

## ORIGINAL ARTICLE

# Modeling additive and non-additive effects in a hybrid population using genome-wide genotyping: prediction accuracy implications

J-M Bouvet<sup>1</sup>, G Makouanzi<sup>2</sup>, D Cros<sup>1</sup> and Ph Vignerol<sup>1,2</sup>

Hybrids are broadly used in plant breeding and accurate estimation of variance components is crucial for optimizing genetic gain. Genome-wide information may be used to explore models designed to assess the extent of additive and non-additive variance and test their prediction accuracy for the genomic selection. Ten linear mixed models, involving pedigree- and marker-based relationship matrices among parents, were developed to estimate additive (A), dominance (D) and epistatic (AA, AD and DD) effects. Five complementary models, involving the gametic phase to estimate marker-based relationships among hybrid progenies, were developed to assess the same effects. The models were compared using tree height and 3303 single-nucleotide polymorphism markers from 1130 cloned individuals obtained via controlled crosses of 13 *Eucalyptus urophylla* females with 9 *Eucalyptus grandis* males. Akaike information criterion (AIC), variance ratios, asymptotic correlation matrices of estimates, goodness-of-fit, prediction accuracy and mean square error (MSE) were used for the comparisons. The variance components and variance ratios differed according to the model. Models with a parent marker-based relationship matrix performed better than those that were pedigree-based, that is, an absence of singularities, lower AIC, higher goodness-of-fit and accuracy and smaller MSE. However, AD and DD variances were estimated with high s.e.s. Using the same criteria, progeny gametic phase-based models performed better in fitting the observations and predicting genetic values. However, DD variance could not be separated from the dominance variance and null estimates were obtained for AA and AD effects. This study highlighted the advantages of progeny models using genome-wide information.

*Heredity* (2016) **116**, 146–157; doi:10.1038/hdy.2015.78; published online 2 September 2015

## INTRODUCTION

Quantitative genetic models partition genetic variance into additive, dominance and epistatic components (Falconer and Mackay, 1996). This partitioning has been widely analyzed in plant genetics and breeding using robust statistical methods based on linear mixed model theory. Some studies have shown very little involvement of dominance variance (Gallais, 1990), whereas others have found that this factor can contribute substantially, for example, in annual plants (Wardyn *et al.*, 2007) or tree species (Bouvet *et al.*, 2009). In addition, combining dominance and additive effects in models can improve genotype prediction, for example, in genomic selection (Denis and Bouvet, 2013). Epistatic effects—although a potential source of bias in estimating additive and dominance effects—are often overlooked in experimental quantitative genetic studies (Lynch and Walsh, 1998). Research is now underway to determine its importance with regard to complex traits (Mäki-Tanila and Hill, 2014). Different strategies have been implemented through quantitative genetics models, quantitative trait loci (QTL) or mutation-QTL experiments (Mackay, 2014), but there is still no consensus on its importance. For plants, depending on the population and trait considered, some studies have shown that epistasis can significantly contribute to total genetic variance, for example, in maize (Dudley and Johnson, 2009), rice (Luo *et al.*, 2009)

or cotton (Li *et al.*, 2014), while also improving the prediction accuracy (Dudley and Johnson, 2009; Hu *et al.*, 2011). Meanwhile, other studies have shown that epistasis has a little impact regarding the genetic architecture of some traits in maize (Buckler *et al.*, 2009) or in the prediction accuracy (Lorenzana and Bernardo, 2009).

Recently, with the advent of high-throughput molecular technology, numerous markers distributed throughout the whole genome have been used to develop new models of genetic variation (Legarra *et al.*, 2009). Additive and non-additive genetic relationship matrices can, for instance, be constructed from genome-wide single-nucleotide polymorphism (SNP) markers to disentangle confounding factors for the estimation of genetic variance, leading to more accurate models than those based on pedigree (Lee *et al.*, 2010; Su *et al.*, 2012; Muñoz *et al.*, 2014).

This new approach in estimating relationships among individuals has markedly broadened the prospects for modeling genetic effects. This is especially the case regarding animal breeding, when purebred animals are evaluated for their crossbred performance, where models that take the allele origin into account are implemented. Some rely on the relationship among purebred parents (Lo *et al.*, 1997), which are referred to as ‘parent models’ in the following, while others are based on the relationship among crossbred progenies estimated by their

<sup>1</sup>CIRAD, Genetic Improvement and Adaptation of Mediterranean and Tropical Plants (AGAP) Research Unit, Montpellier, France and <sup>2</sup>Centre de Recherche sur la Durabilité et la Productivité des Plantations Industrielles, Pointe-Noire, Republic of the Congo

Correspondence: Dr J-M Bouvet, CIRAD-BIOS, UMR: AGAP, Campus of Lavalette TA A-108/01, Cedex 5, Montpellier 34398, France.

E-mail: jean-marc.bouvet@cirad.fr

Received 15 April 2015; revised 3 July 2015; accepted 22 July 2015; published online 2 September 2015

gametotype (haplotype of each parental gamete) (Ibáñez-Escriche *et al.*, 2009; Kinghorn *et al.*, 2010; Zeng *et al.*, 2013), which are referred to as 'progeny models' hereafter. Although both parent and progeny models have been fully studied in assessing the performance of crossbreeds, they have mainly been focused on genetic additive effects. However, it could be important to take the magnitude of non-additive effects into consideration in some breeding contexts.

In plants, significant progress has been achieved in estimating additive and non-additive variance within a single population, but modeling the genetic variation requires further investigations with hybrids resulting from the cross of two different species or from the cross of two different populations of the same species presenting different allelic frequencies. The gene origin should be taken into account when estimating variance in hybrid progenies, as shown by Stuber and Cockerham (1966). Massmann *et al.* (2013) and Technow *et al.* (2012), with dense genotyping in maize, implemented such an approach to estimate additive and dominance effects, but further investigations are needed to address pending questions in hybrid populations. Does a marker-based relationship matrix perform better than one that is pedigree-based for estimating non-additive variance? Are variance estimation and prediction of genetic value improved by using parent or progeny models?

Beyond the advantages of using appropriate models, hybrids derived from a pure species cross or from two different populations of a single species are often used in crop and perennial species breeding to capture the genetic gain resulting from heterosis (Gallais, 2009). Estimating additive and non-additive variance components is crucial to understand adaptive trait expression and guide breeding programs for long rotation species.

In that setting, the objectives of this study were as follows: (i) to test the performance of pedigree- or marker-based models, (ii) to analyze the influence of parent- and progeny-based relationship matrices in estimating additive and non-additive effects (dominance, epistasis) in hybrid populations, and (iii) to assess the impact of including non-additive effects on the model prediction accuracy.

Data of a first-generation hybrid population, resulting from an *Eucalyptus urophylla* × *Eucalyptus grandis* cross achieved under field conditions and genotyped with SNP, were used to model the environmental and genetic effects, estimate the additive and non-additive variance and analyse the prediction efficacies of the different models.

## MATERIALS AND METHODS

### Field experimental data

Thirteen *E. urophylla* females and 9 *E. grandis* males were crossed by controlled pollination to generate 69 full-sib families and 1415 progenies. The males and females were selected in progeny/provenance trials set up in the Republic of the Congo. Each parent tree was selected in a different family. Previous studies have shown that panmixy is the mating pattern in continuous natural *E. urophylla* populations (Tripiana *et al.*, 2007) with a preferential out-crossing system for both species (Horsley and Johnson, 2007). We thus considered that the coefficient of inbreeding and the coefficients of relationship among parents within each species were null. Each of the 1415 progenies was replicated three times using cutting and a clonally replicated progeny test was planted with 1415 clones (4415 trees at a stocking density of 833 trees ha<sup>-1</sup>). The field experiment involved a complete block design with three replications. Around 25 trees replicated in three blocks represented each full-sib family. At 32 months, the total number of trees used in this study was reduced to 3596, representing 1130 clones, owing to natural mortality and elimination of non-genotyped trees. The trial site was located in the Republic of the Congo, east of Pointe-Noire (11°59'21" E, 4°45'51" S). Rainfall averaged 1200 mm year<sup>-1</sup>. The soils were

characterized by low water retention and a very low level of organic matter as well as a poor cationic exchange capacity.

The total tree height (HT) at 32 months was used for model comparison. It averaged 12.02 m and had a high among-individual coefficient of variation (30%).

### Molecular data

Among the 20 000 SNP identified, 3303 were finally selected based on the repeatability, a minor allele frequency >2.5% and a rate of missing data per marker <5%. The 1152 individuals (13 females, 9 males and 1130 progenies) were genotyped with the 3303 SNPs using the DArTseq technology (Wenzl *et al.*, 2004). Haplotype phasing and missing-data inference were done with the localized haplotype clustering method developed in Beagle version 4.0 (Browning and Browning, 2007). Pedigree information among the 1152 individuals, that is, the parent-offspring, half-sib and full-sib relationships, was used by the Beagle algorithm to estimate haplotypes and missing data (Browning and Browning, 2013).

### Genetic model

The models developed in the next section are derived from the Stuber and Cockerham (1966) model. The authors considered two genetically divergent populations (species or populations of the same species but with different allele frequencies) and the hybrid population generated by the crosses of random individuals from the two parent populations. By considering the hybrid population with the gene effect depending on the parent origin, that is, female parents from population F (*E. urophylla* in our experiment) and male parents from population M (*E. grandis* in our experiment), assuming linkage equilibrium in each parent population and limiting epistatic effects to first-order epistasis, the total genetic variance was partitioned as follows:

$$\sigma_G^2 = \sigma_{aF}^2 + \sigma_{aM}^2 + \sigma_{dFM}^2 + \sigma_{aFaF}^2 + \sigma_{aMaM}^2 + \sigma_{aFaM}^2 + \sigma_{aFdFM}^2 + \sigma_{aMdFM}^2 + \sigma_{dFMdFM}^2 \text{ [hybrid variance model]}$$

where  $\sigma_{aF}^2$  and  $\sigma_{aM}^2$  are the female and male additive variances of the hybrid population due to alleles from females (males) crossed with males (females);  $\sigma_{dFM}^2$  is the dominance variance due to the female × male cross;  $\sigma_{aFaF}^2$  and  $\sigma_{aMaM}^2$  are the female (male) additive × additive epistatic variances;  $\sigma_{aFdFM}^2$  and  $\sigma_{aMdFM}^2$  are the female (male) epistatic additive × dominance variances and  $\sigma_{dFMdFM}^2$  is the dominance × dominance epistatic variance involving alleles from the M and F parent populations.

The genetic covariance among two individual hybrids  $x$  and  $y$ :

$$\text{cov}(x, y)_G = \varphi_f \sigma_{aF}^2 + \varphi_m \sigma_{aM}^2 + \varphi_f \varphi_m \sigma_{dFM}^2 + \varphi_f^2 \sigma_{aFaF}^2 + \varphi_m^2 \sigma_{aMaM}^2 + \varphi_f \varphi_m \sigma_{aFaM}^2 + \varphi_f^2 \varphi_m \sigma_{aFdFM}^2 + \varphi_f \varphi_m^2 \sigma_{aMdFM}^2 + \varphi_f^2 \varphi_m^2 \sigma_{dFMdFM}^2 \text{ [hybrid covariance model]}$$

where  $\varphi_f$  and  $\varphi_m$  denote the coancestry coefficient among females of population F and males of population M, respectively; for example, for half-sib hybrids with the same female from F and different males from M:  $\varphi_f = (1 + \Psi_f/2)$  and  $\varphi_m = 0$  and for full sib hybrids, with the same female from F and the same male from M:  $\varphi_f = (1 + \Psi_f/2)$  and  $\varphi_m = (1 + \Psi_m/2)$ .  $\Psi_f$  and  $\Psi_m$  are the coefficients of inbreeding of the parent in population F and M that are considered null in this study.

