



**HAL**  
open science

## Genomic prediction in a multi-generation *Eucalyptus globulus* breeding population

Geoffrey Haristoy, Laurent Bouffier, Luis Fontes, Luis Leal, Jorge Paiva,  
João-Pedro Pina, Jean-Marc Gion

► **To cite this version:**

Geoffrey Haristoy, Laurent Bouffier, Luis Fontes, Luis Leal, Jorge Paiva, et al.. Genomic prediction in a multi-generation *Eucalyptus globulus* breeding population. *Tree Genetics and Genomes*, 2023, 19 (1), pp.8. 10.1007/s11295-022-01579-2 . hal-04011581

**HAL Id: hal-04011581**

**<https://hal.inrae.fr/hal-04011581v1>**

Submitted on 30 Sep 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

[Click here to view linked References](#)

1 **Title:** Genomic prediction in a multi-generation *Eucalyptus globulus* breeding population

2

3 **Authors:**

4 Geoffrey Haristoy (**first author**)

5 Affiliation: CIRAD, UMR AGAP Institut, F-34398 Montpellier, France

6 Email: [geoffrey.haristoy@cirad.fr](mailto:geoffrey.haristoy@cirad.fr)

7 ORCID: 0000-0001-9914-7712

8

9 Laurent Bouffier

10 Affiliation: INRAE, UMR BIOGECO, F-33610 Cestas, France

11 Email: [laurent.bouffier@inrae.fr](mailto:laurent.bouffier@inrae.fr)

12 ORCID: 0000-0001-7493-5077

13

14 Luis Fontes

15 Affiliation: Altri Florestal, 2510-582 Olho Marinho, Portugal

16 Email: [luis.fontes@altri.pt](mailto:luis.fontes@altri.pt)

17

18 Luis Leal

19 Affiliation: Altri Florestal, 2510-582 Olho Marinho, Portugal

20 Email: [luis.leal@altri.pt](mailto:luis.leal@altri.pt)

21

22 Jorge A. P. Paiva

23 Affiliations: Associação CECOLAB (Collaborative Laboratory Towards Circular Economy), 3405-155 Oliveira

24 do Hospital, Portugal ; iBET, 2780-157 Oeiras, Portugal

25 Email: [jorge.pinto.paiva@gmail.com](mailto:jorge.pinto.paiva@gmail.com)

26 ORCID: 0000-0003-3162-4396

27

28 João-Pedro Pina

29 Affiliation: Altri Florestal, 2510-582 Olho Marinho, Portugal

30 Email: [pina.jp@gmail.com](mailto:pina.jp@gmail.com)

31

32 Jean-Marc Gion (**corresponding author and last author**)

33 Affiliation: CIRAD, UMR AGAP Institut, F-34398 Montpellier, France

34 Email: [jean-marc.gion@cirad.fr](mailto:jean-marc.gion@cirad.fr)

35 ORCID: 0000-0003-3958-7796

36

37 **Keywords:** *Eucalyptus globulus*, genomic selection, GBLUP, progeny validation, pedigree error, breeding  
38 programme

39

40 **Statements and Declarations:** The authors have no conflict of interest to declare.

41

42 **Abstract**

43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58

Genomic selection is a promising approach for reducing the length of the selection cycle in forest tree breeding. Its efficiency must be evaluated across generations for this purpose, but such studies have been performed for multi-generational breeding programmes in only a few forest tree species to date. We analysed a subset of the *Eucalyptus globulus* breeding population from the Portuguese company Altri Florestal. In total, 412 genotypes from three successive breeding generation were genotyped with 14,716 SNP markers. A comparison of pedigree-based and marker-based relationship coefficients allowed to correct several documented pedigree errors. Deregressed breeding values were estimated from phenotypic records for growth traits (height and diameter) and survival for 31 field trials distributed in one breeding zone in Portugal, and used as pseudo-phenotypes for genomic prediction models. Accuracy was assessed by cross-validation according to two main scenarios: i) a scenario based on a five random fold number, not taking generation into account ; ii) scenarios investigating progeny validation using parental generations to predict the progenies. Accuracy was highest after pedigree correction, and ranged from 0.46 to 0.60 for the first scenario, from -0.56 to 0.72 for parent/progeny scenarios, and from 0.34 to 0.78 when progenies were added to the calibration population. This genomic selection study provides promising insight for the Altri Florestal *Eucalyptus* breeding programme.

