

Relatedness and Genotype \times Environment Interaction Affect Prediction Accuracies in Genomic Selection: A Study in Cassava

Delphine Ly, Martha Hamblin,* Ismail Rabbi, Gedil Melaku, Moshood Bakare, Hugh G. Gauch Jr., Richardson Okechukwu, Alfred G.O. Dixon, Peter Kulakow, and Jean-Luc Jannink

ABSTRACT

Before implementation of genomic selection, evaluation of the potential accuracy of prediction can be obtained by cross-validation. In this procedure, a population with both phenotypes and genotypes is split into training and validation sets. The prediction model is fitted using the training set, and its accuracy is calculated on the validation set. The degree of genetic relatedness between the training and validation sets may influence the expected accuracy as may the genotype \times environment (G \times E) interaction in those sets. We developed a method to assess these effects and tested it in cassava (*Manihot esculenta* Crantz). We used historical phenotypic data available from the International Institute of Tropical Agriculture Genetic Gain trial and performed genotyping by sequencing for these clones. We tested cross-validation sampling schemes preventing the training and validation sets from sharing (i) genetically close clones or (ii) similar evaluation locations. For 19 traits, plot-basis heritabilities ranged from 0.04 to 0.66. The correlation between predicted and observed phenotypes ranged from 0.15 to 0.47. Across traits, predicting for less related clones decreased accuracy from 0 to 0.07, a small but consistent effect. For 17 traits, predicting for different locations decreased accuracy between 0.01 and 0.18. Genomic selection has potential to accelerate gains in cassava and the existing training population should give a reasonable estimate of future prediction accuracies.

D. Ly, Montpellier Supagro, 2, place Pierre Viala, 34060 Montpellier Cedex 02, France; M.T. Hamblin and J-L Jannink, Dep. of Plant Breeding and Genetics, Cornell Univ., Ithaca, NY 14853; H.G. Gauch Jr., Dep. of Crop and Soil Sciences, Cornell Univ., Ithaca, NY 14853; I.Y. Rabbi, M. Gedil, M. Bakare, R. Okechukwu, and P. Kulakow, International Institute for Tropical Agriculture, PMB 5320, Oyo Road, Ibadan, Nigeria; A.G.O. Dixon, Sierra Leone Agricultural Research Institute, Tower Hill, P.M. B. 1313, Freetown, Sierra Leone; J-L Jannink, USDA-ARS, R.W. Holley Center for Agriculture and Health, Ithaca, NY 14853. Received 21 Nov. 2012. *Corresponding author (mth3@cornell.edu).

Abbreviations: AMMI, additive main effect and multiplicative interaction; a_{top10} , mean relatedness of the top10 individuals in the validation set to those in the training set; AYT, advanced yield trial; BLUE, best linear unbiased estimator; BLUP, best linear unbiased predictor; CMD, cassava mosaic disease; CMDI, cassava mosaic disease incidence; CV-CR, cross-validation close relatives; CV-GE, cross-validation genotype \times environment; CV-Random, random cross-validation; CV-Random_Half, cross-validation scheme in which a randomly chosen half of the observations are used; CV-noCR, cross-validation no close relatives; DM, root dry matter content; GS, genomic selection; G \times E, genotype \times environment; GBS, genotyping by sequencing; MAS, marker-assisted selection; MCBBI, mean cassava bacterial blight incidence; PYT, preliminary yield trial; RCBD, randomized complete-block design; RKHS, reproducing kernel Hilbert spaces; SNP, single nucleotide polymorphism; top10, 10 most closely related individuals; UYT, uniform yield trial.

THE REVOLUTION in sequencing technologies has enabled fast and relatively inexpensive genome information (Metzker, 2010). The increase in DNA-marker information available is considerable, leading to the development of a new approach

Published in Crop Sci. 53:1312–1325 (2013).

doi: 10.2135/cropsci2012.11.0653

Freely available online through the author-supported open-access option.

© Crop Science Society of America | 5585 Guilford Rd., Madison, WI 53711 USA

All rights reserved. No part of this periodical may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or any information storage and retrieval system, without permission in writing from the publisher. Permission for printing and for reprinting the material contained herein has been obtained by the publisher.

to marker-assisted selection (MAS) called genomic selection (GS) (Meuwissen et al., 2001; Goddard and Hayes, 2007; Heffner et al., 2009; Hayes et al., 2009a; Lorenz et al., 2011). The GS approach has been developed to use all markers across the genome, instead of only those with larger effects as in traditional MAS, to predict the performance of individuals (Meuwissen et al., 2001; Jannink et al., 2010). With a sufficient accuracy, selection can be done based on the predictions only, for any trait. Therefore, GS offers the possibility to accelerate breeding cycles. Because prediction requires from the selection candidates only genotypic data, the selection of the seedling happens at an early stage (Heffner et al., 2009). In addition, selecting on the basis of predicted breeding values of individuals rather than their phenotypic records may also make the choice of the parents more accurate.

Our study focused on cassava, which, unlike other crops for which GS has been evaluated, is a strongly outcrossing species, characterized by monoecism and protogyny. This outcrossing characteristic is shared with cattle (*Bos taurus*), a species for which GS has been shown to work effectively (VanRaden et al., 2009; Hayes et al., 2009a). Indeed, a cassava GS study using a relatively small training population and relatively low-density markers has reported reasonable prediction accuracies (Oliveira et al., 2012). Nevertheless, GS accuracies in cassava need further empirical testing.

The realized gains of a genomic selection program will depend on the quality of the predictions, which can be assessed by estimating the prediction accuracies. Genomic selection studies on empirical data generally use cross-validation to estimate prediction accuracies (Kohavi, 1995; Goddard and Hayes, 2007; Lorenz et al., 2011). In plants, the use of cross-validation studies on inbred cultivars has been useful (Melchinger et al., 2004; Schön et al., 2004; Crossa et al., 2010; Riedelsheimer et al., 2012; Massman et al., 2012; Windhausen et al., 2012). Cross-validation is meant to estimate the accuracy with which predictions can be made for selection candidates based on models developed in the training population, by treating a portion of the training population as selection candidates. There are some important differences between cross-validation and the prediction of breeding values in selection candidates, particularly with respect to two factors: the relatedness between individuals and genotype \times environment (G \times E) interaction. In particular, a random cross-validation (CV-Random) might split data into the training and the validation sets so that the information for close relatives or for locations is unrealistically similar in the two sets.

It has been shown that the additive genetic relationship of the training data influences the breeding-value accuracies of the selection candidates (Habier et al., 2007, 2010; Clark et al., 2012). Animals that shared close relationship to the training dataset had highest prediction accuracies (Habier et al., 2007, 2010; Clark et al., 2012;

Pszczola et al., 2012; Pérez-Cabal et al., 2012; Cleveland et al., 2012). In a study on U.S. dairy cattle, Pérez-Cabal et al. (2012) emphasized that the type of relatedness between the training and validation sets also influenced the prediction accuracies. Cleveland et al. (2012) have pointed out that using validation approaches that take into account relatedness between populations can correct for potential overestimation of genomic breeding-value accuracies. In the case of cassava, there are a number of factors that may affect the relatedness of clones in the training population.

Because cassava in Africa originated from South America (Jones, 1959), African cassava germplasm experienced a genetic bottleneck (Kawuki et al., 2011). Furthermore, selection, by frequently using specific elite parents in breeding programs, could make cassava clones in Africa relatively genetically similar. In addition, because of the way landrace germplasm has been collected, there are situations when virtually identical clones may be given different names. This relatedness between clones may have a strong impact on the assessment of the efficiency of GS in cassava. With a random k -fold cross-validation, the genetic relatedness of individuals between the training and validation sets might be higher than that between this whole population, that is, the training population, and the individuals of the next generation of the breeding cycle, that is, the selection candidates.

A second factor that may influence prediction accuracies is G \times E interaction. Genotype and environment effects are not independent: a phenotypic response to a change in environments depends on genotype and vice versa (Comstock and Moll, 1963). With a random k -fold cross-validation, the data in training and validation sets are likely to have been evaluated in the same locations. In that case, G \times E interaction would generate a common error component between the predictions and the clone estimates based on the observations (Lorenz et al., 2011, 2012; Burgueño et al., 2012). Consequently, G \times E interaction may be a confounding factor that upwardly biases the prediction accuracy.

The objectives of our study were to assess the impact of random k -fold cross-validation on the overestimation of prediction accuracies attributable (i) to the relatedness of the individuals between the training and the validation sets and (ii) to the G \times E interaction. The data for the study came from the International Institute for Tropical Agriculture (IITA) Genetic Gain population, a large collection of historically important clones, maintained at Ibadan, Nigeria.

MATERIALS AND METHODS

Phenotypic Trials

Historical phenotypic evaluation data from several types of trials have been used in the training population. All of these trials were conducted by the cassava breeding program at IITA, Ibadan, Nigeria.

Table 1. Description of the traits of interest.

