed: An R package for visualization and annotation of phylogenetic trees with their covariates and other associated data. *Methods Ecol. Evol.* **8**, 28–36 (2017).
62. Charrad, M., Ghazzali, N., Boiteau, V. & Niknafs, A. NbClust: An R package for determining the relevant number of clusters in a data set. *J. Stat. Softw.* **61**, 1–36 (2014).
63. Elshire, R. J. *et al.* A robust, simple genotyping-by-sequencing (GBS) approach for high diversity species. *PLoS ONE* **6**, e19379 (2011).
64. Glaubitz, J. C. *et al.* Tassel-GBS: A high capacity genotyping by sequencing analysis pipeline. *PLoS ONE* **9**, e90346 (2014).
65. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
66. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).
67. Granato, I. S. *et al.* snpReady: A tool to assist breeders in genomic analysis. *Mol. Breed.* **38**, 102 (2018).
68. Cover, T. & Hart, P. Nearest neighbor pattern classification. *IEEE Trans. Inf. Theory* **13**, 21–27 (1967).
69. Le Guen, V., Doaré, F., Weber, C. & Seguin, M. Genetic structure of Amazonian populations of *Hevea brasiliensis* is shaped by hydrographical network and isolation by distance. *Genet. Genomes* **5**, 673–683 (2009).
70. Browning, S. R. & Browning, B. L. Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *Am. J. Hum. Genet.* **81**, 1084–1097 (2007).
71. Wickham, H. *ggplot2: Elegant graphics for data analysis* (Springer, 2016).
72. Aono, A. H. *et al.* A joint learning approach for genomic prediction in polyploid grasses. *Sci. Rep.* **12**, 1–17 (2022).
73. Gianola, D. Priors in whole-genome regression: The Bayesian alphabet returns. *Genetics* **194**, 573–596 (2013).