### Statistical model

We developed two types of model based on the equation by Stuber and Cockerham (1966). The first, called the parent model, was based on the relationship among female ( $f=13$ ) and male ( $m=9$ ) parents of the hybrid progenies:

$$y = X\beta + Z_{\text{col}}\text{col} + Z_{r:b}r : b + Z_p \text{plot} + Z_f a_f + Z_m a_m + Z_d d + Z_i i + \varepsilon \text{ [parent - model]}$$

where  $y$  is the vector of the phenotypic variable (tree height at 32 months),  $\beta$  is the vector of fixed effects due to the general mean and blocks,  $\text{col} \sim N(0, \sigma_{\text{col}}^2 \mathbf{I}_d)$  is the vector of random spatial environmental effects due to the field design column, with  $\sigma_{\text{col}}^2$  being the variance related to the spatial effects,  $\mathbf{I}_d$  is the identity matrix,  $r : b \sim N(0, \sigma_{r:b}^2 \mathbf{I}_d)$  is the vector of random spatial

environmental effects due to field design row by block interaction, with  $\sigma_{r;b}^2$  being the variance related to the spatial effects,  $\mathbf{plot} \sim N(\mathbf{0}, \sigma_{\text{plot}}^2 \mathbf{Id})$  is a vector of random spatial environmental effects common to the individuals planted in the same square plot, with  $\sigma_{\text{plot}}^2$  being the variance related to the spatial effects,  $\mathbf{a}_f \sim N(\mathbf{0}, \sigma_{\text{af}}^2 \mathbf{A}_f)$  is a vector of random additive effects due to *E. urophylla* females, with  $\mathbf{A}_f[f,f]$  being the coancestry coefficient matrix among females  $\{f_i\}$  estimated from the pedigree ( $\mathbf{A}_{fP}$ ) or marker ( $\mathbf{A}_{fG}$ ), with  $\sigma_{\text{af}}^2$  being the additive variance of the hybrid population due to alleles from females crossed with males,  $\mathbf{a}_m \sim N(\mathbf{0}, \sigma_{\text{am}}^2 \mathbf{A}_m)$  is the vector of random additive effects due to *E. grandis* males, with  $\mathbf{A}_m[m,m]$  being the coancestry coefficient matrix among males  $\{m_i\}$  estimated from the pedigree ( $\mathbf{A}_{mP}$ ) or marker ( $\mathbf{A}_{mG}$ ), with  $\sigma_{\text{am}}^2$  being the additive variance of the hybrid population due to alleles from males crossed with females,  $\mathbf{d} \sim N(\mathbf{0}, \sigma_{\text{d}}^2 \mathbf{D})$  is a vector of random dominance effects due to the female  $\times$  male cross,  $\mathbf{D}[f \times m, f \times m]$  is estimated from the pedigree ( $\mathbf{D}_P$ ) or marker ( $\mathbf{D}_G$ ), with  $\sigma_{\text{d}}^2$  being the dominance variance of the hybrid population due to alleles from males crossed with females (see the estimation in the next section),  $\mathbf{i} \sim N(\mathbf{0}, \sigma_{\text{i}}^2 \mathbf{E})$  is the term representing random epistatic effects modeled by either the sum of three additive  $\times$  additive effects  $\mathbf{aa}_f \sim N(\mathbf{0}, \sigma_{\text{aaf}}^2 \mathbf{E}_{\text{AAf}})$ ,  $\mathbf{aa}_m \sim N(\mathbf{0}, \sigma_{\text{aam}}^2 \mathbf{E}_{\text{AAm}})$  and  $\mathbf{a}_f \mathbf{a}_m \sim N(\mathbf{0}, \sigma_{\text{afam}}^2 \mathbf{E}_{\text{AfAm}})$  or the sum of additive  $\times$  dominance effects  $\mathbf{a}_f \mathbf{d} \sim N(\mathbf{0}, \sigma_{\text{afd}}^2 \mathbf{E}_{\text{AfD}})$ ,  $\mathbf{a}_m \mathbf{d} \sim N(\mathbf{0}, \sigma_{\text{amd}}^2 \mathbf{E}_{\text{AmD}})$  or by the dominance  $\times$  dominance effects  $\mathbf{dd} \sim N(\mathbf{0}, \sigma_{\text{dd}}^2 \mathbf{E}_{\text{DD}})$  and  $\mathbf{e} \sim N(\mathbf{0}, \sigma_{\text{e}}^2 \mathbf{Id})$  is the vector of residual effects. The epistatic matrices  $\mathbf{E}_{\text{AAf}}[f,f]$ ,  $\mathbf{E}_{\text{AAm}}[m,m]$ ,  $\mathbf{E}_{\text{AfAm}}[f \times m, f \times m]$ ,  $\mathbf{E}_{\text{AfD}}[f \times m, f \times m]$ ,  $\mathbf{E}_{\text{AmD}}[f \times m, f \times m]$  and  $\mathbf{E}_{\text{DD}}[f \times m, f \times m]$  were calculated with the Hadamard and Kronecker products (formulas detailed in the next section).

$\mathbf{X}$ ,  $\mathbf{Z}_{\text{col}}$ ,  $\mathbf{Z}_{r;b}$ ,  $\mathbf{Z}_p$ ,  $\mathbf{Z}_m$ ,  $\mathbf{Z}_f$ ,  $\mathbf{Z}_{mf}$  and  $\mathbf{Z}_i$  are the incidence matrices connecting the fixed and random effects to the data. The coancestry matrices  $\mathbf{A}_{fG}$ ,  $\mathbf{A}_{mG}$  and  $\mathbf{D}_G$  were estimated using the formulas defined in the next section. Based on the generic model, 10 parent models combining additive, dominance and epistatic effects and marker or pedigree coancestry matrices were developed (Table 1). Additive  $\times$  additive, additive  $\times$  dominance and dominance  $\times$  dominance epistatic effects were estimated in separate models because the restricted maximum likelihood (REML) algorithm failed to converge when all effects were included in a single model.

The second type of models, called the progeny model, was developed using relationships among the ( $c = 1130$ ) hybrid progenies. They differed from the parent models because they were defined using the female- and male-origin haplotypes of each progeny inferred by long-range phasing. Such models were implemented in previous studies (Ibáñez-Escriche *et al.*, 2009; Kinghorn *et al.*, 2010; Zeng *et al.*, 2013) while only considering additive effects in the model. We extended this approach by including dominance and epistatic effects.

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_{\text{col}}\mathbf{col} + \mathbf{Z}_{r;b}\mathbf{r};\mathbf{b} + \mathbf{Z}_p\mathbf{plot} + \mathbf{Z}_c\mathbf{a}_f + \mathbf{Z}_c\mathbf{a}_m + \mathbf{Z}_c\mathbf{d} + \mathbf{Z}_c\mathbf{i} + \boldsymbol{\varepsilon} \text{ [progeny - model]}$$

The fixed effects and environmental random effects were defined as for the first model. The genetic effects are defined by:  $\mathbf{a}_f \sim N(\mathbf{0}, \sigma_{\text{af}}^2 \mathbf{A}_{fG}^H)$ ,  $\mathbf{a}_m \sim N(\mathbf{0}, \sigma_{\text{am}}^2 \mathbf{A}_{mG}^H)$ ,  $\mathbf{d} \sim N(\mathbf{0}, \sigma_{\text{d}}^2 \mathbf{D}_{fG}^H)$  and  $\mathbf{i} \sim N(\mathbf{0}, \sigma_{\text{i}}^2 \mathbf{E}_{fG}^H)$ , with  $\mathbf{A}_{fG}^H[c,c]$  and  $\mathbf{A}_{mG}^H[c,c]$  the molecular-based female and male additive relationship matrices, respectively, and  $\mathbf{D}_{fG}^H[c,c]$  the dominance relationship matrix. The  $\mathbf{i} \sim N(\mathbf{0}, \sigma_{\text{i}}^2 \mathbf{E}^H)$  is the term representing random epistatic effects modeled by either the sum of three additive  $\times$  additive effects  $\mathbf{aa}_f \sim N(\mathbf{0}, \sigma_{\text{aaf}}^2 \mathbf{E}_{\text{AAf}}^H)$ ,  $\mathbf{aa}_m \sim N(\mathbf{0}, \sigma_{\text{aam}}^2 \mathbf{E}_{\text{AAm}}^H)$  and  $\mathbf{a}_f \mathbf{a}_m \sim N(\mathbf{0}, \sigma_{\text{afam}}^2 \mathbf{E}_{\text{AfAm}}^H)$ , or the sum of additive  $\times$  dominance effects  $\mathbf{a}_f \mathbf{d} \sim N(\mathbf{0}, \sigma_{\text{afd}}^2 \mathbf{E}_{\text{AfD}}^H)$ ,  $\mathbf{a}_m \mathbf{d} \sim N(\mathbf{0}, \sigma_{\text{amd}}^2 \mathbf{E}_{\text{AmD}}^H)$ , or by the dominance  $\times$  dominance effects  $\mathbf{dd} \sim N(\mathbf{0}, \sigma_{\text{dd}}^2 \mathbf{E}_{\text{DD}}^H)$ . The  $\mathbf{e} \sim N(\mathbf{0}, \sigma_{\text{e}}^2 \mathbf{Id})$  is the vector of residual effects. The epistatic matrices all had the same dimensions  $[c,c]$  and were estimated using the male and female haplotype of each progeny with the Hadamard products (formulas detailed in the next section). Five progeny models combining additive, dominance and epistatic effects were developed (Table 1).

The best linear unbiased predictors (BLUP related to the genetic effects) were computed by solving the mixed model equations. The variance component estimation based on the REML method and the BLUP calculations were done using the ASReml version 3 program (Gilmour *et al.*, 2006) implemented in R software (R Development Core Team, 2011).

Narrow- and broad-sense heritabilities were defined by  $h^2 = \sigma_{\text{a}}^2 / \sigma_{\text{p}}^2$  and  $H^2 = \sigma_{\text{g}}^2 / \sigma_{\text{p}}^2$ , and dominance and epistatic variance ratios by  $d^2 = \sigma_{\text{d}}^2 / \sigma_{\text{p}}^2$ ,  $i^2 = \sigma_{\text{i}}^2 / \sigma_{\text{p}}^2$ , with  $\sigma_{\text{a}}^2$ ,  $\sigma_{\text{d}}^2$ ,  $\sigma_{\text{g}}^2$  and  $\sigma_{\text{p}}^2$  being the total additive, epistatic, genetic

and phenotypic variances, respectively, that varied according to the model. Details on variance components are given in Table 1. The delta method was used to obtain a first-order approximate s.e. for a nonlinear function of a vector of random variables with known or estimated covariance matrix (Oehlert, 1992).

### Pedigree and marker relationship matrices

For the parent models, the pedigree coancestry coefficients (or kinship)  $\varphi_f$  and  $\varphi_m$  were estimated based on the pedigree of the female and male parent population. As the pedigree was unknown within the *E. urophylla* and *E. grandis* parents, the  $\mathbf{A}_{fP}$  and  $\mathbf{A}_{mP}$  matrices were simply the identity matrix multiplied by 0.5 (equal to the self coancestry of a non-inbred individual),  $\mathbf{D}_P$  was the identity matrix with 0.25, epistatic  $\mathbf{E}_{\text{AAf}}$ ,  $\mathbf{E}_{\text{AAm}}$ ,  $\mathbf{E}_{\text{AfAm}}$  were diagonals with 0.25,  $\mathbf{E}_{\text{AfD}}$ ,  $\mathbf{E}_{\text{AmD}}$  were diagonals with 0.125 and  $\mathbf{E}_{\text{DD}}$  was a diagonal with 0.0625. The molecular marker-based coancestry between two individuals, defined as the probability that two alleles at the locus taken from each individual are equal, that is, identical by state, was defined using Van Raden's estimator (Van Raden, 2008). The female  $\mathbf{A}_{fG}[f,f]$  and male  $\mathbf{A}_{mG}[m,m]$  coancestry matrices were calculated as follows:

$$\mathbf{A}_{fG} = 1/2 \frac{(\mathbf{M}_f - \mathbf{P}_f)(\mathbf{M}_f - \mathbf{P}_f')}{2 \sum_{i=1}^s p_{fi}(1 - p_{fi})}$$

$$\mathbf{A}_{mG} = 1/2 \frac{(\mathbf{M}_m - \mathbf{P}_m)(\mathbf{M}_m - \mathbf{P}_m')}{2 \sum_{i=1}^s p_{mi}(1 - p_{mi})}$$

where  $\mathbf{M}_f[f,s]$  ( $\mathbf{M}_m[m,s]$ ) is the matrix of female (male) marker information (coded as 0, 1 and 2 for, respectively, AA, Aa and aa) and  $s$  the total number of SNPs (3303 SNPs).  $\mathbf{P}_f[f,s]$  ( $\mathbf{P}_m[m,s]$ ) contains frequencies  $p_{fi}$  ( $p_{mi}$ ) of the second allele at each locus such that column  $i$  of  $\mathbf{P}_f$  ( $\mathbf{P}_m$ ) is  $2p_{fi}$  ( $2p_{mi}$ ).