Type of trait	Abbreviation	Name of trait	Description	
Agronomic	SPROUT	Sprouting	Proportion of stakes germinated scored 1 mo after planting	
	VIGOR	Initial vigor	Degree of initial vigor of the establishment scored 1 mo after planting. It is scored from 3, which corresponds to a low vigor, to 7, which is high.	
	HI	Harvest index	Ratio of fresh root weight divided by total biomass	
	DM	Root dry matter content	Percentage dry matter storage root. It measures root dry weight as the percentage of 100 g of the root tubers	
	RTWT	Fresh weight of storage root	Total fresh weight of storage roots harvested per plot measured in kilograms	
	FYLD	Fresh root yield	Fresh weight of harvested roots expressed in tonnes per hectares per plant at harvest	
	DYLD	Dry yield	Dry weight of harvested roots derived by multiplying fresh storage root yield by dry matter content expressed in tonnes per hectares	
	SHTWT	Fresh shoot weight	Total fresh weight of harvested foliage and stems in kilograms per plot	
	TYLD	The top yield	The total fresh weight of harvested foliage and stems expressed in tonnes per hectare	
	RTNO	Root number	Number of storage roots per plot at harvest	
	NOHAV	Plant stands harvested	Counts the number of plant stand at harvest	
	Morphological measured by visual rating	NKLG	Root neck length	Usually scored on a scale of 0 (absent or sessile), 3 (short), 5 (medium), and 7 (long)
		ROTNO	Rotted storage roots	Counts the number of the rotted root per plot at the time of harvest
	Biotic stresses	CMDS	Cassava mosaic disease severity	Cassava mosaic disease (CMD) severity rated on a scale from 1 (no symptoms) to 5 (extremely severe)
CMDI		Cassava mosaic disease incidence	Cassava mosaic disease incidence is the proportion of plants showing CMD symptoms	
CBBS		Cassava bacterial blight severity	Cassava bacterial blight (CBB) severity rated on a scale from 1 (no symptoms) to 5 (extremely severe)	
CBBI		Cassava bacterial blight incidence	Cassava bacterial blight incidence is the proportion of plants showing CBB symptoms	
CADS		Cassava anthracnose disease severity	Cassava anthracnose disease (CAD) severity rated on a scale from 1 (no symptoms) to 5 (extremely severe)	
CADI		Cassava anthracnose disease incidence	Cassava anthracnose disease incidence is the proportion of plants showing CAD symptoms	
CGM		Cassava green mite	Cassava green mite measures the severity rated on a scale from 1 (no symptoms) to 5 (extremely severe)	

The Genetic Gain is a collection of historically important clones that were selected across four decades, from the 1970s to 2007 (Maziya-Dixon et al., 2007; Okechukwu and Dixon, 2008). A small fraction of the clones are landraces and clones from East Africa with uncertain cloning dates. Most additions to the Genetic Gain population came from clones advanced to multi-environment uniform yield trials. The design of the Genetic Gain trial, the major trial type, usually consists of five-plant plots, with no borders, in a single row. Most but not all of the Genetic Gain nurseries are replicated twice. Sometimes the plants are grown at a different density, for example, 0.5 m apart within rows, if the land area available is limited. The plots are always planted in an incomplete block design with two checks per block.

The second most common trial type, called the uniform yield trial (UYT), contains clones that are at an advanced stage in the breeding process. Compared to the Genetic Gain trial, only 15 to 30 genotypes are evaluated in a single trial because UYTs are formed after several stages of selection. Often genotypes are grouped by particular types of traits, such as multiple pest resistance, high dry matter content, or poundability. This type of trial has larger, bordered plots: generally six rows of six plants spaced 1 m apart, in four replications, planted in a randomized complete-block design (RCBD) with two checks. Because of the borders, only 16 plants per plot are harvested. Some variation in plot size for these trials occurred as they were conducted across a period of 12 yr. Uniform yield trials are almost always multilocation and multiyear trials, most commonly conducted across 2 yr and five locations.

Two other types of trial, preliminary yield trial (PYT) and advanced yield trial (AYT), represented less than 10% of the observations. Those trials were conducted earlier than the UYTs in the breeding process, so their design was intermediate between the Genetic Gain and the UYT designs. In the PYT, there are usually 10-plant plots grown in a single 10-m-long row, in one location with two replications. In recent years, the AYT plot design has been the same as for the PYT, but there are usually four replications and one location. The design of both PYT and AYT was an RCBD with two checks.

The data were collected from 2000 to 2011 in 13 locations in Nigeria: Abuja (8.99° N, 7.51° E), Ibadan (7.40° N, 3.90° E), Ilorin (8.50° N, 4.53° E), Ikenne (6.7° N, 3.5° E), Jos (9.94° N, 8.85° E), Kano (12° N, 8.5° E), Mallam Madori (12.3° N, 9.7° E), Mokwa (9.3° N, 5.0° E), Ubiaja (6.66° N, 6.38° E), Onne (4.74° N, 7.15° E), Shonga (9.14° N, 5.1° E), Warri (5.52° N, 5.75° E), and Zaria (10.98° N, 7.76° E).

Eleven agronomic traits and two morphological traits were measured (Table 1). The agronomic trait “plant stands harvested” was used only as a covariate in the statistical models for other traits with which it was correlated (see below). Seven of the traits are related to four biotic stresses: cassava mosaic disease (CMD), caused by a virus from the *Begomovirus* genus that belongs to the Geminiviridae, vectored by the whitefly; cassava bacterial blight, caused by *Xanthomonas axonopodis* pv. *Manihotis*; cassava anthracnose disease, caused by *Colletotrichum gloeosporioides*; and cassava green mite, *Mononychellus tanajoa*.

Genotyping

Using DNeasy Plant Mini Kits (Qiagen), DNA was extracted from 645 clones from the 2011 Genetic Gain trial at IITA and was quantified using PicoGreen. The genotype data were generated using genotyping by sequencing (GBS) (Elshire et al., 2011). Six 95-plex and one 75-plex *Pst*I libraries were constructed and sequenced on Illumina HiSeq, one lane per library.

Single nucleotide polymorphisms (SNPs) were extracted from the raw data by using the TASSEL pipeline version 3.0 (Glaubitz et al., 2012) installed in the Computational Biology Application Suite for High Performance Computing at Cornell University, with alignment to the *Manihot esculenta* reference genome (<http://www.phytozome.net>; accessed 30 Sept. 2011). Single nucleotide polymorphisms were filtered by the following criteria: no more than 80% missing data by clone, no more than 50% missing data by SNP, amount of missing data consistent with read depth, and genotype frequencies consistent with allele frequencies. The final data set consisted of 2069 SNPs scored in 626 clones, with a mean heterozygosity of 0.28 and mean missingness of 17.6%. Because the cassava reference genome is not assembled into chromosomes, it is not possible to show the distribution of SNPs across the genome. However, an overlapping set of GBS SNPs from a *Pst*I library have been genetically mapped and are well distributed across the 18 linkage groups of the cassava genome (I.Y. Rabbi, M.T. Hamblin, M. Gedil, P. Kulakow, A.S. Ikpan, D. Ly, and J.L. Jannink, personal communication, 2012). The missing genotypic data were imputed using a classification method called random forest (Breiman, 2001; Poland et al., 2012).

Statistical Models for Phenotypic Data

Several statistical models were used. The first group of models was used to calculate broad-sense heritabilities on a single plot basis and generate best linear unbiased predictors (BLUPs) for data curation. The second group of models was used to generate best linear unbiased estimators (BLUEs) as an intermediate step to make predictions. The difference between these groups of models was whether they considered the clone effect as random or fixed. For the first group of models, to calculate heritabilities and generate BLUPs, we used mixed models that considered the available clones as a random sample. Mixed models were performed using the *lme4* package in R (R Development Core Team, 2008).

Data were available for 12 yr and 13 locations, but not all locations were evaluated every year. Each combination of a particular year and a particular location was considered as an environment. Within each environment, several types of trials were conducted. Within each trial, clones were usually replicated in blocks. Clones measured in only one environment or only one trial were excluded. For most traits, the model was

$$y_{i,j,k,l} = \mu + \beta_i + t_{j(i)} + r_{k(i,j)} + c_l + \varepsilon_{i,j,k,l} \begin{cases} i = 1, \dots, 100 \\ l = 1, \dots, 603 \end{cases} \quad [1]$$

in which $y_{i,j,k,l}$ was the phenotype, μ was the overall mean, β_i was the fixed effect of the combination of year and location, with i varying from 1 to 100 for 100 combinations of year and location, $t_{j(i)}$ was the random effect of the trial within an environment with a normal distribution $N(0, \sigma^2_T)$, $r_{k(i,j)}$ was the random effect of the replication (or block) within an

environment with a normal distribution $N(0, \sigma^2_R)$, c_l was the effect of a clone considered random with a normal distribution $N(0, \sigma^2_C)$, without considering the additive relationship matrix as the variance–covariance matrix, and with l varying from 1 to 603 for 603 clones, and $\varepsilon_{i,j,k,l}$ was the residual considered as random and following a normal distribution $N(0, \sigma^2)$. The assumption of homogeneity of clonal variance derives from assuming all clones were sampled from the same conceptual population of the IITA breeding program. Therefore, even though different trials sampled different sets of clones, the variance was assumed consistent across trials. The assumption of homogeneity of error variance is no doubt incorrect (e.g., Edwards and Jannink, 2006), but it is assumed for expediency, as in many studies.

Some traits, such as the number of storage roots, the total fresh weight of harvested foliage and stems, and the total fresh weight of storage roots harvested, depended on the number of harvested plants: the correlation between those traits and the number of plants harvested was higher than 0.6. Because of this dependency, the number of harvested plants was taken into account in the model as a fixed effect. For these traits, the model was

$$y_{i,j,k,l,m} = \mu + \beta_i + t_{j(i)} + r_{k(i,j)} + c_l + \delta x_{m(i,j,k,l)} + \varepsilon_{i,j,k,l,m} \begin{cases} i = 1, \dots, 100 \\ l = 1, \dots, 603 \end{cases} \quad [2]$$

with the same notations as above and in which $y_{i,j,k,l,m}$ represents the phenotype and $x_{m(i,j,k,l)}$ the number of plants harvested in plot m , δ is a regression coefficient, and $\varepsilon_{i,j,k,l,m}$ is the residual considered as random and following a normal distribution $N(0, \sigma^2)$. We estimated the heritability as the ratio of the clonal variance to the sum of the clonal variance and the residuals variance.