74. Cuevas, J. *et al.* Genomic prediction of genotype \times environment interaction kernel regression models. *Plant Genome* **9**, plant-genome2016-03 (2016).
75. Pérez, P. & de los Campos, G. Genome-wide regression and prediction with the *bgrr* statistical package. *Genetics* **198**, 483–495 (2014).
76. Granato, I. *et al.* *Bgge*: A new package for genomic-enabled prediction incorporating genotype \times environment interaction models. *G3 Genes Genomes Genet.* **8**, 3039–3047 (2018).
77. Freund, Y. & Schapire, R. E. A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. Syst. Sci.* **55**, 119–139 (1997).
78. Popescu, M.-C., Balas, V. E., Perescu-Popescu, L. & Mastorakis, N. Multilayer perceptron and neural networks. *WSEAS Trans. Circuits Syst.* **8**, 579–588 (2009).
79. Breiman, L. Random forests. *Mach. Learn.* **45**, 5–32 (2001).
80. Shawe-Taylor, J. & Cristianini, N. An introduction to support vector machines and other kernel-based learning methods, vol. 204 (Volume, 2000).
81. Pedregosa, F. *et al.* Scikit-learn: Machine learning in python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011).
82. Aono, A. H. *et al.* Machine learning approaches reveal genomic regions associated with sugarcane brown rust resistance. *Sci. Rep.* **10**, 1–16 (2020).
83. Chen, T. & Guestrin, C. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd International Conference on Knowledge Discovery and Data Mining*, 785–794 (2016).
84. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
85. Goodstein, D. M. *et al.* Phytozome: A comparative platform for green plant genomics. *Nucleic Acids Res.* **40**, D1178–D1186 (2012).
86. Botstein, D. *et al.* Gene ontology: Tool for the unification of biology. *Nat. Genet.* **25**, 25–9 (2000).
87. Chollet, F. *et al.* Keras. <https://keras.io> (2015).
88. Da Silva, I. N., Spatti, D. H., Flauzino, R. A., Liboni, L. H. B. & dos Reis Alves, S. F. *Artificial Neural networks* Vol. 39 (Springer, 2017).
89. Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. Preprint [arXiv:1412.6980](https://arxiv.org/abs/1412.6980) (2014).
90. Bengio, Y. Rmsprop and equilibrated adaptive learning rates for nonconvex optimization. [cor arXiv:abs/1502.04390](https://arxiv.org/abs/1502.04390) (2015).
91. Huang, G.-B. Learning capability and storage capacity of two-hidden-layer feedforward networks. *IEEE Trans. Neural Netw.* **14**, 274–281 (2003).
92. Stathakis, D. How many hidden layers and nodes?. *Int. J. Remote Sens.* **30**, 2133–2147 (2009).
93. de Mendiburu, F. & de Mendiburu, M. F. Package ‘agricolae’. R Package, Version 1–2 (2019).
94. O’Connor, K., Hayes, B. & Topp, B. Prospects for increasing yield in macadamia using component traits and genomics. *Genet. Genomes* **14**, 7 (2018).
95. Cros, D. *et al.* Genomic selection prediction accuracy in a perennial crop: Case study of oil palm (*elaeis guineensis* jacq.). *Theor. Appl. Genet.* **128**, 397–410 (2015).
96. Chandrashekar, T. *et al.* An analysis of growth and drought tolerance in rubber during the immature phase in a dry subhumid climate. *Exp. Agric.* **34**, 287–300 (1998).
97. Zhang, C., Stratopoulos, L. M. F., Pretzsch, H. & Rötzer, T. How do *tilia cordata* greenspire trees cope with drought stress regarding their biomass allocation and ecosystem services?. *Forests* **10**, 676 (2019).
98. Dijkman, M. J. *et al.* Hevea, thirty years of research in the far east. Hevea, Thirty years of research in the Far East. (1951).
99. Gonçalves, P. d. S., Rossetti, A. G., Valois, A. C. C. & VIEGAS, I. Estimativas de correlações genéticas e fenotípicas de alguns caracteres quantitativos em clones jovens de seringueira (*hevea* spp). Embrapa Amazônia Ocidental-Artigo em periódico indexado (ALICE) (1984).
100. Chanroj, V., Rattanawong, R., Phumichai, T., Tangphatsornruang, S. & Ukoskit, K. Genome-wide association mapping of latex yield and girth in amazonian accessions of *hevea brasiliensis* grown in a suboptimal climate zone. *Genomics* **109**, 475–484 (2017).
101. Khan, M. A. *et al.* Analysis of *qt*-allele system conferring drought tolerance at seedling stage in a nested association mapping population of soybean [*g* *lycine max* (L.) *merr.*] using a novel *gwas* procedure. *Planta* **248**, 947–962 (2018).
102. Kumar, S. *et al.* Genome-enabled estimates of additive and nonadditive genetic variances and prediction of apple phenotypes across environments. *G3 Genes Genomes Genet.* **5**, 2711–2718 (2015).
103. Grattapaglia, D. Status and perspectives of genomic selection in forest tree breeding. In *Genomic Selection for Crop Improvement*, 199–249 (Springer, 2017).
104. Heslot, N., Yang, H.-P., Sorrells, M. E. & Jannink, J.-L. Genomic selection in plant breeding: A comparison of models. *Crop Sci.* **52**, 146–160 (2012).
105. Gianola, D., Weigel, K. A., Krämer, N., Stella, A. & Schön, C.-C. Enhancing genome-enabled prediction by bagging genomic blup. *PLoS One* **9**, e91693 (2014).
106. Zingaretti, L. M. *et al.* Exploring deep learning for complex trait genomic prediction in polyploid outcrossing species. *Front. Plant Sci.* **11**, 25 (2020).
107. Waldmann, P., Pfeiffer, C. & Mészáros, G. Sparse convolutional neural networks for genome-wide prediction. *Front. Genet.* **11**, 25 (2020).
108. Liu, Y. *et al.* Phenotype prediction and genome-wide association study using deep convolutional neural network of soybean. *Front. Genet.* **10**, 1091 (2019).
109. Zhang, A. *et al.* Effect of trait heritability, training population size and marker density on genomic prediction accuracy estimation in 22 bi-parental tropical maize populations. *Front. Plant Sci.* **8**, 1916 (2017).
110. Liu, X. *et al.* Factors affecting genomic selection revealed by empirical evidence in maize. *Crop J.* **6**, 341–352 (2018).
111. Raymond, B., Bouwman, A. C., Schrooten, C., Houwing-Duistermaat, J. & Veerkamp, R. F. Utility of whole-genome sequence data for across-breed genomic prediction. *Genet. Sel. Evol.* **50**, 1–12 (2018).
112. Long, N., Gianola, D., Rosa, G. J., Weigel, K. A. & Avendaño, S. Machine learning classification procedure for selecting snps in genomic selection: Application to early mortality in broilers. *J. Anim. Breed. Genet.* **124**, 377–389 (2007).
113. Yin, B. *et al.* Using the structure of genome data in the design of deep neural networks for predicting amyotrophic lateral sclerosis from genotype. *Bioinformatics* **35**, i538–i547 (2019).
114. Bermingham, M. L. *et al.* Application of high-dimensional feature selection: Evaluation for genomic prediction in man. *Sci. Rep.* **5**, 1–12 (2015).
115. Li, B. *et al.* Genomic prediction of breeding values using a subset of snps identified by three machine learning methods. *Front. Genet.* **9**, 237 (2018).
116. Inácio, Í. S. C. G. F. & Alves, M. F. C. Increasing accuracy and reducing costs of genomic prediction by marker selection. *Euphytica* **215**, 18 (2019).
117. Ramzan, F., Gültas, M., Bertram, H., Cavero, D. & Schmitt, A. O. Combining random forests and a signal detection method leads to the robust detection of genotype-phenotype associations. *Genes* **11**, 892 (2020).
118. Luo, Z., Yu, Y., Xiang, J. & Li, F. Genomic selection using a subset of snps identified by genome-wide association analysis for disease resistance traits in aquaculture species. *Aquaculture* **539**, 736620 (2021).