59 **Introduction**

60

61 *Eucalyptus globulus* (Tasmanian blue gum) is an evergreen broadleaf tree species endemic to southern  
62 Australia. This forest tree species is one of the most planted hardwood species in temperate regions worldwide. Its  
63 economic importance is mainly due to the suitability of its wood for pulp and paper production. This species is  
64 also characterized by rapid growth, challenging vegetative propagation by cuttings since it is a rooting recalcitrant  
65 species and its ability to adapt to harsh environmental conditions. *E. globulus* covers an area of 0.84 million  
66 hectares in Portugal (ICNF 2019), and 0.64 million hectares in Spain (MAPA 2019). Since the 1960s, *E. globulus*  
67 breeding programmes have been developed in Portugal, to support intensive silviculture. The Portuguese company  
68 Altri Florestal selects varieties of this species, mostly on the basis of wood productivity-related traits, such as  
69 growth, and traits related to adaptation (survival). The Altri Florestal breeding programme is still in the early stages  
70 of *E. globulus* domestication with only three generations so far, like most of the advanced genetic materials  
71 currently available in forestry (Jones et al. 2006; Borralho et al. 2007).

72 Breeding programmes were initially based on a few key genetic trials, but have gradually expanded, with  
73 an ever-increasing number of trials and phenotyped trees. The genetic performance of trees is commonly evaluated  
74 on the basis of genetic co-variances in known relatives arising from the pedigree, according to the individual mixed  
75 model (Henderson 1950, 1975). This method is based on a numerator relationship matrix (A matrix) derived from  
76 the pedigree, which provides information about the proportion of alleles expected to be identical by descent  
77 between two individuals (Mrode 2013). Such models have already proved effective for the estimation of genetic  
78 components, especially in cases of unbalanced data, a situation frequently encountered in forest tree breeding  
79 (Borralho 1995; Jarvis et al. 1995), but their application is hindered by approximations. Indeed, the base population  
80 is assumed to consist of unrelated founders, although some cryptic relatedness may exist (Powell et al. 2010). In  
81 addition, pedigree information is rarely fully documented, particularly in open crossing strategies (Klápště et al.  
82 2014). Finally, various identity or pedigree errors may occur during the breeding process, from the greenhouse to  
83 the field (pollination, seedlings, cuttings, plantation, etc.), compromising the genetic evaluation (Ericsson 1999).  
84 Such errors are cumulative over generations, and the earlier they occur in breeding cycles, the more likely they are  
85 to have a significant and damaging long-term impact on genetic evaluation. Recent decades have seen considerable  
86 advances in molecular biology, leading to the development of new tools for forest tree breeders. Pedigree-based  
87 relationships are based on expectations of the sharing of genomic material between individuals, but high-  
88 throughput genotyping has made it possible to estimate the proportion of alleles common to individuals (actual or  
89 realised relationship) precisely, including the within-family variability arising from Mendelian sampling (Hill and  
90 Weir 2011). This information can be summarised in a genomic relationship matrix (G matrix) (VanRaden 2008),  
91 and can be used in a genomic selection (GS) strategy through the so-called GBLUP (Genomic Best Linear  
92 Unbiased Prediction) methodology (Meuwissen et al. 2001), which involves replacing the A matrix with a G matrix  
93 in the individual mixed model. The deviation of A and G matrices is one of the key factors promoting the use of  
94 GS to obtain more accurate breeding values (Hayes et al. 2009b). GS exploits the linkage disequilibrium (LD)  
95 between high-throughput molecular data and targeted traits. Based on a calibration population of several hundred  
96 phenotyped and genotyped individuals, a predictive model is built to predict the genetic values of genotyped  
97 individuals. *Eucalyptus* breeding should benefit greatly from GS, as this approach makes early selection possible.  
98 Indeed, age-age correlations for both height and diameter are generally low-to-moderate between the ages of one

99 to three years, delaying progeny evaluation and clonal trials (Salas et al. 2014). *Eucalyptus* breeding programmes  
1 100 therefore require about 12-16 years, from initial recombination to clonal selection and operational deployment  
2  
3 101 (Rezende et al. 2014).