Statistical Models for Genomic Predictions

This study used a two-step approach to make the genomic predictions. The first step consisted in generating BLUEs from all the phenotypic observations, so that each clone had a single phenotypic value for each trait. This reduced the computation time in the subsequent prediction step. By using BLUEs instead of BLUPs, there was no shrinkage attributable to the treatment of clones as random effects (Garrick et al., 2009).

The two models described below were used to generate BLUEs from the curated data. As in the models generating BLUPs, the statistical models to make BLUEs depended on whether a given trait was correlated with the number of harvested plants. For traits not correlated to the number of harvested plants, the model was

$$y_{i,j,k,l} = \mu + \beta_i + t_{j(i)} + r_{k(i,j)} + \chi_l + \varepsilon_{i,j,k,l} \begin{cases} i = 1, \dots, 93 \\ l = 1, \dots, 580 \end{cases} \quad [3]$$

For traits correlated to the number of harvested plants, the model was

$$y_{i,j,k,l,m} = \mu + \beta_i + t_{j(i)} + r_{k(i,j)} + \chi_l + \delta x_{m(i,j,k,l)} + \varepsilon_{i,j,k,l,m} \begin{cases} i = 1, \dots, 93 \\ l = 1, \dots, 580 \end{cases} \quad [4]$$

In both cases, χ_l is the effect of the clones (which is here fixed, unlike in models [1] and [2]); l varies from 1 to 580 because

these models are applied to curated clones and i varies from 1 to 93 because these models are applied to curated environments.

The BLUEs and their genotypic data in the training population were used to make genomic predictions of the validation population, using the R package rrBLUP (Endelman, 2011), which considers marker-based relationships as random effect covariates (Endelman, 2011). The statistical model to generate those genomic predictions is a mixed model, described below in matrix notation:

$$\mathbf{y} = \mu + \mathbf{Z}\mathbf{u} + \mathbf{e} \quad [5]$$

in which \mathbf{y} is the vector of phenotypes, μ is the population mean, \mathbf{u} is a vector of the genotypic values considered a random effect and following a normal distribution $N(0, \mathbf{K}\sigma_u^2)$ in which \mathbf{K} is the realized additive relationship matrix (Endelman, 2011), \mathbf{Z} is an incidence matrix for \mathbf{u} , and \mathbf{e} refers to the vector of random residuals following a normal distribution $N(0, \mathbf{I}\sigma^2)$. We measured the prediction accuracy $r(\hat{a}, y)$ as the correlation between the estimated breeding value, accounting only for additive effects (\hat{a}), and the BLUE (y).

Data Curation

In a first analysis, we sought to identify if any of the 13 locations was particularly different from the others in relative clonal performance. We removed such outlier locations to avoid excessive G×E interaction that would reduce accuracies. The software MATMODEL 3.0 (Gauch and Furnas, 1991) was used to perform the additive main effect and multiplicative interaction (AMMI) analysis. The AMMI analysis integrates additive components to explain main effects and multiplicative components to account for interactions (Zobel et al., 1988). This analysis was done on the UYT datasets, year by year, from 2006 to 2008, years for which the data were the most balanced. In addition, we developed three curation methods to identify potential clone labeling errors, that is, cases where the genotypic data did not correspond to phenotypic data from the same clone. Especially in a trial such as Genetic Gain that conserves historical data, one labeling error in the past could propagate across time and reduce the accuracy of genomic predictions.

The three curation methods all used a similar ad hoc approach: (i) regress phenotypic observations on predictors derived from independent phenotypes or from genotypes, (ii) weight the residuals of this regression inversely to trait heritabilities or accuracies and sum their absolute values across traits, and (iii) identify clones whose total residual scores are extreme relative to the global distribution. Clones that appeared to be outliers in at least two of the methods were removed from subsequent analyses.

In the first curation method, BLUPs for the Genetic Gain dataset were regressed on the BLUPs for the other datasets (UYT, PYT, and AYT). The assumption here was that BLUPs of the two types of data should be similar, unless labeling errors had occurred. This method examined phenotypic data only whereas the following methods also considered the genotypic data.

The second curation method was based on the expectation that genetically similar clones should also be phenotypically similar. For each pair of clones, we examined their relatedness in the relationship matrix \mathbf{A} . The \mathbf{A} matrix was computed by using the rrBLUP function Amat (Endelman, 2011), which

uses the matrix multiplication $\mathbf{W}\mathbf{W}'/c$, in which $W_{ik} = G_{ik} + 1 - 2p_k$, in which G_{ik} is the genotype for the i th individual at the k th marker (coded as -1 , 0 , and 1 for one homozygote, the heterozygote, and the other homozygote, respectively), p_k is the frequency of one of the alleles, and the normalization constant is $c = 2\sum_k p_k(1 - p_k)$. For each pair of clones, we also calculated the difference between their BLUPs. A linear regression of the BLUP difference on the genetic relationship was performed for each pair of clones, for each trait.

The third curation method evaluated the residuals of predictions of clone effects calculated using reproducing kernel Hilbert spaces (RKHS) regression estimates (de los Campos et al., 2009) with a Gaussian kernel function (Endelman, 2011). We chose RKHS for this purpose because it fits the training set phenotypes with a high coefficient of determination (Heslot et al., 2012). We reasoned, therefore, that large residuals from RKHS predictions within the training set could be indicative of problems in the correspondence of phenotypes to genotypes.

Cross-Validation Schemes

Accounting for Relatedness between Training and Validation Sets

Clones were assigned to clusters based on genotypic data using the k -means clustering algorithm, a method that attempts to minimize the distance between points in a cluster and the center of that cluster. We performed the k -means method using the Hartigan and Wong (1979) algorithm on marker data and generated for each trait $N = (\text{number of nonmissing individuals for a given trait}/5)$ clusters from 10 random initial cluster centers.

Two cross-validation schemes were designed to evaluate the influence of relatedness on prediction accuracy. The first scheme (cross-validation no close relatives [CV-noCR]) avoided closely related clones; that is, clones from the same cluster were not allowed in both the training and the validation sets. Therefore, CV-noCR (Fig. 1b) assigned to the validation population a sample of clusters, such that they represented 20% of the whole population. In contrast, the second cross-validation scheme (cross-validation close relatives [CV-CR]) forced close relatives between the training and validation sets (Fig. 1). That is, CV-CR (Fig. 1c) always distributed individuals in clusters to both training and validation sets. Note that in the CV-CR scheme, clusters of only one individual were never in the validation population and were not predicted. This procedure was repeated five times and we considered for each clone the mean of the five predictions.

The relationship between the training set and the validation set was measured by identifying for each individual of the validation set the 10 most closely related individuals ("top10") in the training set (Clark et al., 2012). Each validation set was then characterized by the mean relatedness of the top10 individuals. Different training-validation sets were considered until all individuals had been predicted. For a given cross-validation scheme, we considered the a_{top10} statistic, which is the mean relatedness of the top10 individuals in the validation set to those in the training set, averaged across the different training-validation sets, and finally across the five repetitions.

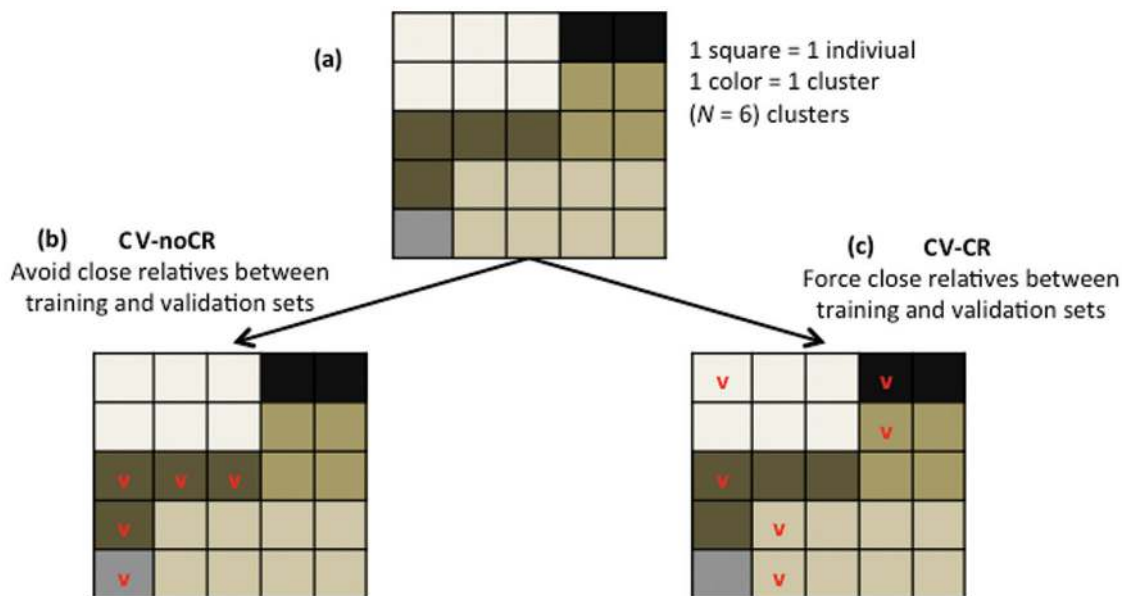


Figure 1. Cross-validation schemes taking into account the effects of relatedness. (a) The big square represents a sample of 25 individuals, where each individual is represented by a little square. Cluster membership is indicated by color. (b and c) A red “v” in a square indicates that the corresponding individual is in the validation set, and the remaining individuals are in the training set. CV-CR, cross-validation close relatives; CV-noCR, cross-validation no close relatives.