The dominance relationship matrices  $\mathbf{D}_P$  and  $\mathbf{D}_G$  were defined as follows. For two individual hybrids X and Y resulting from the cross of female  $f_i$  and male  $m_i$  and female  $f_j$  and male  $m_j$ , respectively, the dominance relationship coefficient was defined as the product of coancestry coefficients between females and males  $D_{X,Y} = \varphi_{f_i,f_j} \times \varphi_{m_i,m_j}$ .  $\mathbf{D}$  can be calculated by the Kronecker product  $\mathbf{A}_f \otimes \mathbf{A}_m$  after removing lines and columns corresponding to the absent crosses. In the parent models, according to the Stuber and Cockerham (1966) models, the relationship matrices due to the first-degree epistatic terms were computed using the Hadamard and Kronecker products. They were defined as follows: for additive  $\times$  additive terms  $\mathbf{E}_{\text{AAf}} = \mathbf{A}_f \# \mathbf{A}_f$ ,  $\mathbf{E}_{\text{AAm}} = \mathbf{A}_m \# \mathbf{A}_m$ ,  $\mathbf{E}_{\text{AfAm}} = \mathbf{A}_f \otimes \mathbf{A}_m$ , for the dominance  $\times$  dominance term  $\mathbf{E}_{\text{DD}} = \mathbf{D} \# \mathbf{D}$  and for the additive  $\times$  dominance terms  $\mathbf{E}_{\text{AfD}} = \mathbf{A}_f \# \mathbf{A}_f \otimes \mathbf{A}_m$  and  $\mathbf{E}_{\text{AmD}} = \mathbf{A}_m \# \mathbf{A}_m \otimes \mathbf{A}_f$ .

For the progeny models, the female and male additive molecular marker-based coancestry matrices  $\mathbf{A}_{fG}^H[c,c]$  and  $\mathbf{A}_{mG}^H[c,c]$  were derived from the haplotypes of each progeny using Van Raden's estimator (Van Raden, 2008):

$$\mathbf{A}_{fG}^H = 1/2 \frac{(\mathbf{H}_f - \mathbf{P}_f^H)(\mathbf{H}_f - \mathbf{P}_f^H)'}{\sum_{i=1}^s p_{fi}(1 - p_{fi})}$$

$$\mathbf{A}_{mG}^H = 1/2 \frac{(\mathbf{H}_m - \mathbf{P}_m^H)(\mathbf{H}_m - \mathbf{P}_m^H)'}{\sum_{i=1}^s p_{mi}(1 - p_{mi})}$$

where  $\mathbf{H}_f[c,s]$  and  $\mathbf{H}_m[c,s]$  are the matrices of female and male haplotype marker information (coded as 0 and 1 for A and a, respectively), respectively, and  $c$  being the number of clones and  $s$  the total number of SNPs (3303 SNPs).  $\mathbf{P}_f^H$  ( $\mathbf{P}_m^H$ ) contained frequencies  $p_{fi}$  ( $p_{mi}$ ) of the allele at each locus such that column  $i$  of  $\mathbf{P}_f^H$  ( $\mathbf{P}_m^H$ ) is  $1p_{fi}$  ( $1p_{mi}$ ).

The dominance relationship matrix was defined using the Hadamard product:  $\mathbf{D}_{fG}^H = \mathbf{A}_{fG}^H \# \mathbf{A}_{fG}^H$ . The relationship matrices due to the first-degree epistatic terms were computed using the Hadamard product, for additive  $\times$  additive terms  $\mathbf{E}_{\text{AAf}}^H = \mathbf{A}_{fG}^H \# \mathbf{A}_{fG}^H$ ,  $\mathbf{E}_{\text{AAm}}^H = \mathbf{A}_{mG}^H \# \mathbf{A}_{mG}^H$ ,  $\mathbf{E}_{\text{AfAm}}^H = \mathbf{A}_{fG}^H \# \mathbf{A}_{mG}^H$ , for the dominance  $\times$  dominance term  $\mathbf{E}_{\text{DD}}^H = \mathbf{D}_{fG}^H \# \mathbf{D}_{fG}^H$  and for the additive  $\times$  dominance terms  $\mathbf{E}_{\text{AfD}}^H = \mathbf{A}_{fG}^H \# (\mathbf{A}_{fG}^H \# \mathbf{A}_{mG}^H)$ ,  $\mathbf{E}_{\text{AmD}}^H = \mathbf{A}_{mG}^H \# (\mathbf{A}_{fG}^H \# \mathbf{A}_{mG}^H)$ .

The correspondence between models and the relationship matrices are given in Table 1.

**Table 1** Main characteristics of the models and the associated relationship matrix

Model type	Code	Relationship matrices related to the model: A additive, D dominance			Variance components	
		Additive	Dominance	Epistasis		
Parent	P_A	$A_{fP}, A_{mP}$			$\sigma_a^2 = \sigma_{af}^2 + \sigma_{am}^2$ $\sigma_g^2 = \sigma_{af}^2 + \sigma_{am}^2$ $\sigma_p^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_{col}^2 + \sigma_{r;b}^2 + \sigma_{plot}^2 + \sigma_e^2/3$ Idem P_A	
	G_A	$A_{fG}, A_{mG}$			$\sigma_a^2 = \sigma_{af}^2 + \sigma_{am}^2$ $\sigma_d^2$ $\sigma_g^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2$ $\sigma_p^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{col}^2 + \sigma_{r;b}^2 + \sigma_{plot}^2 + \sigma_e^2/3$ Idem P_A+D	
	P_A+D	$A_{fP}, A_{mP}$	$D_P$		$\sigma_a^2 = \sigma_{af}^2 + \sigma_{am}^2$ $\sigma_d^2$ $\sigma_i^2 = \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2$ $\sigma_g^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2$ $\sigma_p^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2 + \sigma_{col}^2 + \sigma_{r;b}^2 + \sigma_{plot}^2 + \sigma_e^2/3$ Idem P_A+D+AA	
	G_A+D	$A_{fG}, A_{mG}$	$D_G$		$\sigma_a^2 = \sigma_{af}^2 + \sigma_{am}^2$ $\sigma_d^2$ $\sigma_i^2 = \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2$ $\sigma_g^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2$ $\sigma_p^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2 + \sigma_{col}^2 + \sigma_{r;b}^2 + \sigma_{plot}^2 + \sigma_e^2/3$ Idem P_A+D+AD	
	P_A+D+AA	$A_{fP}, A_{mP}$	$D_P$	$A_{fP}\#A_{fP}, A_{mP}\#A_{mP}$	$\sigma_a^2 = \sigma_{af}^2 + \sigma_{am}^2$ $\sigma_d^2$ $\sigma_i^2 = \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2$ $\sigma_g^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2$ $\sigma_p^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2 + \sigma_{col}^2 + \sigma_{r;b}^2 + \sigma_{plot}^2 + \sigma_e^2/3$ Idem P_A+D+AD	
	G_A+D+AA	$A_{fG}, A_{mG}$	$D_G$	$A_{fG}\#A_{fG}, A_{mG}\#A_{mG}$	$\sigma_a^2 = \sigma_{af}^2 + \sigma_{am}^2$ $\sigma_d^2$ $\sigma_i^2 = \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2$ $\sigma_g^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2$ $\sigma_p^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2 + \sigma_{col}^2 + \sigma_{r;b}^2 + \sigma_{plot}^2 + \sigma_e^2/3$ Idem P_A+D+AD	
	P_A+D+AD	$A_{fP}, A_{mP}$	$D_P$	$A_{fP}\#D_P, A_{mP}\#D_P$	$\sigma_a^2 = \sigma_{af}^2 + \sigma_{am}^2$ $\sigma_d^2$ $\sigma_i^2 = \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2$ $\sigma_g^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2$ $\sigma_p^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2 + \sigma_{col}^2 + \sigma_{r;b}^2 + \sigma_{plot}^2 + \sigma_e^2/3$ Idem P_A+D+AD	
	G_A+D+AD	$A_{fG}, A_{mG}$	$D_G$	$A_{fG}\#D_G, A_{mG}\#D_G$	$\sigma_a^2 = \sigma_{af}^2 + \sigma_{am}^2$ $\sigma_d^2$ $\sigma_i^2 = \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2$ $\sigma_g^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2$ $\sigma_p^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{aaf}^2 + \sigma_{aam}^2 + \sigma_{afam}^2 + \sigma_{col}^2 + \sigma_{r;b}^2 + \sigma_{plot}^2 + \sigma_e^2/3$ Idem P_A+D+AD	
	P_A+D+DD	$A_{fP}, A_{mP}$	$D_P$	$D_P\#D_P$	$\sigma_a^2 = \sigma_{af}^2 + \sigma_{am}^2$ $\sigma_d^2$ $\sigma_i^2 = \sigma_{ddf}^2$ $\sigma_g^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{ddf}^2$ $\sigma_p^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{ddf}^2 + \sigma_{col}^2 + \sigma_{r;b}^2 + \sigma_{plot}^2 + \sigma_e^2/3$ Idem P_A+D+DD	
	G_A+D+DD	$A_{fG}, A_{mG}$	$D_G$	$D_G\#D_G$	$\sigma_a^2 = \sigma_{af}^2 + \sigma_{am}^2$ $\sigma_d^2$ $\sigma_i^2 = \sigma_{ddf}^2$ $\sigma_g^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{ddf}^2$ $\sigma_p^2 = \sigma_{af}^2 + \sigma_{am}^2 + \sigma_d^2 + \sigma_{ddf}^2 + \sigma_{col}^2 + \sigma_{r;b}^2 + \sigma_{plot}^2 + \sigma_e^2/3$ Idem P_A+D+DD	
	Progeny	G_A <sup>H</sup>	$A_{fG}^H, A_{mG}^H$			Idem P_A
		G_A <sup>H</sup> +D <sup>H</sup>	$A_{fG}^H, A_{mG}^H$	$D_{fG}^H$		Idem P_A+D
G_A <sup>H</sup> +D <sup>H</sup> +A <sup>H</sup> A <sup>H</sup>		$A_{fG}^H, A_{mG}^H$	$D_{fG}^H$	$A_{fG}^H\#A_{mG}^H, A_{mG}^H\#A_{fG}^H$	Idem P_A+D+AA	
G_A <sup>H</sup> +D <sup>H</sup> +A <sup>H</sup> D <sup>H</sup>		$A_{fG}^H, A_{mG}^H$	$D_{fG}^H$	$A_{fG}^H\#D_{fG}^H, A_{mG}^H\#D_{fG}^H$	Idem P_A+D+AD	
G_A <sup>H</sup> +D <sup>H</sup> +D <sup>H</sup> D <sup>H</sup>		$A_{fG}^H, A_{mG}^H$	$D_{fG}^H$	$D_{fG}^H\#D_{fG}^H, D_{mG}^H\#D_{fG}^H$	Idem P_A+D+DD	

Abbreviations: F, female; G, marker; H, haplotype; M, male.  $\sigma_a^2, \sigma_d^2, \sigma_i^2, \sigma_g^2$  and  $\sigma_p^2$  being the additive, dominance, epistatic, total genetic and phenotypic variances, respectively.

### Model comparison

Models were compared using the Akaike information criterion (AIC; Akaike, 1974) defined as  $AIC = -2\ln(R) + 2t$ , where  $\ln(R)$  is the log-likelihood of the model and  $t$  is the number of variance parameters. The model with the lowest AIC value presented the best data fit.