4 102 Over the last decade, GS has benefited from many proof-of-concept studies for the genus *Eucalyptus*,  
5  
6 103 with more than 20 publications and promising results, at least as good as those obtained by conventional  
7 104 phenotypic selection in most studies (Lebedev et al. 2020; Ahmar et al. 2021). However, such studies have  
8  
9 105 generally focused on a limited number of progeny trials, whereas most breeding programmes involve at least  
10 106 several dozen, if not hundreds of trials. Furthermore, only a few studies have explored GS across generations,  
11  
12 107 whereas such approach would be required in current breeding programmes based on recurrent selection strategies  
13 108 (Grattapaglia 2017). This study aims to fill this gap by applying GS in the context of the advanced *Eucalyptus*  
14 109 *globulus* breeding programme of the Portuguese company Altri Florestal. Our main objectives were: i) to highlight  
15 109 pedigree errors by comparing pedigree-based (A matrix) and marker-based relationship coefficients (G matrix)  
16 110 and to investigate the consequences of such errors for breeding value prediction; ii) to assess the accuracy of GS  
17 111 for three major traits in forest tree breeding (height, diameter and survival); and iii) to investigate GS accuracy  
18 111 over generations.  
19 112  
20  
21 113  
22  
23 114  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## 115 **Materials and methods**

1 116

### 3 117 1. Biological resources

4 118

5 119 This study was performed on a subset of the *Eucalyptus globulus* breeding population of Altri Florestal,  
6 120 referred to here as POP<sub>TOT</sub>. In total, 81,520 *E. globulus* genotypes were evaluated in 31 progeny and clonal trials  
7 121 (92,679 trees) planted in the coastal region of Portugal between 1986 and 2009. This region, which is one of the  
8 122 four breeding zones defined by Altri Florestal based on temperatures and precipitations, is considered to be an  
9 123 “unrestricted” environment as it is not subjected to major drought or cold events. According to the documented  
10 124 pedigree, POP<sub>TOT</sub> encompassed three breeding generations, with 2.8% of the genotypes (2,250) of unknown  
11 125 parentage belonging to the base population called G0. The next generation, corresponding to the first improved  
12 126 generation (G1), contained 41.6% of the genotypes (33,909), each with at least one documented G0 parent. The  
13 127 remaining 55.6% of the genotypes (45,361 genotypes) had at least one documented grandparent and belonged to  
14 128 the second improved breeding generation (G2). The G0 genotypes were documented as originating from selections  
15 129 from populations native to Australia or growing in various stands in Portugal, USA, Chile, Spain and Uruguay. A  
16 130 subsample of 412 POP<sub>TOT</sub> genotypes corresponding to the genotypes available in clonal archives was selected for  
17 131 high-throughput molecular genotyping: 46 from G0, 292 from G1 and 74 from G2 (Fig. 1). Three full-sib (FS)  
18 132 families accounted for 81% of the genotyped G2 trees: FS1 (15 genotypes), FS2 (29 genotypes) and FS3 (16  
19 133 genotypes).

20 134

### 21 135 2. SNP genotyping

22 136

23 137 Total genomic DNA was extracted from dried leaves from mature trees with an adapted CTAB protocol  
24 138 (Doyle 1991). DNA concentration was determined with a Quantit fluorometer (Invitrogen, Carlsbad, USA).  
25 139 Single-nucleotide polymorphism (SNP) genotyping was performed with the commercial Axiom Euc72K SNP  
26 140 Array (Affymetrix, Santa Clara, USA), anchored on the 11 linkage groups of the reference *Eucalyptus* genome  
27 141 (Myburg et al. 2014). A first set of thresholds was applied with the default settings of Axiom Suite Analysis v5.0.1  
28 142 software (Affymetrix, Santa Clara, USA) and an SNP call rate threshold of 97% (detailed in Table S1). SNPs were,  
29 143 thus, classified as polymorphic (PolyHighResolution and NoMinorHom), monomorphic, or failed SNPs  
30 144 (CallRateBelowThreshold, OTV, Other). The highest quality polymorphic SNPs were then selected according to  
31 145 a second set of thresholds: a Fisher's linear discriminant (FLD) above 5, a homozygote ratio offset (HomRo) above  
32 146 1, a heterozygous strength offset (HetSO) above 0, and a minor allele frequency (MAF) higher than 0.05. Eight  
33 147 genotypes were genotyped twice to assess the repeatability of the genotyping experiment, calculated as the  
34 148 proportion of identical SNP alleles between two samples of a given genotype.