Accounting for the Genotype × Environment Interaction

In the cross-validation genotype × environment (CV-GE) scheme, the observations were split, by location, into two disjoint sets of six locations. Best linear unbiased estimators were calculated separately in each set and according to the models presented above. Five folds of 20% of the BLUEs of one set of locations were used for the validation set so that all clones in this set were predicted. Clones of the training set that also appeared in the validation set were removed from the training set, so that there were no common clones and no common locations represented between the training and the validation sets. This scheme of cross-validation was repeated 15 times with different random sets of locations. Because prediction accuracies of CV-GE used smaller training populations built on only half of the observations, we compared their results to cross-validation schemes in which observations were split randomly into two sets (CV-Random_Half), repeated five times.

Estimation of Genotype × Environment Interaction effects

To calculate the relative magnitude of G×E effects, variance components were estimated using ASReml (Gilmour et al., 2002) on data from Genetic Gain trials. The linear model was $Y_{i,j,k,l} = \mu + \beta_i + r_{k(i)} + c_l + \gamma_{i,l} + \varepsilon_{i,j,k,l}$ with terms defined similarly as for model [1], all fitted as random, with the additional term $\gamma_{i,l}$ accounting for clone × environment interaction. The covariance matrix for the additive clone effects (c_l) was proportional to the realized relationship matrix (\mathbf{A}) calculated in rrBLUP (Endelman, 2011), $c_l \sim N(0, \mathbf{A}\sigma_g^2)$. The covariance matrix for the additive clone × environment interaction effects ($\gamma_{i,l}$) was block diagonal, $\gamma_{i,l} \sim N(0, \mathbf{B}\sigma_{ge}^2)$, for effects estimated in the same environment $\mathbf{B} = \mathbf{A}$ while for effects estimated in different environments $\mathbf{B} = \mathbf{0}$. The ratio $\sigma_{ge}^2 / (\sigma_g^2 + \sigma_{ge}^2)$ reveals the magnitude of G×E interaction effects.

RESULTS

Heritabilities of the Different Traits in the Different Trials

Heritabilities for each trait, calculated using all the phenotypic data from all the trials, are shown in Table 2 (see Methods). For the two traits related to CMD, the heritabilities were high (0.66 and 0.63). Cassava mosaic disease is a trait controlled primarily by a few major genes (Lokko et al., 2005). Agronomic traits such as those related to yield and growth had much lower heritabilities, between 0.11 and 0.28. Traits related to diseases other than CMD had heritabilities lower than the ones for CMD; in the case of bacterial blight severity and incidence, heritabilities were below 0.10. Disease traits were sometimes difficult to score accurately because of the uneven spreading of inoculum or the presence or absence of the disease in a particular season or location. The heritabilities of the morphological traits were quite low (0.07 and 0.12). For root number, in particular, this could be partly explained by the high influence of the age of the plant at harvest.

In this study, the heritabilities for root dry matter content (DM), shoot weight, and fresh yield were considerably lower (between 0.11 and 0.28) than those reported by Oliveira et al. (2012), which were, respectively, 0.67, 0.83 and 0.76. Furthermore, the heritability values that we obtained were lower than expected for many other traits. Calculation of heritabilities by trial type (UYT, AYT, PYT, and Genetic Gain) showed that the heritabilities in the Genetic Gain dataset were generally lower than those in the other datasets, decreasing the overall heritability (Fig. 2). The Genetic Gain trials used smaller plots and

Table 2. Heritabilities of cassava traits of interest, accuracies of prediction, and mean relatedness of the top10 individuals in the validation set to those in the training set (a_{top10}) of different cross-validation (CV) schemes.

Trait [†]	Heritability	CV without close relatives		CV-Random [‡] fivefold		CV with close relatives	
		Accuracy	a_{top10}	Accuracy	a_{top10}	Accuracy	a_{top10}
MCMDI	0.66	0.417	0.231	0.487	0.260	0.474	0.267
MCMDS	0.63	0.462	0.23	0.503	0.261	0.513	0.266
MCADI	0.38	0.177	0.267	0.184	0.201	0.202	0.285
DM	0.28	0.459	0.229	0.482	0.258	0.477	0.268
SPROUT	0.28	0.259	0.23	0.304	0.260	0.306	0.268
HI	0.27	0.431	0.23	0.483	0.259	0.479	0.268
FYLD	0.26	0.358	0.231	0.407	0.261	0.395	0.267
DYLD	0.21	0.231	0.23	0.304	0.259	0.296	0.268
TYLD	0.2	0.195	0.229	0.251	0.260	0.232	0.267
MCGM	0.18	0.47	0.231	0.308	0.259	0.501	0.267
MCADS	0.17	0.145	0.266	0.177	0.201	0.207	0.284
VIGOR	0.17	0.277	0.23	0.494	0.259	0.299	0.268
RTNO	0.14	0.342	0.228	0.399	0.260	0.384	0.267
RTWT	0.14	0.308	0.23	0.352	0.260	0.349	0.267
NKLG	0.12	0.202	0.233	0.190	0.258	0.195	0.267
SHTWT	0.11	0.228	0.231	0.299	0.260	0.297	0.266
MCBBS	0.09	0.266	0.229	0.303	0.260	0.316	0.267
ROTNO	0.07	0.188	0.241	0.211	0.252	0.229	0.276
MCBBI	0.04	0.238	0.229	0.255	0.260	0.26	0.266

[†]MCMDI, mean cassava mosaic disease incidence; MCMDS, mean cassava mosaic disease severity; MCADI, mean cassava anthracnose disease incidence; DM, root dry matter content; SPROUT, sprouting; HI, harvest index; FYLD, fresh root yield; DYLD, dry yield; TYLD, the top yield; MCGM, mean cassava green mite; MCADS, mean cassava anthracnose disease severity; VIGOR, initial vigor; RTNO, root number; RTWT, fresh weight of storage root; NKLG, root neck length; SHTWT, fresh shoot weight; MCBBS, mean cassava bacterial blight severity; ROTNO, rotted storage roots; MCBBI, mean cassava bacterial blight incidence.

[‡]CV-Random, random cross-validation.

were expected to have a larger environmental variance than the other trial types.

In addition to the trial type, we hypothesized that two other factors might contribute to the low heritabilities. First, there might be some outlier locations that produced G×E interactions. Because we aimed to have a training population whose predictions would be generally useful across locations, we identified and removed locations where clones behaved very differently. The AMMI results for UYT data showed that clone effects for Onne location had a strong negative correlation with clone effects from the other locations in 2006 (Fig. 3) and 2007 (not shown), 2 out of the 3 yr that were examined. Indeed, the score for Onne on the first principal component axis of the principal component analysis of the G×E interaction effects was strongly negative relative to the scores of the other locations. Consequently, the Onne location, which is one of the highest rainfall areas of Nigeria, was removed from subsequent analyses.

Finally, because of the long-term historical nature of the Genetic Gain collection, some labeling errors might have occurred, so that genotype and phenotypic data were incorrectly associated with the same clone name. We looked for possible labeling errors using three approaches: (i) comparing types of trial (Fig. 4A), (ii) regressing BLUP differences of all pairs of clones on their genetic relationship (Fig. 4B), and (iii) regressing the predicted genotypic values on the calculated BLUPs (Fig. 4C). For each approach, we plotted the distribution of the residuals and arbitrarily

identified outliers in the distribution tail (Fig. 4A, 4B, and 4C). The three curation methods identified, respectively, 22, 23, and 4 outliers. In total, 23 clones were identified as outliers in at least two out of three curation methods.

The curation method comparing types of trial could only be used if clones had data in the Genetic Gain population and in at least one of the other datasets. Consequently, some potential outlier clones may not have been identified. We removed the clones identified as outliers for the following analyses.

Cross-Validation Schemes to Account for Relatedness

Principal components analysis of the SNP data provided no evidence of genetic structure in the Genetic Gain population (data not shown). However, many pairs of clones were genetically very close, possibly because clones had been renamed by farmers and collected as distinct accessions (Fig. 5). To assess the effect of relatedness on prediction accuracies, we tried different methods to identify individuals that were closely related (while the dendrogram provides a visual representation of the problem, it does not provide a basis for assignment of clones to groups). The *k*-means algorithm created different clusters where many of those clusters contained two clones. These clusters of genetically close individuals were then used to design the training and validation sets in the cross-validation schemes.