119. Pimenta, R. J. G. *et al.* Genome-wide approaches for the identification of markers and genes associated with sugarcane yellow leaf virus resistance. *Sci. Rep.* **11**, 1–18 (2021).
120. Francisco, F. R. *et al.* Unravelling rubber tree growth by integrating gwas and biological network-based approaches. *Front. Plant Sci.* **2719**, 12 (2021).
121. Nadeem, M. A. *et al.* Dna molecular markers in plant breeding: Current status and recent advancements in genomic selection and genome editing. *Biotechnol. Biotechnol. Equip.* **32**, 261–285 (2018).
122. Smith, D. R. The design of divide and conquer algorithms. *Sci. Comput. Prog.* **5**, 37–58 (1985).
123. Frosyniotis, D., Stafylopatis, A. & Likas, A. A divide-and-conquer method for multi-net classifiers. *Pattern Anal. Appl.* **6**, 32–40 (2003).
124. Mohamad, M. Divide and conquer approach in reducing ann training time for small and large data. *J. Appl. Sci.* **13**, 133–139 (2013).
125. Feng, J., Wang, L., Yu, H., Jiao, L. & Zhang, X. Divide-and-conquer dual-architecture convolutional neural network for classification of hyperspectral images. *Remote Sens.* **11**, 484 (2019).
126. Sakhakarmi, S. & Park, J. W. Multi-level-phase deep learning using divide-and-conquer for scaffolding safety. *Int. J. Environ. Res. Public Health* **17**, 2391 (2020).
127. Fu, W., Breininger, K., Schaffert, R., Ravikumar, N. & Maier, A. A divide-and-conquer approach towards understanding deep networks. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 183–191 (Springer, 2019).
128. Intanagonwiwat, C. The divide-and-conquer neural network: its architecture and training. In 1998 IEEE International Joint Conference on Neural Networks Proceedings. IEEE World Congress on Computational Intelligence (Cat. No. 98CH36227), vol. 1, 462–467 (IEEE, 1998).
129. Linhart, Y. B. & Grant, M. C. Evolutionary significance of local genetic differentiation in plants. *Annu. Rev. Ecol. Syst.* **27**, 237–277 (1996).
130. Würschum, T. Mapping qtl for agronomic traits in breeding populations. *Theor. Appl. Genet.* **125**, 201–210 (2012).
131. Ogut, F., Bian, Y., Bradbury, P. J. & Holland, J. B. Joint-multiple family linkage analysis predicts within-family variation better than single-family analysis of the maize nested association mapping population. *Heredity* **114**, 552–563 (2015).
132. Hirschhorn, J. N. *et al.* Genomewide linkage analysis of stature in multiple populations reveals several regions with evidence of linkage to adult height. *Am. J. Hum. Genet.* **69**, 106–116 (2001).
133. Pressoir, G. & Berthaud, J. Patterns of population structure in maize landraces from the central valleys of Oaxaca in Mexico. *Heredity* **92**, 88–94 (2004).
134. Mastrangelo, A. M., Marone, D., Laidò, G., De Leonardis, A. M. & De Vita, P. Alternative splicing: Enhancing ability to cope with stress via transcriptome plasticity. *Plant Sci.* **185**, 40–49 (2012).
135. Wei, H. *et al.* Alternative splicing complexity contributes to genetic improvement of drought resistance in the rice maintainer huhan2b. *Sci. Rep.* **7**, 1–13 (2017).
136. Szakonyi, D. & Duque, P. Alternative splicing as a regulator of early plant development. *Front. Plant Sci.* **9**, 1174 (2018).
137. Roldán-Arjona, T., Ariza, R. R. & Córdoba-Cañero, D. Dna base excision repair in plants: An unfolding story with familiar and novel characters. *Front. Plant Sci.* **10**, 1055 (2019).
138. Murphy, T. M. What is base excision repair good for?: Knockout mutants for fpg and ogg glycosylase genes in Arabidopsis. *Physiol. Plant.* **123**, 227–232 (2005).
139. Sheela, K. G. & Deepa, S. N. Review on methods to fix number of hidden neurons in neural networks. *Math. Probl. Eng.* **2013** (2013).
140. Kolmogorov, A. N. On the representation of continuous functions of many variables by superposition of continuous functions of one variable and addition. In *Doklady Akademii Nauk*, vol. 114, 953–956 (Russian Academy of Sciences, 1957).