35 149

### 36 150 3. Comparison of pedigree-based and genomic relationship matrices

37 151

38 152 A dedicated R script was used to check the uniqueness of each genotype based on SNPs. Pairs of genotypes  
39 153 with different identities and more than 99% SNP alleles in common were considered synonymous. In this case,  
40 154 the highest call rate profile was retained, and pairs of synonymous genotypes were renamed under the same identity

41

42

43

44

45

155 label in the pedigree and field measurement files. The parents were considered to be unknown if the initially  
 156 documented parents of two synonymous genotypes were different. After identity correction, the final set of unique  
 157 genotypes was named POP<sub>GS</sub>. An initial expected additive numerator relationship matrix ( $A_I$  matrix) was calculated  
 158 for POP<sub>GS</sub> with the R package kinship2 (Sinnwell et al. 2014) based on the documented pedigree. A genomic  
 159 relationship matrix ( $G$  matrix) was also calculated from the observed allele frequencies (VanRaden 2008), with  
 160 the R package AGHMatrix (Amadeu et al. 2016):

$$G = \frac{(M-P)(M-P)'}{2 \sum_{i=1}^m p_i(1-p_i)} \quad (\text{Eq. 1})$$

161 where  $M$  is a matrix of dimension  $n \times m$  ( $n$  is the number of individuals and  $m$  the number of loci) giving the  
 162 genotype at each locus  $i$ , coded as 1 for minor allele homozygous, 2 for heterozygous, and 3 for major allele  
 163 homozygous.  $P$  is the matrix of allele frequencies ( $n \times m$ ) for all individuals, which takes the following form  
 164  $2 \times (p_i - 0.5)$  where  $p_i$  is the frequency of the least frequent allele at the considered locus  $i$ . The documented  
 165 parent-progeny relationship (P/P) were checked by counting the number of non-concordant SNPs for each  
 166 documented P/P. P/P highlighting fewer than 115 non-concordant SNPs were considered as “true”, whereas P/P  
 167 harbouring more than 115 non-concordant SNPs were considered “false”. P/P involving non-genotyped parent  
 168 were considered to be “undetermined”. For each progeny involved in an “undetermined” P/P,  $G$  coefficients for  
 169 all documented full and half siblings were compared to the expected  $A_I$  coefficients. If more than 40% of the  
 170 pairwise differences between  $G$  and  $A_I$  coefficients exceeded the threshold of 0.2 (Thumma et al. 2022), the  
 171 pedigree was considered to be inconsistent. For such individuals, the identity of the initial parents was either  
 172 considered to be unknown or was replaced by a new parent identity if allelic patterns were found to be concordant.  
 173 The corrected pedigree was used to generate the corrected numerator relationship matrix  $A_C$ . The  $A_I$ ,  $A_C$  and  $G$   
 174 matrices were visualised with the ggplot2 (Wickham 2016) and ComplexHeatmap (Gu et al. 2016) packages in  
 175 the R statistical environment (Rstudio Team 2021).  
 176  
 177

#### 178 4. Pseudo-phenotype estimates

179  
 180 Trials have been measured at various ages for growth traits and survival, but only the most recent  
 181 measurements were considered for each trial (i.e. ages ranging from 6 to 17 years depending on the trial). The  
 182 diameter over-bark at breast height (DBH) was calculated as the mean of two tree calliper measurements taken at  
 183 right angles. Height (HT) was measured with a telescopic pole. Survival (SV) was equal to 1 for living trees, and  
 184 0 for dead trees. For each trait, estimated breeding values (EBV) were obtained with ASREML 4.0 (Butler et al.  
 185 2017) from a BLUP meta-analysis routinely implemented at Altri Florestal as detailed in Borralho et al. (2018).  
 186 This mixed-model included fixed effects (trial and replicates within trial) and random effects (incomplete blocks  
 187 within replicates, additive genetic effect, full-sib family, clone within full-sib family), as well as pedigree  
 188 relationships across all trees (numerator relationship matrix). The variance components and heritability estimates  
 189 considered for the meta-analysis were described in Table 1 following Borralho et al. (2018). Considering this  
 190 mixed-model, an initial estimated breeding value (EBV<sub>I</sub>) was calculated with the  $A_I$  matrix, and a corrected  
 191 estimated breeding value (EBV<sub>C</sub>) was calculated with the  $A_C$  matrix. EBV accuracy was estimated as follows (Isik  
 192 et al. 2017):