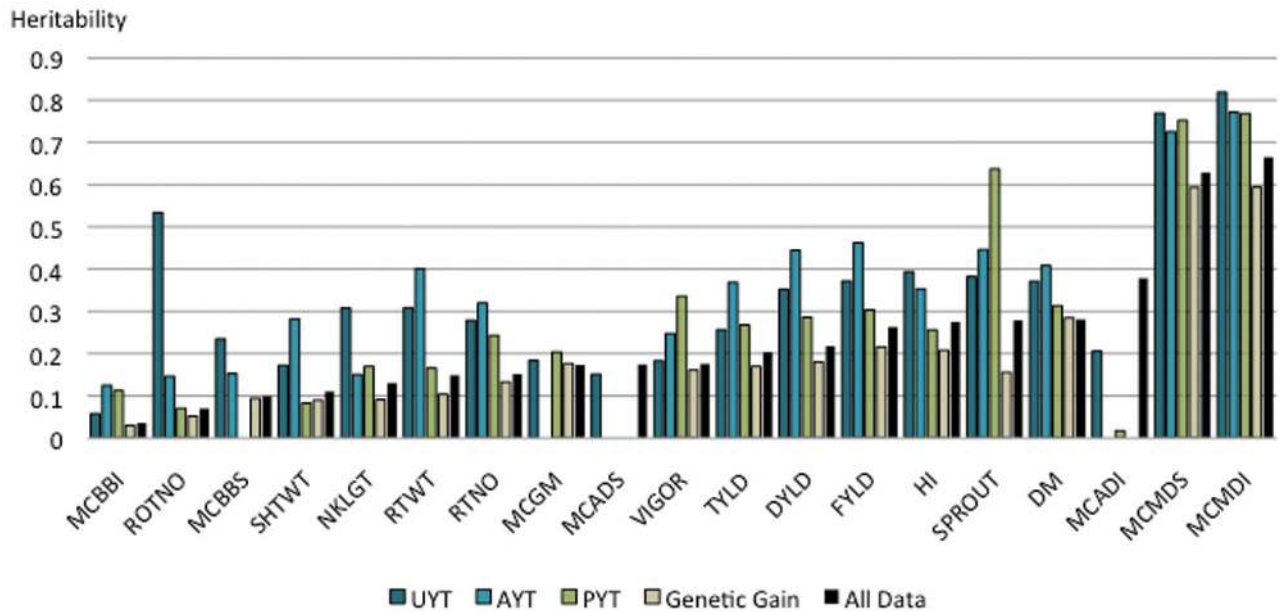


Figure 2. Broad-sense heritabilities of the different traits in the different trial types. MCBBI, mean cassava bacterial blight incidence; ROTNO, rotted storage roots; MCBBS, mean cassava bacterial blight severity; SHTWT, fresh shoot weight; NKLGT, root neck length; RTWT, fresh weight of storage root; RTNO, root number; MCGM, mean cassava green mite; MCADS, mean cassava anthracnose disease severity; VIGOR, initial vigor; TYLD, the top yield; DYLD, dry yield; FYLD, fresh root yield; HI, harvest index; SPROUT, sprouting; DM, root dry matter content; MCADI, mean cassava anthracnose disease incidence; MCMDS, mean cassava mosaic disease severity; MCMDI, mean cassava mosaic disease incidence; UYT, uniform yield trial; AYT, advanced yield trial; PYT, preliminary yield trial.

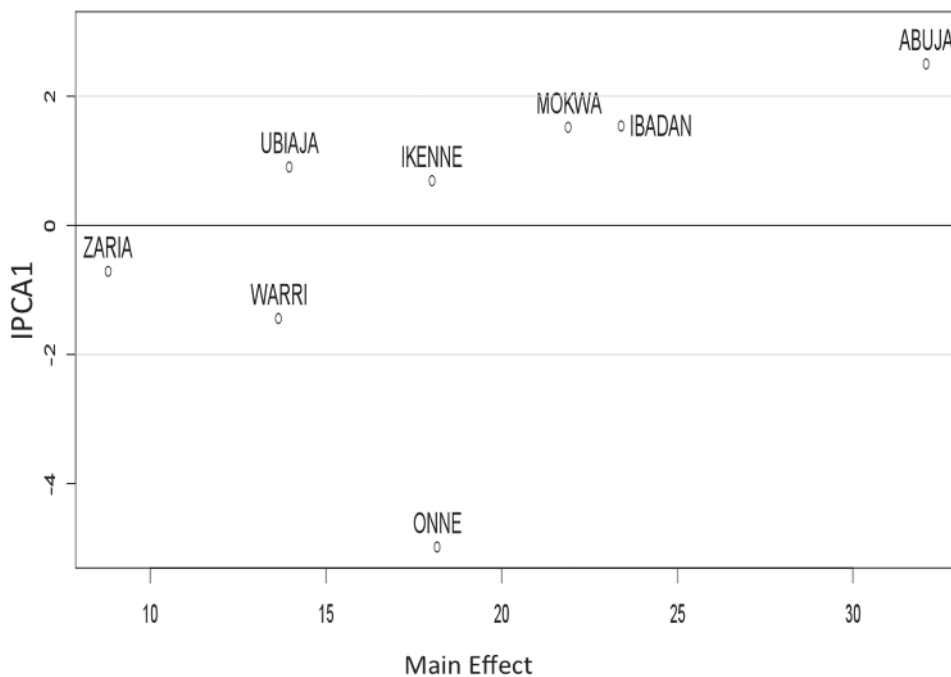


Figure 3. Identification of outlier locations. Additive main effect and multiplicative interaction (AMMI) graph of the main effects and the first axis of the principal component analysis (IPCA1) obtained by AMMI analysis on uniform yield trial data in 2006.

As explained in the Methods, we generated cross-validation sets according to two different schemes: CV-noCR, avoiding close relatives between training and validation sets, and CV-CR forcing close relatives between those sets. We calculated a statistic, a_{top10} , to measure the relatedness between the training and validation sets. As expected, this statistic was lowest

in CV-noCR and highest in CV-CR (Table 2). The variance of the a_{top10} across individuals within the top10 (see Methods) for each cross-validation scheme and for each trait was about 10^{-3} . The a_{top10} of CV-Random, while intermediate, was most often closer to the a_{top10} of CV-CR than to that of CV-noCR. In several cases, the relatedness statistics for the CV-Random and the CV-CR

were equal. This confirmed that the CV-Random scheme created validation sets whose members had close relatives in the training population.

When we compared the prediction accuracies for the different cross-validation schemes, we found that higher accuracies were associated with higher relatedness measurements. The prediction accuracies between the validation schemes CV-Random and CV-CR were very close. Across traits, the CV-noCR scheme showed lower accuracies compared with the other two.

Cross-Validation Scheme to Account for Genotype × Environment Interaction

We tested another cross-validation scheme, CV-GE, to assess the impact of G×E interaction as a confounding factor. We were specifically interested in the G×E interaction explained by locations, because the variability across years cannot be experimentally controlled. Because splitting locations between training and validation sets resulted in a training set containing only about half of the total observations, we compared the CV-GE scheme with a random k -fold scheme that used only a random half of the observations in the training set (CV-Random_Half). As expected, the prediction accuracies decreased for CV-Random_Half relative to CV-Random (Fig. 6). In the CV-GE scheme, the training population used the G×E interaction effects of 6 out of 12 locations (across all 12 yr) to predict the six other locations. In 17 of 19 traits evaluated, accuracies for CV-GE (red crosses in Fig. 6) were lower than accuracies for CV-Random_Half (blue triangles in Fig. 6). The loss of the G×E interaction effects of the six remaining locations reduced the prediction accuracies, indicating that for most traits there were shared genotype × location residuals across training and validation sets, biasing estimated prediction accuracies upward in CV-Random. Furthermore, the a_{top10} values were close to what was obtained using CV-Random or CV-CR. Therefore, the cross-validation avoiding G×E interaction did not avoid close relatives between the training and the validation set and only evaluated the impact of G×E interaction.

DISCUSSION

The curation work presented here was motivated by the low heritabilities observed when using all datasets. The curation work eliminated some clones but did not have much impact on the heritabilities. It appeared that these heritabilities were low because of the low heritabilities in the Genetic Gain data, which might be explained by the field-plot design and by the long-term maintenance of the Genetic Gain collection. Some noise was possibly introduced by phenotyping-protocol variation across time. It is also possible that epigenetic variation between generations of vegetative propagation and/or some somatic mutations might have occurred across time, causing variation in clones' phenotypes