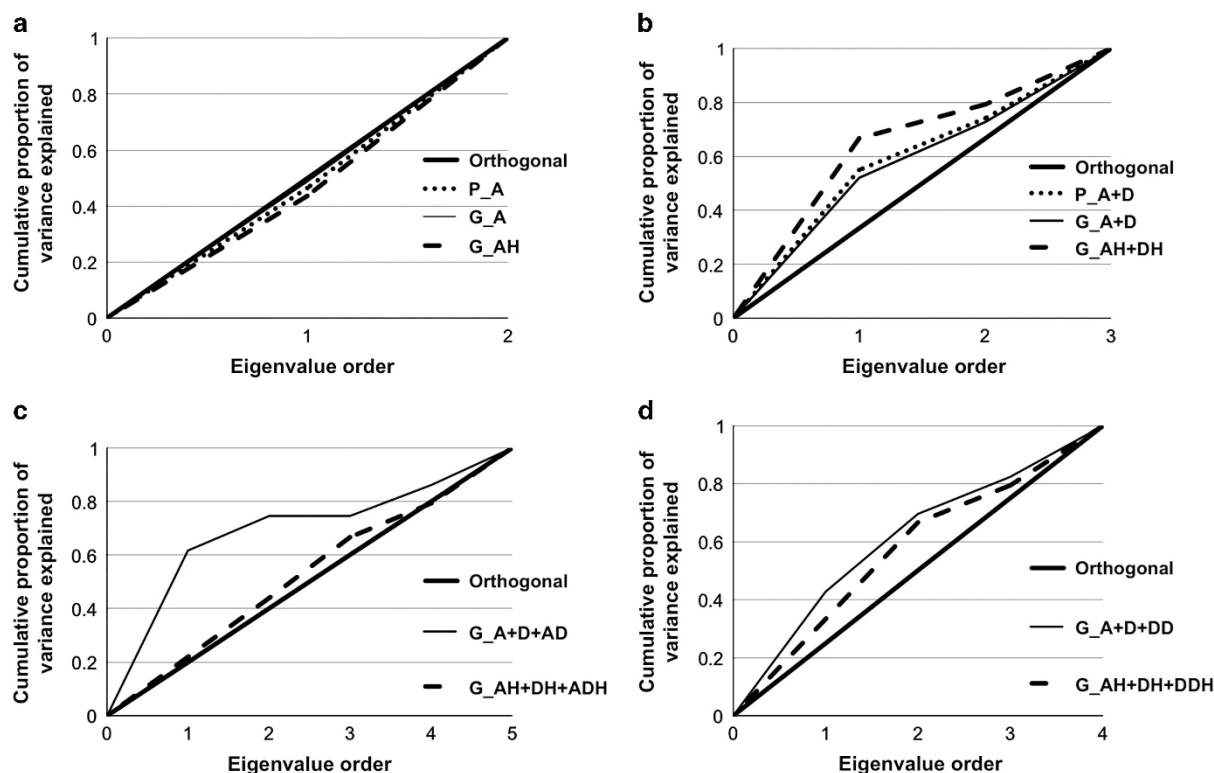
To assess the dependency among variance components, we used the procedure described in Muñoz *et al.* (2014): with COVES being the asymptotic variance-covariance matrix of estimates of variance parameters, the asymptotic sampling correlation matrix of estimates (CORES) was computed as  $CORES = (VARES^{1/2} COVES VARES^{-1/2})$ , where VARES is a diagonal matrix containing the diagonal elements of V (variances of estimates). Analyzing the magnitude of the correlations of the CORES matrix allows assessment of the dependency between the pairwise estimates. To achieve an overall assessment of dependency between the estimates, Muñoz *et al.* (2014) suggested assessing the percentage of variation explained by eigenvalues of the CORES matrix. This was carried out by plotting the cumulative percentage of variance explained by the model vs the eigenvalue order (Figure 1). The eigenvalues were computed using the R statistical software function 'eigen', which computes eigenvalues and eigenvectors of matrices (R Development Core Team, 2011).

Goodness-of-fit was evaluated with the full data set (1130 clones) by assessing the correlation between predicted additive genetic values and phenotypes of individual trees (average of three replicates of each individual clone)  $r(\hat{A}_{full}, \hat{Y}_{full})$  and between predicted total genetic values and phenotypes  $r(\hat{G}_{full}, \hat{Y}_{full})$ .

The prediction accuracy was tested using the pedigree-based relationship matrix for parent models, or the so-called P-BLUP method, and the marker relationship matrix for parent and progeny models, that is, the G-BLUP method (Cros *et al.*, 2015). A cross-validation procedure was implemented to evaluate the prediction accuracy, with the data set divided into nine subsets. The half-sib progeny phenotypes related to one female and the half-sib progeny phenotypes related to one male were removed at each set. As the smallest number of parents corresponded to the nine males, nine subsets were possible. This process allows comparison of the prediction ability of the parent and progeny models. We chose this method to make sure that the progenies in the training and candidate sets would not have any common parents, which would have biased the comparison of parent and progeny models regarding their predictive capacity. Eight of the sets, including trees with phenotype and genotype, were used in turn for training the models to estimate the model parameters, and the remaining set (candidates set), made of trees with only genotypes, was used for testing the prediction accuracy. This procedure led to different sizes of the training and candidate populations, which ranged from 794 to 1013 clones and from 117 to 336 clones, respectively. The prediction accuracy of the models was evaluated on the candidate set by the correlation between the predicted additive and total genetic values using the training set and phenotypes of individual tree (average of three replicates for each individual clone):  $r(\hat{A}_{can}, \hat{Y}_{can})$  and  $r(\hat{G}_{can}, \hat{Y}_{can})$ .

We also used the prediction stability defined by Muñoz *et al.* (2014), which is the correlation between breeding or total genetic values of clones of the





**Figure 1** Proportion of variance explained by eigenvalues of asymptotic correlation matrices of variance estimates by considering: additive effects from the pedigree parent model (P\_A), vs the marker parent model (G\_A), vs the progeny model (G\_A<sup>H</sup>) (a); additive and dominance effects from the pedigree parent model (P\_A+D), vs the marker parent model (G\_A+D), vs the progeny model (G\_A<sup>H</sup>+D<sup>H</sup>) (b); additive, dominance and additive × dominance effects from the marker parent model (G\_A+D+AD), vs the progeny model (G\_A<sup>H</sup>+D<sup>H</sup>+AD<sup>H</sup>) (c); additive, dominance and dominance × dominance effects from the marker parent model (G\_A+D+D#D), vs the progeny model (G\_A<sup>H</sup>+D<sup>H</sup>+DD<sup>H</sup>) (d). The diagonal represents an independent correlation matrix between variance estimates.

candidate population predicted with the full data set and the breeding or total genetic values predicted with the training population:  $r(\hat{A}_{\text{can}}, \hat{A}_{\text{full}})$  and  $r(\hat{G}_{\text{can}}, \hat{G}_{\text{full}})$ . These correlations measured the dependency of the predicted additive or total genetic values on the phenotypes. Associated mean square errors  $\text{MSE}(\hat{A}_{\text{can}}, \hat{A}_{\text{full}})$  and  $\text{MSE}(\hat{G}_{\text{can}}, \hat{G}_{\text{full}})$  were calculated.

An analysis of variance (ANOVA) was performed to test the null hypothesis, that is, an absence of difference among models and subsets. ANOVA was conducted with prediction accuracy, prediction stability and MSE as dependent variables. When the ANOVA detected a significant model effect, a pairwise mean comparison among models was carried out for each variable using Newman–Keuls test (Shaffer, 2007).

## RESULTS

### Variance components

The variance component estimates for the 15 models are presented in Table 2. When comparing the 10 parent models, small-to-moderate changes were observed in the genetic, spatial environmental and residual variances, except for the plot variance  $\sigma_{\text{plot}}^2$ , which markedly decreased by the inclusion of non-additive effects. Epistatic variances could not be estimated with the three pedigree-based parent models, that is, P\_A+D+AA, P\_A+D+AD and P\_A+D+DD, because matrix singularities were detected and the REML algorithm did not converge. This is partly explained by some redundancy in the dominance and epistasis matrices (as described in the pedigree and marker relationship matrices section). Otherwise, epistatic components were estimated with the three equivalent marker-based models, but there were some particularities in the estimates.  $\sigma_{\text{aaf}}^2$  and  $\sigma_{\text{aam}}^2$  were null with G\_A+D+AA, while  $\sigma_{\text{afd}}^2$  and  $\sigma_{\text{amd}}^2$  were quite high (1.022 and 4.863,

respectively) with G\_A+D+AD, but the dominance variance  $\sigma_{\text{d}}^2$  became null. Finally, with G\_A+D+DD, the  $\sigma_{\text{dd}}^2$  epistatic variance was estimated without affecting the other variances, but it had a very high s.e., as also the  $\sigma_{\text{d}}^2$  estimates (see Supplementary Table S1 in Supplementary Material).

With the five progeny models, the spatial environmental and residual variances did not change when including non-additive effects, but the additive variances decreased (Table 2). According to the ASReml version 3 (Gilmour *et al.*, 2006) algorithm output, models that included epistatic effects converged at the boundary for the  $\sigma_{\text{aaf}}^2$ ,  $\sigma_{\text{aam}}^2$ ,  $\sigma_{\text{afd}}^2$ ,  $\sigma_{\text{amd}}^2$  and  $\sigma_{\text{dd}}^2$  epistatic variances. The components were thus considered as null.

The comparison between the parent and progeny models showed a minor change in the additive variance estimate but a marked change for dominance variance, for example, for P\_A+D and G\_A+D, additive and dominance were close ( $\sigma_{\text{af}}^2=1.239$ ,  $\sigma_{\text{am}}^2=0.921$ ,  $\sigma_{\text{d}}^2=2.640$ ) and ( $\sigma_{\text{af}}^2=1.260$ ,  $\sigma_{\text{am}}^2=0.953$ ,  $\sigma_{\text{d}}^2=2.395$ ), respectively, whereas an increase in  $\sigma_{\text{d}}^2$  was noted for G\_A<sup>H</sup>+D<sup>H</sup> ( $\sigma_{\text{af}}^2=1.340$ ,  $\sigma_{\text{am}}^2=0.819$ ,  $\sigma_{\text{d}}^2=4.313$ ). A change was also observed with residual variance, which decreased from parent models ( $\sigma_{\text{e}}^2=9.772$  in average) to progeny models ( $\sigma_{\text{e}}^2=8.029$  on average).

The overall degree of dependency between the model variance estimates is clearly illustrated in Figure 1, which shows, for each model, the cumulative proportion of variance explained by eigenvalues compared with the diagonal representing orthogonal correlation matrix, that is, independence between estimates. Figure 1a shows that the cumulative percentage of genetic variance of the P\_A and G\_A

models was closer to the diagonal than  $G_{A^H}$ , highlighting that these former models led to less correlated variance estimates. A similar trend was observed with  $P_{A+D}$  and  $G_{A+D}$  compared with  $G_{A^H+D^H}$  (Figure 1b). On the other hand, with epistatic models, Figures 1c and d show that the progeny models generated less correlated variance estimates than the parent models. An examination of the asymptotic correlation matrix of variance estimates led to similar conclusions (Supplementary Tables S2A–C in Supplementary Materials).

The variance ratios differed between models (Table 3). The narrow sense heritability ( $h^2$ ) estimated with all of the additive effect models ( $P_A$ ,  $G_A$  and  $G_{A^H}$ ) decreased after including the dominance effect. Adding epistatic effects led to a decrease in both additive and dominance variance when the epistatic estimates differed from zero. As expected, the narrow-sense heritability,  $h^2$  ( $h^2 = 0.136$  in average), was lower than the broad-sense heritability  $H^2$  with a marked difference ( $H^2 = 0.334$  on average). This observation stresses the preponderance of non-additive effects, while the  $(\sigma_d^2 + \sigma_i^2) / \sigma_a^2$  ratio was 122, 108, 226 and 293%, with  $P_{A+D}$ ,  $G_{A+D}$ ,  $G_{A+D+DD}$  and  $G_{A+D+AD}$  parent models, respectively. The same trend was

observed with the progeny model, where the  $(\sigma_d^2 + \sigma_i^2) / \sigma_a^2$  ratio was 199% with  $G_{A^H+D^H}$ .