$$r = \sqrt{1 - \frac{S^2}{(1+F)\sigma^2_A}} \quad (\text{Eq.2})$$

where  $S$  is the standard error of the EBV,  $F$  is the coefficient of inbreeding and  $\sigma_A^2$  is the additive genetic variance. Additional EBV<sub>C</sub> values were calculated with the BLUP meta-analysis based on truncated phenotypic data: i) the EBV<sub>C-T01</sub> values considering only phenotypic data from the G0 and G1 genotypes, and ii) the EBV<sub>C-T2</sub> values considering only phenotypic data from G2 genotypes. EBV<sub>C-T01</sub> and EBV<sub>C-T2</sub> were, thus, estimated with independent phenotypic data sets. Considering EBV as phenotypes in genomic prediction may introduce bias and heterogeneity (Garrick et al. 2009). EBV was therefore deregressed and weighted (dEBV) following Garrick et al. (2009) based on estimates of heritability and without removal of the parent average effect, as many individuals had unknown fathers. The resulting dEBV<sub>I</sub>, dEBV<sub>C</sub>, dEBV<sub>C-T01</sub>, dEBV<sub>C-T2</sub> were used as pseudo-phenotypes for GS.

**Table 1** Variances associated with each random effects and heritability ( $h^2$ ) estimated for HT, DBH and SV.

$\sigma_A^2$ ,  $\sigma_r^2$ ,  $\sigma_f^2$ ,  $\sigma_c^2$  and  $\sigma_e^2$  are the variances associated with the following random effects: additive genetic effect, incomplete block within replicate, full-sib family, clone within full-sib family and residuals, respectively.

Heritability was estimated as:  $h^2 = \frac{\sigma_A^2}{\sigma_r^2 + \sigma_A^2 + \sigma_f^2 + \sigma_c^2 + \sigma_e^2}$

Trait	$\sigma_A^2$	$\sigma_r^2$	$\sigma_f^2$	$\sigma_c^2$	$\sigma_e^2$	$h^2$
HT	0.25	0.19	0.03	0.05	0.72	0.20
DBH	0.15	0.08	0.04	0.04	0.80	0.14
SV	0.45	-	0.05	0.25	3.29	0.11

## 5. Genomic prediction models

Genomic estimated breeding values (GEBV<sub>I</sub>, GEBV<sub>C</sub>, GEBV<sub>C-T01</sub> and GEBV<sub>C-T2</sub>) were estimated for each trait, from dEBV<sub>I</sub>, dEBV<sub>C</sub>, dEBV<sub>C-T01</sub>, and dEBV<sub>C-T2</sub>, respectively. The following GBLUP model was implemented with the BreedR R package (Munoz and Rodriguez 2020):

$$y = X\mu + Za + e \quad (\text{Eq. 3})$$

in which  $y$  is the vector of pseudo-phenotypes (dEBV<sub>I</sub>, dEBV<sub>C</sub>, dEBV<sub>C-T01</sub>, or dEBV<sub>C-T2</sub>),  $\mu$  the population mean,  $a$  the vector of random additive genetic effects and  $e$  the vector of residuals effects.  $X$  and  $Z$  are the incidence matrices for  $\mu$  and  $a$  effects. The vector  $a$  was assumed to follow a normal distribution  $a \sim N(0, G\sigma_a^2)$ , where  $G$  is the realised relationship matrix and  $\sigma_a^2$  the variance of additive effects. The vector  $e$  followed a normal distribution with  $e \sim N(0, I\sigma_e^2)$ , where  $\sigma_e^2$  is the residual variance. The solutions for the random genetic effects of Eq. 3 are the genomic estimated breeding values, GEBV. GBLUP shrunk marker effects uniformly, assuming a centred normal distribution and a common variance for marker effects.