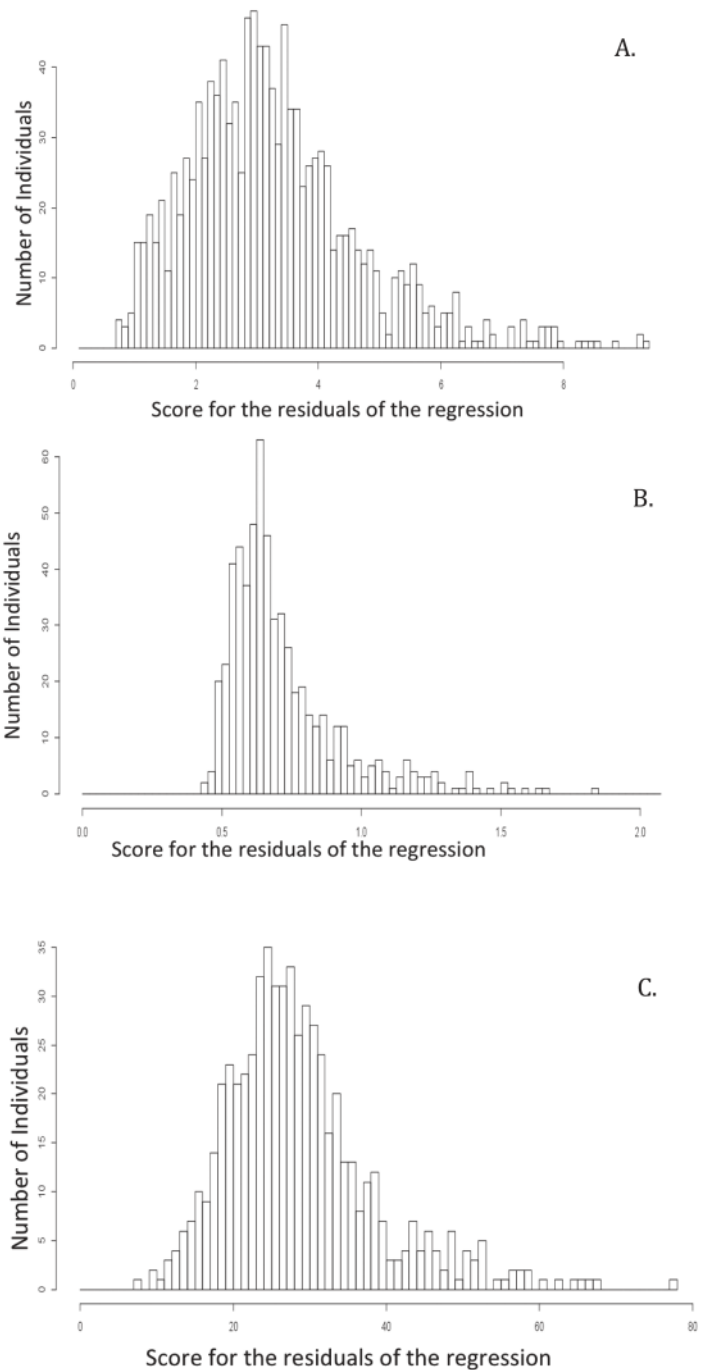


Figure 4. Distribution of the sum of the absolute value of the residuals, across traits. A. Comparing the best linear unbiased predictors between the Genetic Gain and the other trial types, weighted by heritabilities. B. Comparing the phenotypic data to the genotypic data, weighted by heritabilities. C. For predictions using the Gaussian kernel, weighted by accuracies.

(McKey et al., 2010). Indeed, IITA established the Genetic Gain population not for direct breeding purposes but to maintain clones. Rigorous phenotyping of the Genetic Gain population was therefore not a priority although having this data has been critical to start a genomic selection program. The next steps for implementation of GS will have data from larger and more replicated plot designs than are currently

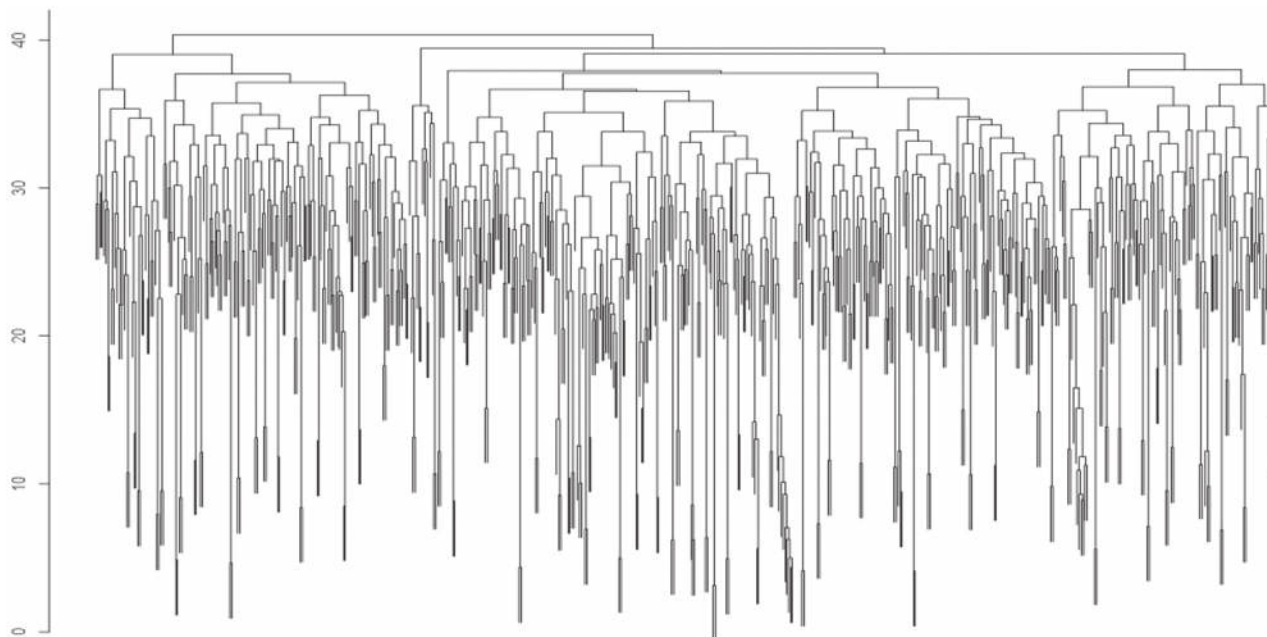


Figure 5. Dendrogram of the genotypic data. A hierarchical clustering using the Euclidean distance between the genotypes was used to represent a dendrogram of the clones, and allowed capturing visually the strongest pattern to represent in our clustering methods.

available. As the quality of the phenotype data improves, so should heritabilities and prediction accuracies (Pszczola et al., 2012).

While traits with higher heritabilities tended to show higher prediction accuracies (Table 2), this trend was fairly weak ($r^2 = 0.26$) and there were some striking exceptions. For example, prediction accuracies for harvest index and cassava mosaic disease incidence (CMDI) were similar although heritability for CMDI was much higher. Some discrepancies were presumably caused by differences in genetic architecture and the proximity of SNP markers to quantitative trait loci. In the case of CMD traits, selection history might have reduced accuracy; CMD has been a target of very strong and directed selection. Because it is affected primarily by major genes, efforts at resistance-allele introgression could cause two clones with divergent genetic backgrounds to have similar disease resistance (or, conversely, clones with similar background could be divergent in their resistance). These cases would make CMD resistance difficult for GS models to predict.

Prediction accuracies are a correlation between BLUEs, which include both additive and nonadditive components, and the estimated breeding values, which only consider the additive component; this causes a systematic underestimation of prediction accuracies. That said, accuracies for most of the traits studied are not currently high enough for GS to be effective in cassava breeding. If the breeding cycle is reduced from 5 to 2 yr, prediction accuracies must be at least 0.4 to match or exceed gain from phenotypic selection. However, as noted above, anticipated improvement in heritabilities will increase accuracy, as will larger training population size and higher marker density.

Testing the CV-noCR and CV-CR cross-validation schemes showed that the more closely the training and the validation sets were related, the higher the prediction accuracies (Table 2). The mean accuracy of CV-Random was quite close to that of CV-CR, showing that the presence of close relatives in the population would bias CV-Random accuracy estimates upward relative to a realistic expectation in a selection program. Accounting for relatedness between the populations might then be a way to correct this potential overestimation (Cleveland et al., 2012). Nevertheless, our study showed that the differences in prediction accuracies between CV-noCR and CV-CR were low. The very closely related pairs of clones did not affect strongly the prediction accuracies, possibly because the number of those pairs was relatively small. We note also that the a_{top10} values did not differ greatly between CV-CR and CV-noCR, suggesting an overall high level of relatedness among clones in the IITA population. Caution should therefore be exercised in using this population for prediction in other populations that might well not be within its domain of inference.

The differences in prediction accuracies that we observed were much lower than those reported in the Pérez-Cabal et al. (2012) and Cleveland et al. (2012) animal studies. This might be because those studies evaluated relatedness based on pedigree whereas our study evaluated relatedness based on genotypic information. Studies showed that prediction accuracies increased when using marker information instead of pedigree if the marker density was high enough (Villanueva et al., 2005; Hayes et al., 2009b). Similarly, high-density markers should be able to generate more distinct training and validation sets than

Prediction accuracies in the different cross validation schemes

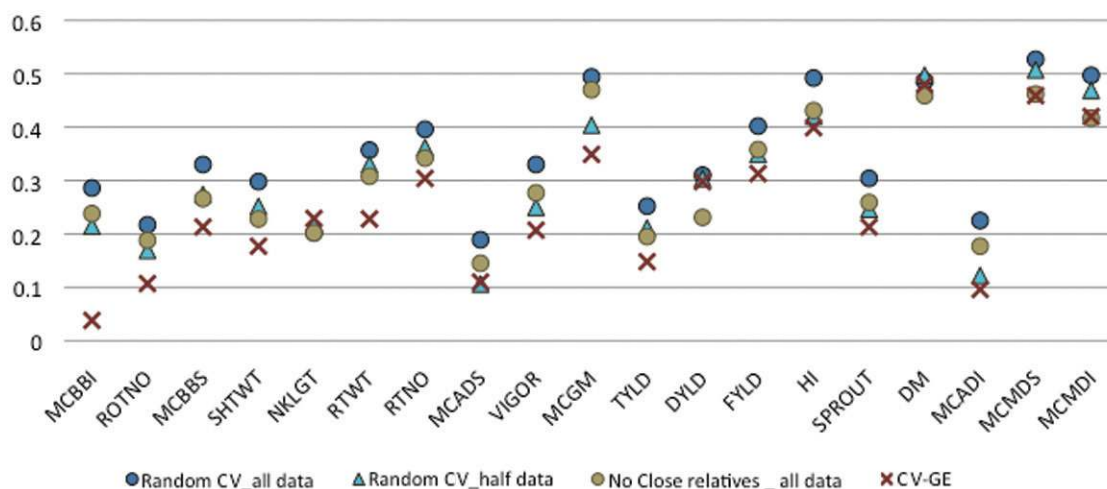


Figure 6. Comparison of the prediction accuracies of different cross-validation schemes. Traits are ranked according to their heritabilities, from lower (on the left) to higher. MCBBI, mean cassava bacterial blight incidence; ROTNO, rotted storage roots; MCBBS, mean cassava bacterial blight severity; SHTWT, fresh shoot weight; NKLG1, root neck length; RTWT, fresh weight of storage root; RTNO, root number; MCADS, mean cassava anthracnose disease severity; VIGOR, initial vigor; MCGM, mean cassava green mite; TYLD, the top yield; DYLD, dry yield; FYLD, fresh root yield; HI, harvest index; SPROUT, sprouting; DM, root dry matter content; MCADI, mean cassava anthracnose disease incidence; MCMDS, mean cassava mosaic disease severity; MCMDI, mean cassava mosaic disease incidence; CV, cross-validation; CV-GE, cross-validation genotype \times environment.

use of pedigree information. However, our study used marker data at a relatively low density, compared with other studies, for example, in dairy cattle or maize (*Zea mays* L.). Therefore, the use of a relatively small number of markers to assess relatedness in the training and validation sets might explain the lower difference in prediction accuracies between CV-noCR and CV-CR compared with the difference assessed in animal studies.