### Model goodness-of-fit and prediction accuracy

The relative model quality was assessed through the AIC. The results showed that the progeny models, with AIC ranging from 8841 to 8847, better fitted the observations than the parent models, with AIC ranging from 8889 to 8955 (Table 2).

The goodness-of-fit was estimated using the full data set (Table 4). With the parent models, the correlations between additive or total genetic values and phenotypes were low, with  $r(\hat{A}_{full}, \hat{Y}_{full})$  varying in the (0.361; 0.368) interval and  $r(\hat{G}_{full}, \hat{Y}_{full})$  in the (0.478; 0.490) interval, while both did not show marked differences between models. Regarding the progeny models, the correlations were higher, showing a better fit to the observations, with  $r(\hat{A}_{full}, \hat{Y}_{full})$  varying in the (0.710; 0.770) interval and  $r(\hat{G}_{full}, \hat{Y}_{full})$  had the same value for all models (0.858) due to the null epistatic variance estimates.

As expected, the prediction accuracy had much lower correlations than goodness-of-fit, with  $r(\hat{A}_{can}, \hat{Y}_{can})$  varying in the (0.191; 0.288) interval and  $r(\hat{G}_{can}, \hat{Y}_{can})$  in the (0.191; 0.294) interval (Table 4).

**Table 2 Environmental and genetic variance components and parameters related to the quality of the different models**

Model type	Code	Source of variation													AIC	Observation
		$\sigma_{r,b}^2$	$\sigma_{col}^2$	$\sigma_{af}^2$	$\sigma_{am}^2$	$\sigma_d^2$	$\sigma_{dd}^2$	$\sigma_{aaf}^2$	$\sigma_{aam}^2$	$\sigma_{afam}^2$	$\sigma_{afd}^2$	$\sigma_{amd}^2$	$\sigma_{plot}^2$	$\sigma_e^2$		
Parent	P_A	0.620	0.029	1.393	1.206								2.605	9.784	8894.230	
	G_A	0.620	0.030	1.165	1.104								2.570	9.784	8890.146	
	P_A+D	0.621	0.027	1.239	0.921	2.640							2.177	9.772	8892.548	
	G_A+D	0.620	0.028	1.260	0.953	2.395							2.191	9.772	8889.140	
	P_A+D+AA	—	—	—	—	—	—	—	—	—	—	—	—	—	8955.555	Singularities
	G_A+D+AA	0.620	0.028	1.260	0.953	2.395	0.000	0.000					2.191	9.772	8893.140	
	P_A+D+DD	—	—	—	—	—	—	—	—	—	—	—	—	—	8952.068	Singularities
	G_A+D+DD	0.619	0.028	1.224	0.899	1.855	2.973						2.181	9.773	8891.045	
	P_A+D+AD	—	—	—	—	—	—	—	—	—	—	—	—	—	8950.993	Singularities
Progeny	G_A+D+AD	0.620	0.028	1.085	0.922	0.000					1.022	4.863	2.181	9.774	8892.846	
	G_A <sup>H</sup>	0.648	0.021	2.143	1.659								2.564	8.390	8847.049	
	G_A <sup>H</sup> +D <sup>H</sup>	0.648	0.024	1.340	0.819	4.313							2.575	8.029	8841.162	
	G_A <sup>H</sup> +D <sup>H</sup> +A <sup>H</sup> A <sup>H</sup>	0.648	0.024	1.340	0.819	4.313		0.000	0.000				2.575	8.029	8846.090	
	G_A <sup>H</sup> +D <sup>H</sup> +D <sup>H</sup> D <sup>H</sup>	0.648	0.024	1.340	0.819	4.313	0.000						2.575	8.029	8843.162	
G_A <sup>H</sup> +D <sup>H</sup> +A <sup>H</sup> D <sup>H</sup>	0.648	0.024	1.340	0.819	4.313					0.000	0.000	2.575	8.029	8844.766		

**Table 3 Heritabilities and variance ratios and their s.e.**

Model	Code	$h^2$	s.e.( $h^2$ )	$d^2$	s.e.( $d^2$ )	$i^2$	s.e.( $i^2$ )	$H^2$	s.e.( $H^2$ )
Parent	P_A	0.166	0.074						
	G_A	0.149	0.068						
	P_A+D	0.124	0.068	0.152	0.087			0.276	0.091
	G_A+D	0.129	0.063	0.139	0.086			0.268	0.095
	P_A+D+AA	—	—	—	—	—	—	—	—
	G_A+D+AA	0.129	0.063	0.139	0.086	0.000	0.000	0.268	0.095
	P_A+D+DD	—	—	—	—	—	—	—	—
	G_A+D+DD	0.109	0.078	0.095	0.156	0.152	0.457	0.356	0.278
	P_A+D+AD	—	—	—	—	—	—	—	—
Progeny	G_A+D+AD	0.098	0.057	0.000	0.000	0.287	0.150	0.385	0.135
	G_A <sup>H</sup>	0.246	0.065						
	G_A <sup>H</sup> +D <sup>H</sup>	0.122	0.051	0.243	0.096			0.365	0.086
	G_A <sup>H</sup> +D <sup>H</sup> +A <sup>H</sup> A <sup>H</sup>	0.122	0.051	0.243	0.096	0.000	0.000	0.365	0.086
	G_A <sup>H</sup> +D <sup>H</sup> +D <sup>H</sup> D <sup>H</sup>	0.122	0.051	0.243	0.096	0.000	0.000	0.365	0.086
G_A <sup>H</sup> +D <sup>H</sup> +A <sup>H</sup> D <sup>H</sup>	0.122	0.051	0.243	0.096	0.000	0.000	0.365	0.086	

**Table 4 Comparison of models for goodness-of-fit, prediction accuracy ability and stability**

Model type	Code	Goodness-of-fit (full data set)		Prediction accuracy (cross-validation procedure)		Prediction stability (cross-validation procedure)			
		$r(\hat{A}_{full}, \hat{Y}_{full})$	$r(\hat{G}_{full}, \hat{Y}_{full})$	$r(\hat{A}_{can}, \hat{Y}_{can})^a$ P = 0.002 <sup>b</sup>	$r(\hat{G}_{can}, \hat{Y}_{can})$ P = 0.000	$r(\hat{A}_{can}, \hat{A}_{full})$ P = 0.484	$r(\hat{G}_{can}, \hat{G}_{full})$ P = 0.551	MSE( $\hat{A}_{can}, \hat{A}_{full}$ ) <sup>a</sup> P = 0.014	MSE( $\hat{G}_{can}, \hat{G}_{full}$ ) P = 0.918
Parent	P_A	0.367		0.198a <sup>c</sup>		0.713a		0.478ab	
	G_A	0.368		0.226ab		0.783a		0.385ab	
	P_A+D	0.361	0.490	0.191a	0.191a	0.720a	0.523a	0.372ab	0.809a
	G_A+D	0.364	0.478	0.225ab	0.202a	0.794a	0.598a	0.353ab	0.711a
	P_A+D+AA	—	—	—	—	—	—	—	—
	G_A+D+AA	0.364	0.478	0.203a	0.199a	0.710a	0.594a	0.404ab	0.711a
	P_A+D+DD	—	—	—	—	—	—	—	—
	G_A+D+DD	0.363	0.480	0.222ab	0.205a	0.795a	0.600a	0.327ab	0.719a
	P_A+D+AD	—	—	—	—	—	—	—	—
G_A+D+AD	0.366	0.481	0.215a	0.200a	0.779a	0.597a	0.315ab	0.722a	
Progeny	G_A <sup>H</sup>	0.770		0.266ab		0.665a		0.574b	
	G_A <sup>H</sup> +D <sup>H</sup>	0.710	0.858	0.288b	0.294b	0.752a	0.548a	0.236a	0.754a
	G_A <sup>H</sup> +D <sup>H</sup> +A <sup>H</sup> A <sup>H</sup>	0.710	0.858	0.288b	0.294b	0.752a	0.548a	0.236a	0.754a
	G_A <sup>H</sup> +D <sup>H</sup> +D <sup>H</sup> D <sup>H</sup>	0.710	0.858	0.288b	0.294b	0.752a	0.548a	0.236a	0.754a
	G_A <sup>H</sup> +D <sup>H</sup> +A <sup>H</sup> D <sup>H</sup>	0.710	0.858	0.288b	0.294b	0.752a	0.548a	0.236a	0.754a
	G_A <sup>H</sup> +D <sup>H</sup> +A <sup>H</sup> D <sup>H</sup>	0.710	0.858	0.288b	0.294b	0.752a	0.548a	0.236a	0.754a
	G_A <sup>H</sup> +D <sup>H</sup> +A <sup>H</sup> D <sup>H</sup>	0.710	0.858	0.288b	0.294b	0.752a	0.548a	0.236a	0.754a

<sup>a</sup>Correlations and mean square errors are calculated as the average over nine subsets of the cross-validation procedure.

<sup>b</sup>Probability associated with Fisher test of the analysis of variance comparing the models.

<sup>c</sup>Letters a and b correspond to the groups defined by Newman–Keuls test for mean comparison. The alpha risk is 5%.

With parent models,  $r(\hat{A}_{can}, \hat{Y}_{can})$  had higher estimates with marker than with pedigree-based relationships, for example,  $r(\hat{A}_{can}, \hat{Y}_{can}) = 0.198$  with P\_A and  $r(\hat{A}_{can}, \hat{Y}_{can}) = 0.226$  with G\_A (Table 4). Regarding the progeny models, the predictive ability estimates were higher than the parent models, that is, on average  $r(\hat{A}_{can}, \hat{Y}_{can}) = 0.277$ . However, although ANOVA showed a global significant model effect ( $P = 0.002$ ), only two overlapping groups were detected by the Newman–Keuls test (Table 4). The predictive ability for total genetic value showed higher estimates for the progeny models, that is,  $r(\hat{G}_{can}, \hat{Y}_{can}) = 0.294$  on average, than for the parent models, that is,  $r(\hat{G}_{can}, \hat{Y}_{can}) = 0.199$ . ANOVA ( $P = 0.000$ ) and the mean comparison test (Table 4) clearly differentiated the two groups. As shown by the boxplot graphs (Figures 2a and b), all the models presented a marked variation in the prediction accuracy between the nine subsets of the cross-validation process, and no clear differences were detected among models. Parent model using a marker-based pedigree relationship matrix (P\_A and P\_A+D) even showed negative values.

The prediction stability, which measured the dependency of the additive or total genetic values on the phenotypes, showed better performance for the parent models, with  $r(\hat{A}_{can}, \hat{A}_{full}) = 0.756$  and  $r(\hat{G}_{can}, \hat{G}_{full}) = 0.582$  on average, than for the progeny models, with  $r(\hat{A}_{can}, \hat{A}_{full}) = 0.734$  and  $r(\hat{G}_{can}, \hat{G}_{full}) = 0.548$  (Table 4). However, ANOVA did not reveal significant differences among models ( $P = 0.484$  and  $P = 0.551$ ).