## 6. Cross-validation scenarios

Nine cross-validation scenarios were tested to assess the accuracy of GS for each trait (Table 2). In the random cross-validation scenario (S0), the three generations were all included in the calibration population (CP) and the validation population (VP), whereas the others cross-validation scenarios were designed for the investigation of GS accuracy over breeding generations by using different configurations of pseudo-phenotypes for both the CP and the VP.



**Table 2** Description of the nine GS scenarios (S0, S1<sub>a</sub>, S1<sub>b</sub>, S2<sub>a</sub>, S2<sub>b</sub>, S2<sub>c</sub>, S2<sub>d</sub>, S3<sub>a</sub> and S3<sub>b</sub>) regarding the pseudo-phenotype used, the number of genotypes in the calibration (CP) and validation (VP) populations (the number of genotypes indicated is after pedigree correction), the number of iterations performed, and the corresponding accuracy

\* randomly selected in POP<sub>GS</sub>; \*\* 20% of the G1 from POP<sub>GS</sub>; \*\*\* 20% of the three main G2 FS, and 20% of all the other G2 families; \*\*\*\* mean scenario accuracy over the 100 iterations.

Scenarios	CP		VP		Iterations	Accuracy
	Pseudo-phenotypes	Genotypes	Pseudo-phenotypes	Genotypes		
S0	dEBV <sub>C</sub>	321 G0/G1/G2 *	dEBV <sub>C</sub>	80 G0/G1/G2 *	100	$\rho(\text{GEBV}_C, \text{dEBV}_C)$ ****
S1 <sub>a</sub>	dEBV <sub>C</sub>	44 G0	dEBV <sub>C</sub>	284 G1	-	$\rho(\text{GEBV}_C, \text{dEBV}_C)$
S1 <sub>b</sub>	dEBV <sub>C</sub>	44 G0, 57 G1 **	dEBV <sub>C</sub>	227 G1	100	$\rho(\text{GEBV}_C, \text{dEBV}_C)$ ****
S2 <sub>a</sub>	dEBV <sub>C</sub>	44 G0, 284 G1	dEBV <sub>C</sub>	73 G2	-	$\rho(\text{GEBV}_C, \text{dEBV}_C)$
S2 <sub>b</sub>	dEBV <sub>C</sub>	44 G0, 284 G1	dEBV <sub>C-T2</sub>	73 G2	-	$\rho(\text{GEBV}_C, \text{dEBV}_{C-T2})$
S2 <sub>c</sub>	dEBV <sub>C-T01</sub>	44 G0, 284 G1	dEBV <sub>C</sub>	73 G2	-	$\rho(\text{GEBV}_{C-T01}, \text{dEBV}_C)$
S2 <sub>d</sub>	dEBV <sub>C-T01</sub>	44 G0, 284 G1	dEBV <sub>C-T2</sub>	73 G2	-	$\rho(\text{GEBV}_{C-T01}, \text{dEBV}_{C-T2})$
S3 <sub>a</sub>	dEBV <sub>C</sub>	44 G0, 284 G1, 14 G2 ***	dEBV <sub>C</sub>	59 G2	100	$\rho(\text{GEBV}_C, \text{dEBV}_C)$ ****
S3 <sub>b</sub>	dEBV <sub>C</sub>	44 G0, 284 G1, 14 G2 ***	dEBV <sub>C-T2</sub>	59 G2	100	$\rho(\text{GEBV}_C, \text{dEBV}_{C-T2})$ ****