Moreover, the size of the population and its diversity influence the impact of close relatives in the training set. Indeed Pérez-Cabal et al. (2012) and Cleveland et al. (2012) selected the individuals for their studies so that their populations were not only quite large but also quite diverse. Because cassava is propagated vegetatively, farmers do not propagate all clones at the same rate, and in the long term some clones are likely to be lost (McKey et al., 2010). It is not clear whether the genetic bottleneck of African cassava, because of its introduction from the Amazonian region (Jones, 1959), has been more severe than that experienced by domesticated cattle. These factors might cause a lower diversity compared with animals and thus reduce the impact of close relatives in the training set on prediction accuracies.

Even if, for cassava breeding programs, the relatedness between the training and the validation sets does not have as much impact as for some animal cases, it may still be worth taking into account that it could affect prediction accuracy. This analysis should consider the level of relatedness expected between training population and selection candidates in the initial generations of the breeding cycle. This contrasts with the relatedness between training and validation sets in the CV-noCR

scheme, both of which included individuals from the same generations, that is, all the years of the Genetic Gain program. A high level of relatedness within the training population might depend on the germplasm used at the beginning and the selection history of the population. The results of our study suggest that it would be interesting in a genomic selection breeding program to compare the a_{top10} assessed by a k -fold random cross-validation within the training population and the actual a_{top10} between the training population to the selection candidates, when the training population contains parents of the validation population. We measured the mean relatedness between parent and offspring for nine clones for which pedigree data were available and obtained a value of 0.30 ± 0.01 . If the relatedness assessed in a k -folds cross-validation were higher than the mean relatedness expected between parent and offspring, we might overestimate the prediction accuracy. In our study, the a_{top10} using the CV-noCR scheme (Table 2) was lower than this value estimated for parents and offspring, so the CV-noCR scheme probably did not overestimate prediction accuracy. Note that 0.30 is not the parent-offspring relatedness expected from methods of calculation using pedigree (which would be 0.5). Coefficients calculated from marker data are not comparable to those calculated from pedigree (e.g., they can be negative; Endelman and Jannink, 2012).

When we analyzed the effect of G \times E interaction, we found that the prediction accuracies of CV-GE were generally lower than those obtained by CV-Random_Half (Fig. 6). When clones were evaluated in the same locations for both the training and validation sets, G \times E

interaction, which has been shown to have a significant effect on traits related to yield in Nigerian cassava (Aina et al., 2009), seemed to be a strong confounding factor that leads to overestimation of prediction accuracies. There was considerable variation in the magnitude of this effect for different traits; for example, accuracy for mean cassava bacterial blight incidence (MCBBI) dropped almost to zero while accuracy for DM was unaffected. To test whether the loss of accuracy correlated to the magnitude of G×E interaction effects, we estimated an additive × environment interaction effect for each trait (see Methods). The ratio of additive × environment variance compared with additive-genetic plus additive × environment variance varied from 0.84 for MCBBI to 0.10 for DM, consistent with the hypothesis that traits with a smaller ratio should show a smaller reduction in accuracy in the G×E interaction cross-validation scheme. Across all the traits, the linear correlation between accuracy reduction and magnitude of G×E interaction was significantly positive ($P = 0.04$).

The presence of G×E interaction effects reduces our ability to make predictions when selecting for locations where no evaluations have been done previously. In those cases, the phenotypic observations are likely not to be as correlated to the predictions as the random cross-validation prediction accuracies. When expanding a genomic selection breeding program to new locations, to have a better estimation of our prediction ability, cross-validations schemes should aim at reducing the overestimation caused by G×E interaction by using training and validation sets that do not share common locations. Furthermore, given the impact of the G×E interaction on the prediction accuracies, instead of removing the G×E interaction effect, exploiting it would be a worthwhile goal. It may be worth delineating mega-environments to make predictions within them, thus exploiting the narrow adaptations of genotypes in those mega-environments (Gauch, 1997; Annicchiarico, 2002).

CONCLUSIONS

Prediction accuracies obtained by random cross-validations, used to evaluate the prospects for success of genomic selection, will be overestimated if there are close relatives in the training population. Relatedness should therefore be examined in a genomic selection breeding process to better evaluate prediction accuracies and should be considered in designing the training population. Genotype × environment interactions also contribute to overestimation of prediction accuracies and should be considered when expanding a breeding program to new experimental sites. Prediction accuracies need improvement if GS is to outperform phenotypic selection on a per-year basis; these improvements are expected as training populations increase in size and are less dependent on historical phenotype data.

Acknowledgments

DL and JLJ designed the study and interpreted the results. DL performed the statistical analyses. AGOD, PK, IYR, RO, and MB selected and assembled the genetic gain collection and contributed to phenotypic evaluation. GM oversaw cassava tissue collection and DNA preparation. MTH was responsible for the GBS data and bioinformatic analyses. HGG performed the AMMI analysis. DL wrote the paper with assistance from MTH and JLJ. This work was supported by the project “Genomic Selection: The next frontier for rapid gains in maize and wheat improvement”, through funds from The Bill and Melinda Gates Foundation. We thank Deniz Akdemir and Jeff Endelman for help with statistical analyses, Lisa Blanchard for preparation of GBS libraries, and Jeff Glaubitz, Rob Elshire, and Qi Sun for help with the GBS bioinformatics pipeline.

References

- Aina, O.O., A.G.O. Dixon, I. Paul, and E.A. Akinrinde. 2009. G × E interaction effects on yield and yield components of cassava (landraces and improved) genotypes in the savanna regions of Nigeria. *Afr. J. Biotechnol.* 8:4933–4945.
- Annicchiarico, P. 2002. Case study: Durum wheat. In: Genotype × environment interactions: Challenges and opportunities for plant breeding and cultivar recommendations. FAO, Rome, Italy. p. 89–104.
- Breiman, L.E.O. 2001. Random forests. *Statistics* 45:5–32.
- Burgueño, J., G. de los Campos, K. Weigel, and J. Crossa. 2012. Genomic prediction of breeding values when modeling genotype × environment interaction using pedigree and dense molecular markers. *Crop Sci.* 52:707–719. doi:10.2135/cropsci2011.06.0299
- Clark, S.A., J.M. Hickey, H.D. Daetwyler, and J.H.J. van der Werf. 2012. The importance of information on relatives for the prediction of genomic breeding values and the implications for the makeup of reference data sets in livestock breeding schemes. *Genet. Sel. Evol.* 44:4. doi:10.1186/1297-9686-44-4
- Cleveland, M.A., J.M. Hickey, and S. Forni. 2012. A common dataset for genomic analysis of livestock populations. *G3* 2:429–35.
- Comstock, R.E., and R.H. Moll. 1963. Genotype-environment interactions. In: W.D. Hanson and H.F. Robinson, editors, *Statistical genetics and plant breeding*. National Academy of Sciences-National Research Council (U.S.) Publ. 982. National Academies, Washington, DC. p. 164–196.
- Crossa, J., G.D.L. Campos, P. Pérez, D. Gianola, J. Burgueño, J.L. Araus, D. Makumbi, R.P. Singh, S. Dreisigacker, J. Yan, V. Arief, M. Banziger, and H.-J. Braun. 2010. Prediction of genetic values of quantitative traits in plant breeding using pedigree and molecular markers. *Genetics* 186:713–724. doi:10.1534/genetics.110.118521
- de los Campos, G., D. Gianola, and G.J.M. Rosa. 2009. Reproducing kernel Hilbert spaces regression: A general framework for genetic evaluation. *J. Anim. Sci.* 87:1883–1887. doi:10.2527/jas.2008-1259
- Edwards, J.W., and J.-L. Jannink. 2006. Bayesian modeling of heterogeneous error and genotype by environment interaction variances. *Crop Sci.* 46:820–833. doi:10.2135/cropsci2005.0164
- Elshire, R.J., J.C. Glaubitz, Q. Sun, J.A. Poland, K. Kawamoto, E.S. Buckler, and S.E. Mitchell. 2011. A robust, simple