The MSE of prediction of additive values presented, on average for all the models, lower estimates than obtained for the total genetic values, that is,  $MSE(\hat{A}_{can}, \hat{A}_{full}) = 0.346$  and  $MSE(\hat{G}_{can}, \hat{G}_{full}) = 0.734$  (Table 4). For the additive value prediction, although the progeny model G\_A<sup>H</sup>+D<sup>H</sup> had a lower MSE, Fisher test was barely significant, while two overlapping groups were detected by Newman–Keuls test (Table 4). No significant differences between MSE were observed for the prediction of total genetic values. MSE distributions clearly

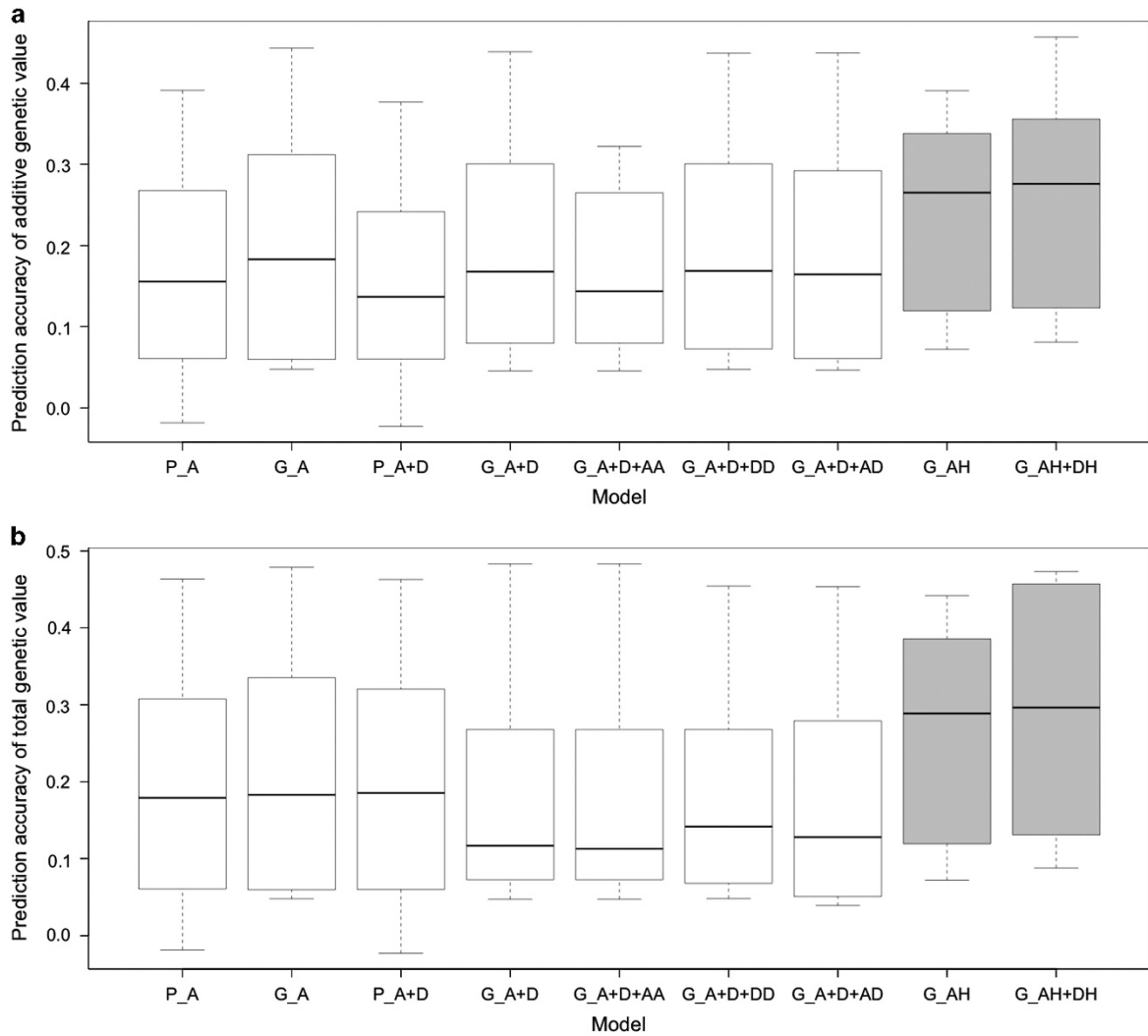
separated the parent from the progeny models, with the latter exhibiting a much narrower distribution (Figures 3a and b).

## DISCUSSION

In this study, we used a first-generation hybrid population resulting from the cross of two parent species, to analyze the ability of the models to separate additive and non-additive effects and to predict genetic values. We compared two sets of models, the first based on pedigree and marker-based parent relationship matrices, with the second using hybrid progeny haplotypes to define marker-based relationship matrices. The model of Stuber and Cockerham (1966) has been implemented with parent pedigree or marker-based relationship matrices in previous studies (for example, Lorenzana and Bernardo, 2009; Massman *et al.*, 2013; Cros *et al.*, 2015), whereas one goal of our study was to develop a new set of models using progeny haplotypes to better estimate non-additive and particularly epistatic effects.

### Modeling genetic variance in hybrid populations

Our approach relied on the variance partition of Stuber and Cockerham (1966) while considering two parental populations and complemented the recent study of Muñoz *et al.* (2014) in which variance was modeled when considering a single population. According to Stuber and Cockerham (1966), distinguishing the parental origin of alleles into models should theoretically lead to variance components different from considering alleles originating from a single population; in the first case we have two allele frequency distributions at each locus and in the second case a unique distribution. We verified this assumption by estimating additive and dominance variances while considering a single population in our design (Appendix Table A1, in annex). Variance components, estimated with a single population approach, were lower for the pedigree-based



**Figure 2** Comparison of the distribution of prediction accuracy for the breeding value ( $r(\hat{A}_{can}, \hat{Y}_{can})$ ) (a), total genetic value ( $r(\hat{G}_{can}, \hat{Y}_{can})$ ) (b), resulting from the cross-validation procedure among the pedigree-based parent models (P\_A, P\_A+D), marker-based parent models (G\_A, G\_A+D, G\_A+D+AA, G\_A+D+DD, G\_A+D+AD) and progeny models (G\_A<sup>H</sup>, G\_A<sup>H</sup>+D<sup>H</sup>).

relationship P\_A+D model ( $\sigma_a^2=1.959$  and  $\sigma_d^2=1.097$ ) and much lower with the marker-based relationship G\_A+D model ( $\sigma_a^2=0.08$  and  $\sigma_d^2=0.01$ ) as compared with the estimates presented in Table 2.

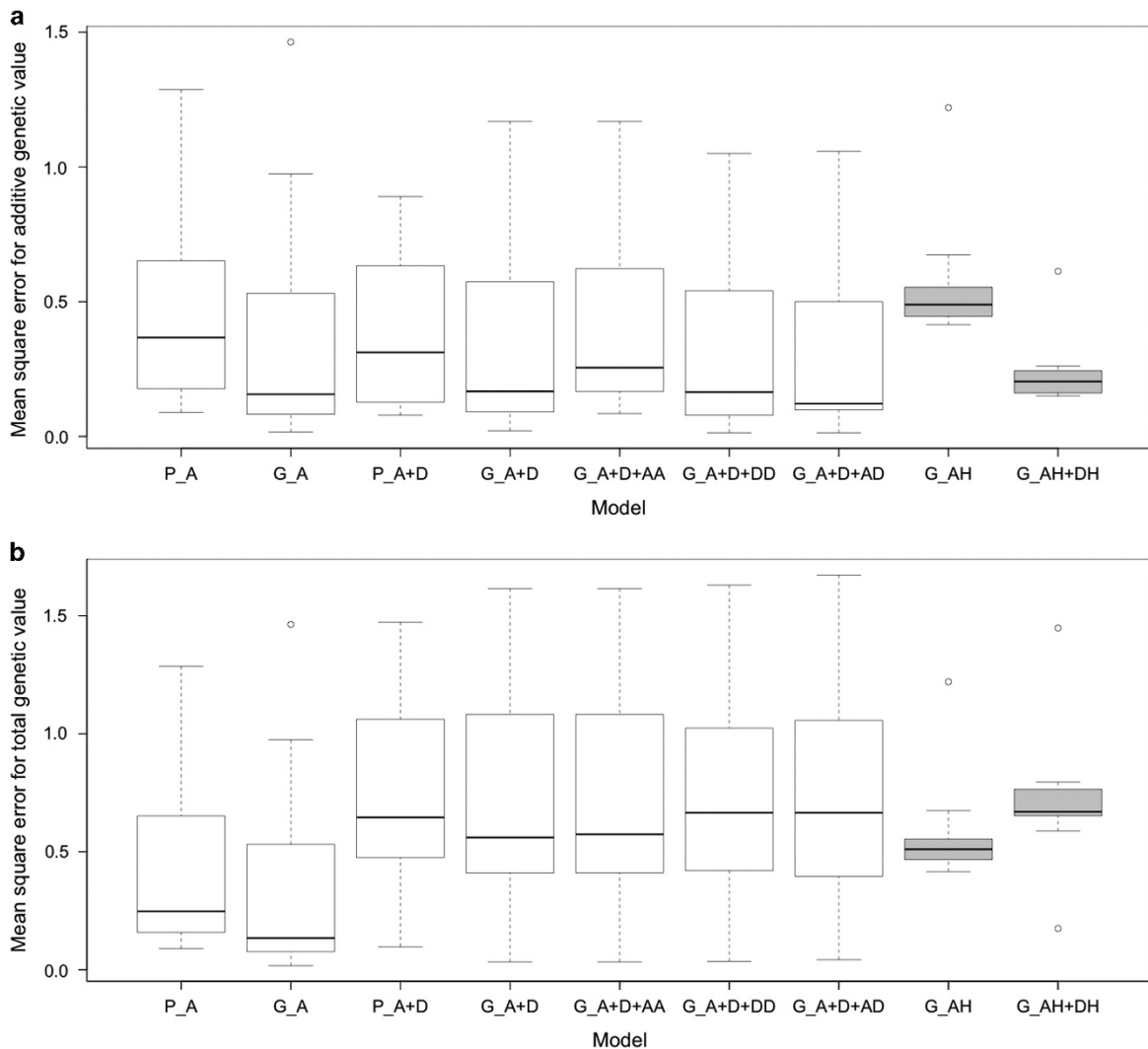
With the dual population approach, a first set of nested models based on the parent relationship was implemented. Variance components were affected by the use of marker-based ( $A_{fG}$ ,  $A_{mG}$ ,  $D_G$ ) or pedigree-based relationship matrices ( $A_{fP}$ ,  $A_{mP}$ ,  $D_P$ ). The additive variance decreased from P\_A to G\_A (16% for  $\sigma_{af}^2$  and 8% for  $\sigma_{am}^2$ ) but increased from P\_A+D to G\_A+D (1.6% for  $\sigma_{af}^2$  and 3.4% for  $\sigma_{am}^2$ ), while the dominance variance decreased (9%) (Table 2). However, these variations were much smaller than the additive and dominance variance s.es. (Supplementary Table S1 in Supplementary Material), reflecting the moderate impact of using marker- or pedigree-based relationship matrices. This result could first be explained by the ability of  $A_{fG}$ ,  $A_{mG}$  and  $D_G$  to reflect the assumption of an absence of genetic relationship between *E. urophylla* parents and between *E. grandis* parents. This assumption was suggested by the fact that parents are progenies of trees selected in natural populations and the mating pattern in *Eucalyptus* is instead panmixy. Second, the set of markers used for establishing  $A_{fG}$ ,  $A_{mG}$  and  $D_G$ , resulting from strong selection

(3303 SNPs were selected among 20000), was of sufficiently good quality to correctly reflect the relationship. Our findings were consistent with those obtained in previous studies comparing the use of pedigree- and marker-based relationship matrices in modeling (Forni *et al.*, 2011; Wang *et al.*, 2014).

Using the dual population approach, we developed the so-called progeny models based on the female and male haplotypes of each hybrid progeny. The relationship estimation was based on the classic formula of Van Raden (2008), that is, using SNPs of each haplotype independently to estimate the coancestry coefficient. It differs from approaches constructing haplotype segments surrounding putative QTLs to estimate the relationship coefficient (Makgahlela *et al.*, 2013). Our progeny models estimated variances with smaller errors as compared with parent models (Table 3). However, they did not perform better in estimating the non-additive effects (epistatic variance estimates were all null). As already mentioned, this result could be attributed to an overparametrization of the model.

Our results showed that including non-additive effects in parent and progeny models decreased the additive variance (Table 2) and consequently the narrow-sense heritability (Table 3). This trend was





**Figure 3** Comparison of the distribution of MSEs for breeding values ( $MSE(\hat{A}_{can}, \hat{A}_{full})$ ) (a) and total genetic values ( $MSE(\hat{G}_{can}, \hat{G}_{full})$ ) (b) resulting from the cross-validation procedure among the pedigree-based parent models (P\_A, P\_A+D), marker-based parent models (G\_A, G\_A+D, G\_A+D+AA, G\_A+D+DD, G\_A+D+AD) and progeny models (G\_A<sup>H</sup>, G\_A<sup>H</sup>+D<sup>H</sup>).

expected from a theoretical standpoint (Gallais, 1990; Falconer and Mackay, 1996) as additive variance in the additive-only model actually captures non-additive variation, as confirmed experimentally in other studies (Su *et al.*, 2012; Muñoz *et al.*, 2014; Sun *et al.*, 2014).