141. Hecht-Nielsen, R. Kolmogorov's mapping neural network existence theorem. In *Proceedings of the International Conference on Neural Networks*, vol. 3, 11–14 (IEEE Press New York, 1987).
142. Wang, S.-C. Artificial neural network. In *Interdisciplinary Computing in Java Programming*, 81–100 (Springer, 2003).
143. Thomas, A. J., Petridis, M., Walters, S. D., Gheytaasi, S. M. & Morgan, R. E. Two hidden layers are usually better than one. In *International Conference on Engineering Applications of Neural Networks*, 279–290 (Springer, 2017).
144. Hornik, K. *et al.* Multilayer feedforward networks are universal approximators. *Neural Netw.* **2**, 359–366 (1989).
145. Hornik, K. Some new results on neural network approximation. *Neural Netw.* **6**, 1069–1072 (1993).
146. Kurková, V. Kolmogorov's theorem and multilayer neural networks. *Neural Netw.* **5**, 501–506 (1992).
147. Chester, D. L. Why two hidden layers are better than one. In *Proceedings of IJCNN*, Washington, DC, vol. 1, 265–268 (1990).
148. Sontag, E. D. Feedback stabilization using two-hidden-layer nets. In *1991 American Control Conference*, 815–820 (IEEE, 1991).
149. Islam, M. M. & Murase, K. A new algorithm to design compact two-hidden-layer artificial neural networks. *Neural Netw.* **14**, 1265–1278 (2001).
150. Kurková, V. & Sanguineti, M. Can two hidden layers make a difference? In *International Conference on Adaptive and Natural Computing Algorithms*, 30–39 (Springer, 2013).
151. Jiang, N., Zhang, Z., Ma, X. & Wang, J. The lower bound on the number of hidden neurons in multi-valued multi-threshold neural networks. In *2008 Second International Symposium on Intelligent Information Technology Application*, vol. 1, 103–107 (IEEE, 2008).
152. Guo, Z. *et al.* The impact of population structure on genomic prediction in stratified populations. *Theor. Appl. Genet.* **127**, 749–762 (2014).
153. Wang, Q. *et al.* Effects of marker density and population structure on the genomic prediction accuracy for growth trait in Pacific white shrimp *Litopenaeus vannamei*. *BMC Genet.* **18**, 1–9 (2017).
154. Berro, I., Lado, B., Nalin, R. S., Quincke, M. & Gutiérrez, L. Training population optimization for genomic selection. *Plant Genome* **12**, 190028 (2019).
155. Stewart-Brown, B. B., Song, Q., Vaughn, J. N. & Li, Z. Genomic selection for yield and seed composition traits within an applied soybean breeding program. *G3 Genes Genomes Genet.* **9**, 2253–2265 (2019).
156. Goodnight, C. J. Population differentiation and the correlation among traits at the population level. *Am. Nat.* **133**, 888–900 (1989).
157. Merilä, J. & Crnokrak, P. Comparison of genetic differentiation at marker loci and quantitative traits. *J. Evol. Biol.* **14**, 892–903 (2001).
158. Bolnick, D. I. *et al.* Why intraspecific trait variation matters in community ecology. *Trends Ecol. Evol.* **26**, 183–192 (2011).

Acknowledgements

The authors gratefully acknowledge the Fundação de Amparo à Pesquisa do Estado de São Paulo (FAPESP) for Ph.D. fellowships to FF (2018/18985-7) and AA (2019/03232-6) and for a research internship abroad (BEPE) scholarship to AA (2019/26858-8), the Coordenação de Aperfeiçoamento do Pessoal de Nível Superior (CAPES)

for financial support (Computational Biology Program and CAPES-Agropolis Program), and the Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) for research fellowships to AS and PG.

Author contributions

A.A. and F.F. performed all the analyses and wrote the manuscript; P.G., E.J. and V.G. conducted the field experiments; L.S., R.F., M.Q., G.G. and A.S. conceived the project. All authors reviewed, read and approved the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-022-20416-z>.

Correspondence and requests for materials should be addressed to A.P.S.

Reprints and permissions information is available at www.nature.com/reprints.

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



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022