- genotyping-by-sequencing (GBS) approach for high diversity species. *PLoS ONE* 6:e19379. doi:10.1371/journal.pone.0019379
- Endelman, J.B. 2011. Ridge regression and other kernels for genomic selection with R package rrBLUP. *Plant Gen.* 4:250–255. doi:10.3835/plantgenome2011.08.0024
- Endelman, J.B., and J.-L. Jannink. 2012. Shrinkage estimation of the realized relationship matrix. *G3* 2:1405–1413.
- Garrick, D.J., J.F. Taylor, and R.L. Fernando. 2009. Deregressing estimated breeding values and weighting information for genomic regression analyses. *Genet. Sel. Evol.* 41:55. doi:10.1186/1297-9686-41-55
- Gauch, H.G. 1997. Identifying mega-environments and targeting genotypes. *Crop Sci.* 37:311–326. doi:10.2135/cropsci1997.0011183X003700020002x
- Gauch, H.G., and R.E. Furnas. 1991. Statistical analysis of yield trials with MATMODEL. *Agron. J.* 83:916–920. doi:10.2134/agronj1991.00021962008300050027x
- Gilmour, A.R., B.J. Gogel, B.R. Cullis, S.J. Welham, and R. Thompson. 2002. ASReml user guide release 1.0. VSN International Ltd., Hemel Hempstead, UK.
- Glaubitz, J., T. Casstevens, R. Elshire, J. Harriman, and E.S. Buckler. 2012. TASSEL 3.0 genotyping by sequencing (GBS) pipeline documentation. Edward S. Buckler, USDA-ARS, Ithaca, NY. <http://www.maizegenetics.net/tassel/docs/TasselPipelineGBS.pdf> (accessed 6 Oct. 2011).
- Goddard, M.E., and B.J. Hayes. 2007. Genomic selection. *J. Anim. Breed. Genet.* 124:323–330. doi:10.1111/j.1439-0388.2007.00702.x
- Habier, D., R.L. Fernando, and J.C.M. Dekkers. 2007. The impact of genetic relationship information on genome-assisted breeding values. *Genetics* 177:2389–2397.
- Habier, D., J. Tetens, F.-R. Seefried, P. Lichtner, and G. Thaller. 2010. The impact of genetic relationship information on genomic breeding values in German Holstein cattle. *Genet. Sel. Evol.* 42:5. doi:10.1186/1297-9686-42-5
- Hartigan, J.A., and M.A. Wong. 1979. Algorithm AS 136: A K-means clustering algorithm. *Appl. Stat* 28:100–108.
- Hayes, B.J., P.J. Bowman, A.J. Chamberlain, and M.E. Goddard. 2009a. Genomic selection in dairy cattle: Progress and challenges. *J. Dairy Sci.* 92:433–443. doi:10.3168/jds.2008-1646
- Hayes, B.J., P.M. Visscher, and M.E. Goddard. 2009b. Increased accuracy of artificial selection by using the realized relationship matrix. *Genet. Res.* 91:47–60. doi:10.1017/S0016672308009981
- Heffner, E.L., M.E. Sorrells, and J.-L. Jannink. 2009. Genomic selection for crop improvement. *Crop Sci.* 49:1–12. doi:10.2135/cropsci2008.08.0512
- Heslot, N., H.-P. Yang, M.E. Sorrells, and J.-L. Jannink. 2012. Genomic selection in plant breeding: A comparison of models. *Crop Sci.* 52:146–160. doi:10.2135/cropsci2011.06.0297
- Jannink, J.-L., A.J. Lorenz, and H. Iwata. 2010. Genomic selection in plant breeding: From theory to practice. *Brief. Funct. Gen.* 9:166–177. doi:10.1093/bfgp/elq001
- Jones, W. 1959. *Manioc in Africa*. Stanford Univ. Press, Stanford, CA.
- Kawuki, R.S., M. Ferguson, M.T. Labuschagne, L. Herselman, J. Orone, I. Ralimanana, M. Bidiaka, S. Lukombo, M.C. Kanyange, G. Gashaka, G. Mkamilo, J. Gethi, and H. Obiero. 2011. Variation in qualitative and quantitative traits of cassava germplasm from selected national breeding programmes in sub-Saharan Africa. *Field Crops Res.* 122:151–156. doi:10.1016/j.fcr.2011.03.006
- Kohavi, R. 1995. A study of cross-validation and bootstrap for accuracy estimation and model selection. In: *Proceedings of the 14th International Joint Conference on Artificial Intelligence*, Montreal, QC, Canada. 20–25 Aug. 1995. Morgan Kaufmann Publishers Inc. San Francisco, CA. 2:1137–1143.
- Lokko, Y., E.Y. Danquah, S.K. Offei, A.G.O. Dixon, and M.A. Gedil. 2005. Molecular markers associated with a new source of resistance to the cassava mosaic disease. *Afr. J. Biotechnol.* 4:873–881.
- Lorenz, A.J., S. Chao, F.G. Asoro, E.L. Heffner, T. Hayashi, H. Iwata, K.P. Smith, M.E. Sorrells, and J. Jannink. 2011. Genomic selection in plant breeding: Knowledge and prospects. *Adv. Agron.* 110:77–123. doi:10.1016/B978-0-12-385531-2.00002-5
- Lorenz, A.J., K.P. Smith, and J.-L. Jannink. 2012. Potential and optimization of genomic selection for *Fusarium* head blight resistance in six-row barley. *Crop Sci.* 52:1609–1621. doi:10.2135/cropsci2011.09.0503
- Massman, J.M., A. Gordillo, R.E. Lorenzana, and R. Bernardo. 2012. Genomewide predictions from maize single-cross data. *Theor. Appl. Genet.* 126:13–22. doi:10.1007/s00122-012-1955-y
- Maziya-Dixon, B., A.G.O. Dixon, and A.-R.A. Adebawale. 2007. Targeting different end uses of cassava: Genotypic variations for cyanogenic potentials and pasting properties. *Int. J. Food Sci. Technol.* 42:969–976. doi:10.1111/j.1365-2621.2006.01319.x
- McKey, D., M. Elias, B. Pujol, and A. Duputié. 2010. The evolutionary ecology of clonally propagated domesticated plants. *New Phytol.* 186:318–332. doi:10.1111/j.1469-8137.2010.03210.x
- Melchinger, A.E., H.F. Utz, and C.C. Schön. 2004. QTL analyses of complex traits with cross validation, bootstrapping and other biometric methods. *Euphytica* 137:1–11. doi:10.1023/B:EUPH.0000040498.48379.68
- Metzker, M.L. 2010. Sequencing technologies – The next generation. *Nat. Rev. Genet.* 11:31–46. doi:10.1038/nrg2626
- Meuwissen, T.H., B.J. Hayes, and M.E. Goddard. 2001. Prediction of total genetic value using genome-wide dense marker maps. *Genetics* 157:1819–1829.
- Okechukwu, R.U., and A.G.O. Dixon. 2008. Genetic gains from 30 years of cassava breeding in Nigeria for storage root yield and disease resistance in elite cassava genotypes. *J. Crop Improv.* 22:181–208. doi:10.1080/15427520802212506
- Oliveira, E.J., M.D.V. Resende, V. Silva Santos, C.F. Ferreira, G.A.F. Oliveira, M.S. Silva, L.A. Oliveira, and C.I. Aguilar-Vildoso. 2012. Genome-wide selection in cassava. *Euphytica* 187:263–276. doi:10.1007/s10681-012-0722-0
- Pérez-Cabal, M.A., A.I. Vazquez, D. Gianola, G.J.M. Rosa, and K.A. Weigel. 2012. Accuracy of genome-enabled prediction in a dairy cattle population using different cross-validation layouts. *Front. Genet.* 3:27. doi:10.3389/fgene.2012.00027
- Poland, J., J. Endelman, J. Dawson, J. Rutkoski, S. Wu, S. Dreisigacker, J. Crossa, H. Sanchez-villeda, and M. Sorrells. 2012. Genomic selection in wheat breeding using genotyping-by-sequencing. *Plant Gen.* 5:103–113. doi:10.3835/plantgenome2012.06.0006
- Pszczola, M., T. Strabel, H.A. Mulder, and M.P.L. Calus. 2012. Reliability of direct genomic values for animals with different relationships within and to the reference population. *J. Dairy Sci.* 95:389–400. doi:10.3168/jds.2011-4338

- R Development Core Team. 2008. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <http://www.R-project.org> (accessed 15 Mar. 2012).
- Riedelsheimer, C., A. Czedik-Eysenberg, C. Grieder, J. Lisec, F. Technow, R. Sulpice, T. Altmann, M. Stitt, L. Willmitzer, and A.E. Melchinger. 2012. Genomic and metabolic prediction of complex heterotic traits in hybrid maize. *Nat. Genet.* 44:217–220. doi:10.1038/ng.1033
- Schön, C.C., H.F. Utz, S. Groh, B. Truberg, S. Openshaw, and A.E. Melchinger. 2004. Quantitative trait locus mapping based on resampling in a vast maize testcross experiment and its relevance to quantitative genetics for complex traits. *Genetics* 167:485–498. doi:10.1534/genetics.167.1.485
- VanRaden, P.M., C.P. Van Tassell, G.R. Wiggans, T.S. Sonstegard, R.D. Schnabel, J.F. Taylor, and F.S. Schenkel. 2009. Invited review: Reliability of genomic predictions for North American Holstein bulls. *J. Dairy Sci.* 92:16–24. doi:10.3168/jds.2008-1514
- Villanueva, B., J. Fernández, M.A. Toro, and J. Ferna. 2005. Benefits from marker-assisted selection under an additive polygenic genetic model. *J. Anim. Sci.* 83:1747–1752.
- Windhausen, V.S., G.N. Atlin, J.M. Hickey, J. Crossa, J.-L. Jannink, M.E. Sorrells, B. Raman, J.E. Cairns, A. Tarekegne, K. Semagn, Y. Beyene, P. Grudloyma, F. Technow, C. Riedelsheimer, and A.E. Melchinger. 2012. Effectiveness of genomic prediction of maize hybrid performance in different breeding populations and environments. *G3* 2:1427–1436.
- Zobel, R.W., M.J. Wright, and H.G. Gauch. 1988. Statistical analysis of a yield trial. *Agron. J.* 80:388–393. doi:10.2134/ajonj1988.00021962008000030002x