#### Estimating non-additive effects

One of the goals of our study was to evaluate the extent of epistatic variance in a hybrid population. We used the Hadamard product to define the relationship coefficients of epistatic variance–covariance matrices. This way of calculating the coefficients supposes that the population is in linkage equilibrium and the markers are not linked (Schnell, 1963). To assess the potential bias of our approach, we estimated the linkage disequilibrium (LD) using the ‘genetics’ R package (R Development Core Team, 2011) and the LD estimators ‘ $r^2$ ’ from Hill and Robertson (1968). The mean and s.d of  $r^2$  were  $r^2 = 0.03$  (0.11) for *E. urophylla* and  $r^2 = 0.05$  (0.12) for *E. grandis*. These values showed very low LD, indicating that Hadamard matrix multiplication should not substantially bias epistatic variance estimates. Regarding the linkage between markers, Schnell (1963) showed that the recombination rate should be taken into account in the

epistatic coefficient matrices. Omitting this parameter upward biased the epistatic variance components. For *Eucalyptus*, the recombination rate is not well known and there are wide variations in estimates (Silva-Junior and Grattapaglia, 2015), which makes bias evaluation difficult.

Our analyses of the data used for the model did not lead to conclusive results. With parent models, non-null epistatic variances were estimated with G\_A+D+DD and G\_A+D+AD; however, the s.e.s. of  $\sigma_{dd}^2$ ,  $\sigma_{afd}^2$  and  $\sigma_{amd}^2$  were very substantial, that is, twofold the variance estimates (Supplementary Table S1 in Supplementary Material). With the progeny models, we estimated null epistatic variances and a very high dominance variance. Null epistatic variance may result from a small contribution of this effect to the genetic variance of growth traits, as observed in other forest tree studies (Lepointevin *et al.*, 2011; Araújo *et al.*, 2012) and/or because we implemented an overfitted model, which could lead to null variance estimates. Increasing the number of parents and reducing the imbalance in the mating design could be a way to better estimate non-null epistatic variances and reduce their s.e.

As shown in previous forest tree studies (Araújo *et al.*, 2012), with an experimental design involving cloned full-sibs and pedigree-based models, modeling cannot separate actual additive and dominance variance from epistatic variance. This design upwardly biased the estimation of true additive and dominance variances and downwardly biased epistatic variance. Lee *et al.* (2010) and Muñoz *et al.* (2014) have shown that the use of a marker-based relationship can disentangle non-additive effects from additive and dominance effects. Our models, using a marker-based relationship matrix, had this property (Figures 1a–d), but they were not able to clearly estimate the epistatic variance.

In addition, with the Stuber and Cockerham (1966) model, we could not separate the pure dominance effect from the female-additive  $\times$  male-additive effect because the relationship matrices  $\{\varphi_m, \varphi_f\}$  used to estimate  $\sigma_d^2$  and  $\sigma_{afam}^2$  were identical. Hence, the dominance variance was inflated and the additive  $\times$  additive variance was downwardly biased.

Although epistatic variance cannot be clearly estimated, our results have shown that non-additive variance explained a significant part of the total genetic variance, that is, the same extent of additive variance with parent models and twofold higher additive variance with progeny models (Table 3). In the case of height at 32 months, our results stressed the importance of non-additive effects and confirmed the findings of previous studies of Baltunis *et al.* (2009) with radiata pine and Araújo *et al.* (2012) with *E. globulus* dealing with growth traits. However, they differed from earlier forest tree studies using full-sib families with clonal replicates and analyzing growth traits, which revealed a very low or limited proportion of non-additive variance, for example, Lepoittevin *et al.* (2011) with maritime pine. Our results could be explained by the inter-species hybrid nature of our material, which can exacerbate heterosis (Melchinger *et al.*, 2007). Some examples have been reported in poplar (Wu, 2000), rice (Zhou *et al.*, 2012) and maize hybrids (Guo *et al.*, 2014).

### Impact of modeling on the prediction accuracy

We tested the prediction accuracy of the models using the G-BLUP approach, which is among the most widely used parametric methods in genomic selection and giving good results (Heslot *et al.*, 2015). Our results revealed that the prediction accuracy estimated in the candidate population (on average  $r(\hat{A}_{can}, \hat{Y}_{can}) = 0.242$  and  $r(\hat{G}_{can}, \hat{Y}_{can}) = 0.242$ ) remained low compared with the goodness-of-fit (on average  $r(\hat{A}_{full}, \hat{Y}_{full}) = 0.513$  and  $r(\hat{G}_{full}, \hat{Y}_{full}) = 0.649$ ) and prediction stability (on average  $r(\hat{A}_{can}, \hat{A}_{full}) = 0.747$  and  $r(\hat{G}_{can}, \hat{G}_{full}) = 0.567$ ) using the full data set. Although comparisons with other studies should be carried out with caution because the prediction accuracies are calculated according to different methods and formulas, the prediction stability values assessed in our study were similar to those estimated in previous forest tree studies for growth traits (Resende *et al.*, 2012; Muñoz *et al.*, 2014; Beaulieu *et al.*, 2014a, b; Gamal El-Dien *et al.*, 2015). The prediction accuracy defined in our study by  $r(\hat{A}_{can}, \hat{Y}_{can})$  and  $r(\hat{G}_{can}, \hat{Y}_{can})$  is supposed to give a more objective genomic selection potential. With our experimental breeding data, these parameters were quite low (on average  $r(\hat{A}_{can}, \hat{Y}_{can}) = 0.242$  and  $r(\hat{G}_{can}, \hat{Y}_{can}) = 0.242$ ), which could be explained by factors influencing the genomic selection performance, such as: heritability, training population size, and LD. In our case, narrow- and broad-sense heritabilities of height at 32 months presented low-to-moderate estimates (Table 3), which partially explained the low genomic selection accuracy. LD, which was very low, that is,  $r^2 = 0.03$  and  $r^2 = 0.05$  for female and male populations, respectively, was likely also a key factor explaining our accuracy. Finally, the training population

size, ranging from 794 to 1013 clones to predict 117 to 336 clones, was probably not sufficiently large, as this factor is critical, especially when non-additive effects are included in the model (Denis and Bouvet, 2013).

Our model comparison for prediction ability showed some unexpected results. Surprisingly, for parent models,  $r(\hat{G}_{can}, \hat{Y}_{can})$  were smaller than  $r(\hat{A}_{can}, \hat{Y}_{can})$ , whereas we expected a higher value because the total genetic value was more related to the phenotype than to the breeding value. This result could be attributed to the less accurate prediction of non-additive effects owing to overparameterization of the model.

Our study has shown that the prediction accuracy and stability were improved by using marker-based instead of pedigree-based relationship matrices (Table 4). The marker-based matrix had the advantage of capturing both the Mendelian segregation within full-sib families and genetic links through unknown common ancestors, which are not available in the known pedigree. This feature has been observed in genomic selection for animal (Chen *et al.*, 2011) or plant breeding (Zapata-Valenzuela *et al.*, 2013; Muñoz *et al.*, 2014; Beaulieu *et al.*, 2014a; Cros *et al.*, 2015). The comparison of parent and progeny models revealed the higher prediction accuracy of the latter, which better assessed the relationship among progenies by capturing the Mendelian segregation. Note that the performance of this model is improved by reducing the error related to haplotype reconstruction, which can be done by using well-known and multi-generation pedigrees when imputing the data (Hayes *et al.*, 2012).

We tested how the inclusion of non-additive genetic effects affected the prediction accuracy and stability. In our parent model, adding dominance or epistatic effects did not improve the prediction, that is,  $r(\hat{A}_{can}, \hat{Y}_{can})$  or  $r(\hat{G}_{can}, \hat{Y}_{can})$  even decreased from  $G_A$  to  $G_{A+D}$  and  $G_{A+D+AA}$ . However, with progeny models, the prediction accuracy regarding the breeding value increased from  $G_{A^H}$  to  $G_{A^H+D^H}$ . Our study generated new findings to supplement previously published results based on experimental data that had diverse conclusions, for example, in plants, Dudley and Johnson (2009), Hu *et al.* (2011) and Wang *et al.* (2012) found that including marker interactions substantially increased the prediction accuracy, but Lorenzana and Bernardo (2009) and Muñoz *et al.* (2014) came to opposite conclusions. With animals, Su *et al.* (2012) and Sun *et al.* (2014) showed that including non-additive effects improved the prediction. According to Denis and Bouvet (2013), with simulated data, the inclusion of non-additive effects improved the prediction accuracy when non-additive effects were preponderant, the training data set size was high and updated over selection cycles to reassess the relationship between markers and QTLs.

### Implications for breeding

In the case of single populations, marker-based estimates of coancestry have proved to be a useful alternative to genealogical information for estimating variance components as well as for predicting genetic values (Muñoz *et al.*, 2014). In hybrid populations, we have shown that using genome-wide information can improve the variance partition, although epistasis detection requires further investigation to effectively evaluate its relative magnitude. Using progeny models with a higher number of parents and a more balanced design could improve their performance. The modeling approach incorporating molecular data could reduce the overestimation of additive variance and improve the estimation of non-additive effects and is crucial for accurately predicting genetic gain in hybrid breeding programs, as previously mentioned, for example, in maize (Parvez *et al.*, 2007), oil palm (Cros *et al.*, 2015) or *Eucalyptus* (Bouvet *et al.*, 2009). Regarding the

prediction accuracy, our study has shown that the inclusion of non-additive effects in genomic selection depends on the modeling approach—as compared with parent models, our progeny models, using parental haplotypes, improved the prediction accuracy in the hybrid population.

## DATA ARCHIVING

Data are available from the Dryad Digital Repository: <http://dx.doi.org/10.5061/dryad.g73t2>

## CONFLICT OF INTEREST

The authors declare no conflict of interest.

## ACKNOWLEDGEMENTS

This study was conducted within the framework of the *Eucalyptus* breeding program carried out by CRPDI in Congo and was part of the WUETREE project funded by the Bioadapt Programme of the French Research National Agency (ANR), while also being funded by the Congolese Government and CIRAD. Field experiments were conducted at CRDPI in the Congo, and molecular analyses were performed under the technical supervision of Alexandre Vaillant in the CIRAD laboratories in Montpellier, France. We thank our colleagues in the Congo for their valuable help in sampling, without whom this work would not have been possible. We are grateful to the three reviewers for their helpful comments that significantly improved the manuscript.

## AUTHOR CONTRIBUTIONS

Field experiment idea and design, implementation and sampling: PhV and GM; study and analysis design: JMB; supervision of writing: JMB; and contribution of ideas, comments, analyses and revision of manuscript versions: PhV, DC and GM.

Akaike H (1974). A new look at the statistical model identification. *Trans Autom Control* **19**: 716–723.

Araújo JA, Borralho NMG, Dehon G (2012). The importance and type of non-additive genetic effects for growth in *Eucalyptus globulus*. *Tree Genet Genomes* **8**: 327–337.

Baltunis BS, Wu HX, Dungey HS, Mullin TJ, Brawner JT (2009). Comparisons of genetic parameters and clonal value predictions from clonal trials and seedling base population trials of radiata pine. *Tree Genet Genomes* **5**: 269–278.

Beaulieu J, Doerksen T, Clément S, MacKay J, Bousquet J (2014a). Accuracy of genomic selection models in a large population of open-pollinated families in white spruce. *Heredity* **113**: 343–352.

Beaulieu J, Doerksen T, MacKay JK, Rainvilleand A, Bousquet J (2014b). Genomic selection accuracies within and between environments and small breeding groups in white spruce. *BMC Genomics* **15**: 1048.

Bouvet J-M, Saya A, Vigneron Ph (2009). Trends in additive, dominance and environmental effects with age for growth traits in *Eucalyptus* hybrid populations. *Euphytica* **165**: 35–54.

Browning BL, Browning SR (2013). Improving the accuracy and efficiency of identity by descent detection in population data. *Genetics* **194**: 459–471.

Browning SR, Browning BL (2007). 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.

Buckler ES, Holland JB, Bradbury PJ, Acharya CB, Brown PJ, Browne C *et al.* (2009). The genetic architecture of maize flowering time. *Science* **325**: 714–718.

Chen CY, Misztal I, Aguilar I, Legarra A, Muir WM (2011). Effect of different genomic relationship matrices on accuracy and scale. *J Anim Sci* **89**: 2673–2679.

Cros D, Denis M, Sánchez L, Cochard B, Flori A, Durant-Gasselín T *et al.* (2015). Genomic selection prediction accuracy in a perennial crop: case study of oil palm *Elaeis guineensis* Jacq. *Theor Appl Genet* **128**: 397–410.

Denis M, Bouvet J-M (2013). Efficiency of genomic selection with models including dominance effect in the context of *Eucalyptus* breeding. *Tree Genet Genomes* **9**: 37–51.

Dudley JW, Johnson GR (2009). Epistatic models improve prediction of performance in corn. *Crop Sci* **49**: 763–770.

Falconer DS, Mackay TFC (1996). *Introduction to Quantitative Genetics*, edn. 4. Longmans Green: Harlow, Essex, UK.

Forni S, Aguilar I, Misztal I (2011). Different genomic relationship matrices for single-step analysis using phenotypic, pedigree and genomic information. *Genet Sel Evol* **43**: 1.

Gallais A (1990). *Théorie de la Sélection en Amélioration des Plantes*. Masson: Paris, France, 588 p.

Gallais A (2009). *Hétérosis et Variétés Hybrides en Amélioration des Plantes*. Quae: Versailles, France, p 376.

Gamal El-Dien O, Ratcliffe B, Kláp J, Chen C, Porth I, El-Kassaby YA (2015). Prediction accuracies for growth and wood attributes of interior spruce in space using genotyping-by-sequencing. *BMC Genomics* **16**: 370.

Gilmour AR, Gogel B, Cullis BR, Thompson R (2006). *ASReml, User Guide. Release 2.0*. VSN International Ltd: Hemel Hempstead, UK.

Guo T, Yang N, Tong H, Pan Q, Yang X, Tang J *et al.* (2014). Genetic basis of grain yield heterosis in an 'immortalized F2' maize population. *Theor Appl Genet* **127**: 2149–2158.

Hayes BJ, Bowman PJ, Daetwyler HD, Kijas JW, Van der Werf JHJ (2012). Accuracy of genotype imputation in sheep breeds. *Anim Genet* **43**: 72–80.

Heslot N, Jannink J-L, Sorrells ME (2015). Perspectives for Genomic Selection Applications and Research in Plants. *Crop Sci* **55**: 1–12.

Hill WG, Robertson A (1968). Linkage disequilibrium in finite populations. *Theor Appl Genet* **38**: 226231.

Horsley TN, Johnson SD (2007). Is *Eucalyptus* cryptically self-incompatible? *Annals of Botany* **100**: 1373–1378.

Hu Z, Li Y, Song X, Han Y, Cai X, Xu S *et al.* (2011). Genomic value prediction for quantitative traits under the epistatic model. *BMC Genet* **12**: 15.

Ibáñez-Escriche N, Fernando RL, Toosi A, Dekkers JCM (2009). Genomic selection of purebreds for crossbred performance. *Genet Sel Evol* **41**: 12.

Kinghorn BP, Hickey JM, Van Der Werf JHJ (2010). Reciprocal recurrent genomic selection for total genetic merit in crossbred individuals. Proceedings of the 9th World Congress on Genetics Applied to Livestock Production. German Society for Animal Science: Leipzig, Germany, Paper 0036.

Lee H, Goddard ME, Visscher PM, Van der Werf JHJ (2010). Using the realized relationship matrix to disentangle confounding factors for the estimation of genetic variance components of complex traits. *Genet Sel Evol* **42**: 22.

Legarra A, Aguilar I, Misztal I (2009). A relationship matrix including full pedigree and genomic information. *J Dairy Sci* **92**: 4656–4663.

Lepoittevin C, Rousseau JP, Guillemin A, Gaurvit C, Besson F, Hubert F *et al.* (2011). Genetic parameters of growth, straightness and wood chemistry traits in *Pinus pinaster*. *Ann Forest Sci* **68**: 873–884.

Li CQ, Song L, Zhao HH, Wang QL, Fu YZ (2014). Identification of quantitative trait loci with main and epistatic effects for plant architecture traits in Upland cotton *Gossypium hirsutum* L. *Plant Breeding* **133**: 390–400.

Lo LL, Fernando RL, Grossman M (1997). Genetic evaluation by BLUP in two-breed terminal crossbreeding systems under dominance. *J Anim Sci* **75**: 2877–2884.

Lorenzana R, Bernardo R (2009). Accuracy of genotypic value predictions for marker-based selection in biparental plant populations. *Theor Appl Genet* **120**: 151–161.

Luo X, Fu Y, Zhang P, Wu S, Tian F, Liu J *et al.* (2009). Additive and over-dominant effects resulting from epistatic loci are the primary genetic basis of heterosis in rice. *J Integr Plant Biol* **51**: 393–408.

Lynch M, Walsh B (1998). *Genetics and Analysis of Quantitative Traits*. Sinauer Associates, Inc.: Sunderland, MN, USA, 980 p.

Mackay TF (2014). Epistasis and quantitative traits: using model organisms to study gene–gene interactions. *Nat Rev Genet* **15**: 22–33.

Makgahlela ML, Knürr T, Aamand GP, Strandén I, Mäntysaari EA (2013). Single step evaluations using haplotype segments. *Interbull Bull* **47**: 217–221.

Massman JM, Gordillo A, Lorenzana RE, Bernardo R (2013). Genome wide predictions from maize single-cross data. *Theor Appl Genet* **126**: 13–22.

Melchinger AE, Utz HF, Piepho H-P, Zeng Z-B, Schön CC (2007). The Role of epistasis in the manifestation of heterosis: a systems-oriented approach. *Genetics* **177**: 1815–1825.

Muñoz PR, Resende Jr MFR, Gezan SA, Deon VRM, de los Campos G, Kirst M *et al.* (2014). Unraveling additive from non-additive effects using genomic relationship matrices. *Genetics* **198**: 1759–1768.

Mäki-Tanila A, Hill WG (2014). Influence of gene interaction on complex trait variation with multilocus models. *Genetics* **198**: 355–367.

Oehlert GW (1992). A note on the delta method. *Am Stat* **46**: 27–29.

Parvez AS, Rather AG, Warsi MZK (2007). Implications of epistasis in maize breeding. *Int J Plant Breeding Genet* **1**: 1–11.

R Development Core Team (2011). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing: Vienna, Austria.

Resende MDV, Resende MFR, Sansaloni CP, Petrolí CD, Missiaggia AA, Aguiar AM *et al.* (2012). Genomic selection for growth and wood quality in *Eucalyptus*: capturing the missing heritability and accelerating breeding for complex traits in forest trees. *New Phytol* **194**: 116–128.

Schnell FW (1963). The covariance between relatives in the presence of linkage. *Stat Genet Plant Breeding NAS-NRC* **982**: 468–483.

Shaffer JP (2007). Controlling the false discovery rate with constraints: the Newman-Keuls test revisited. *Biom J* **47**: 136–143.

Silva-Junior OB, Grattapaglia D (2015). Genome-wide patterns of recombination, linkage disequilibrium and nucleotide diversity from pooled resequencing and single nucleotide polymorphism genotyping unlock the evolutionary history of *Eucalyptus grandis*. *New Phytol*; e-pub ahead of print 16 June 2015; doi:10.1111/nph.13505.

Stuber CW, Cockerham CC (1966). Gene effects and variances in hybrid populations. *Genetics* **54**: 1279–1286.

Su G, Christensen OF, Ostensen T, Henryon M, Lund MS (2012). Estimating additive and non-additive genetic variances and predicting genetic merits using genome-wide dense single nucleotide polymorphism markers. *PLoS One* **7**: e45293.

Sun C, Van Raden PM, Cole BJ, O'Connell JR (2014). Improvement of prediction ability for genomic selection of dairy cattle by including dominance effects. *PLoS One* **9**: e103934.

- Technow F, Riedelsheimer C, Schrag TA, Melchinger AE (2012). Genomic prediction of hybrid performance in maize with models incorporating dominance and population specific marker effects. *Theor Appl Genet* **125**: 1181–1194.
- Tripiana V, Bourgeois M, Verhaegen D, Vigneron P, Bouvet J-M (2007). Combining microsatellites, growth, and adaptive traits for managing in situ genetic resources of *Eucalyptus urophylla*. *Can J Forest Res* **37**: 773–785.
- Van Raden PM (2008). Efficient methods to compute genomic predictions. *J Dairy Sci* **91**: 4414–4423.
- Wang D, El-Basyoni IS, Baenziger PS, Crossa J, Eskridge KM, Dweikat I (2012). Prediction of genetic values of quantitative traits with epistatic effects in plant breeding populations. *Heredity* **109**: 313–319.
- Wang H, Misztal I, Legarra A (2014). Differences between genomic-based and pedigree-based relationships in a chicken population, as a function of quality control and pedigree links among individuals. *J Anim Breed Genet* **131**: 445–451.
- Wardyn BM, Edwards JW, Lamkey KR (2007). The genetic structure of a maize population: the role of dominance. *Crop Sci* **47**: 467–474.
- Wenzl P, Carling J, Kudma D, Jaccoud D, Huttner E, Kleinhofs A *et al.* (2004). Diversity arrays technology DAuT for whole-genome profiling of barley. *Proc Natl Acad Sci USA* **101**: 9915–9992.
- Wu RL (2000). Partitioning of population genetic variance under multiplicative-epistatic gene action. *Theor Appl Genet* **100**: 743–749.
- Zapata-Valenzuela J, Whetten RW, Neale D, McKeand S, Isik F (2013). Genomic estimated breeding values using genomic relationship matrices in a cloned population of loblolly pine. *Genes Genomes Genet* **3**: 909–916.
- Zeng J, Toosi A, Fernando RL, Dekkers JCM, Garrick DJ (2013). Genomic selection of purebred animals for crossbred performance in the presence of dominant gene action. *Genet Sel Evol* **45**: 11.
- Zhou G, Chen Y, Yao W, Zhang C, Xie W, Hua J *et al.* (2012). Genetic composition of yield heterosis in an elite rice hybrid. *Proc Natl Acad Sci USA* **109**: 15847–15852.

Supplementary Information accompanies this paper on Heredity website (<http://www.nature.com/hdy>)

## ANNEX

Estimation of variance components with a single population model.

The single population progeny model was developed based on the relationship among the ( $c = 1130$ ) hybrid progenies.

$$y = X\beta + Z_{col}col + Z_{r;b}r : b + Z_{p}plot + Z_c a + Z_c d + \varepsilon - \text{single population model}$$

The definition of the fixed and environmental random effects was similar to that of the parent model. For the marker-based model using the 3303 SNPs, additive and dominance effects were defined by:  $a \sim N(0, \sigma_a^2 A_G)$  and  $d \sim N(0, \sigma_d^2 D_G)$ . Molecular-based relationship matrices among the 1130 clones were defined using the van Raden (2008) formula for  $A_G[c,c]$  and the formula by Su *et al.* (2012) for  $D_G[c,c]$ . Genetic effects of the pedigree-based model were defined as:  $a \sim N(0, \sigma_a^2 A_P)$  and  $d \sim N(0, \sigma_d^2 D_P)$ , with  $A_P[c,c]$  and  $D_P[c,c]$  calculated using pedigree information.

**Table A1 Environmental and genetic variance components and AIC of the different models considering a single population model**

Variance components	Pedigree-based relationship	Marker-based relationship
	Single P_A+D	Single G_A+D
$\sigma_{r;b}^2$	0.649	0.600
$\sigma_{col}^2$	0.019	0.033
$\sigma_a^2$	1.959	0.080
$\sigma_d^2$	1.097	0.010
$\sigma_{plot}^2$	2.331	3.492
$\sigma_e^2$	8.055	9.179
Loglik	-4416.060	-4444.057
Parameters	6	6
AIC	8844.120	8900.1138