For the S0 scenario, 80% of POP<sub>GS</sub> (321 genotypes) were randomly assigned to the CP, with the remaining 20% (80 genotypes) were used as the VP (5 folds, 100 iterations). For S1<sub>a</sub>, only G0 genotypes were included in the CP (44 genotypes) for prediction of all the G1 genotypes in the VP. S1<sub>b</sub> was similar to S1<sub>a</sub> but 20% of the G1 genotypes (57 genotypes) were added to the CP (100 iterations). For S2<sub>a</sub>, all genotyped G0 and G1 were included in the CP (328 genotypes), and GEBV<sub>C</sub> were predicted in the VP for the 73 G2 genotypes. In S2<sub>b</sub>, prediction accuracy was estimated relative to dEBV<sub>C-T2</sub>. The S2<sub>c</sub> scenario used dEBV<sub>C-T01</sub> in the CP to predict the G2 GEBV<sub>C-T01</sub> of the VP, which was compared to dEBV<sub>C</sub> for accuracy estimation. The S2<sub>d</sub> used the same pseudo-phenotype in the CP to estimate GEBV<sub>C-T01</sub>, which was compared to dEBV<sub>C-T2</sub>. In S3<sub>a</sub> and S3<sub>b</sub>, 14 additional G2 genotypes were added to the G0 and G1 genotypes for the CP (342 genotypes), considering 20% of the three main FS families and 20% of all the other families used to estimate the GEBV<sub>C</sub> for the remaining G2 (59) in the VP (100 iterations). Prediction accuracy was estimated relative to either dEBV<sub>C</sub> (S3<sub>a</sub>) or dEBV<sub>C-T2</sub> (S3<sub>b</sub>). When 100 iterations were considered (S0, S1<sub>b</sub>, S3<sub>a</sub> and S3<sub>b</sub>), the accuracy was calculated as the mean accuracy over the 100 iterations.

254 **Results**

1 255

2 256 1. Genotyping results

3 257

4 258 The first filtering of 412 genotypes with 68,055 SNPs classified 28,242 SNPs as polymorphic SNPs (41.5%),  
 5 259 23,165 as monomorphic (34.0%), and 16,648 as failed (24.5%). The second filtering step selected the highest  
 6 260 quality polymorphic SNPs: 14,716 SNPs uniformly distributed over the 11 chromosomes of the reference genome  
 7 261 (Table 3). This final set of SNPs was used for pedigree correction and genomic predictions. The mean sample call  
 8 262 rate over all genotype samples was 99.6%. Based on the replicated samples, repeatability was estimated at 99.98%.

9 263

10 264 **Table 3** Marker coverage of the 14,716 SNPs, anchored onto the 11 chromosomes (Chr) of the *E. grandis* v2.0  
 11 265 reference genome (Myburg et al. 2014; Bartholomé et al. 2015). SNP positions were retrieved from Affymetrix  
 12 266 documentation to determine if they were located within gene or not. \* Chromosome physical length in Mb from  
 13 267 Myburg et al. (2014) and their corresponding genetic length in cM based on the *Eucalyptus* composite map from  
 14 268 Hudson et al. (2012)

Chr	Number of SNPs	Chr length *		Density		Mean distance between SNPs (kb)	Number of SNPs within gene	Proportion of SNPs within gene (%)
		(Mb)	(cM)	(SNP/Mb)	(SNP/cM)			
1	1,192	40.3	93.8	29.6	12.7	38	805	68
2	1,938	64.2	102.1	30.2	19.0	31	1,253	65
3	1,932	80.1	105.6	24.1	18.3	43	1,249	65
4	786	42.0	80.9	18.7	9.7	51	479	61
5	1,596	74.7	95.9	21.4	16.6	48	1,002	63
6	1,569	54.0	125.3	29.0	12.5	37	1,025	65
7	1,275	52.4	87.7	24.3	14.5	43	872	68
8	1,568	74.3	137.3	21.1	11.4	46	950	61
9	934	39.0	82.9	24.0	11.3	41	569	61
10	955	39.4	97.8	24.2	9.8	39	614	64
11	971	45.5	97.3	21.4	10.0	46	625	64
Genome	14,716	605.9	1,106.5	24.4	13.3	42	9,443	64

15 269

16 270

17 271 The number of SNPs per chromosome (Table 3) ranged from 786 (Chr4) to 1,938 (Chr2) with a mean of 1,338  
 18 272 SNPs per chromosome (mean density of 24.4 SNPs per Mb). The largest distance observed between two  
 19 273 neighbouring SNPs on the same chromosome ranged from 970 kb (Chr11) to 3760 kb (Chr5), and the smallest  
 20 274 distance between SNPs on the same chromosome ranged from 30 bp (Chr2, Chr3, Chr5, Chr9) to 40 bp (Chr11).  
 21 275 Considering the composite linkage map of *Eucalyptus* (Hudson et al. 2012), the number of SNPs per cM ranged  
 22 276 from 9.7 (Chr4) to 19.0 (Chr2). Based on the *E. grandis* reference genome, 64% of the SNPs (9,443 SNPs) were  
 23 277 located within a gene corresponding to 6,273 different genes out of the 36,376 genes estimated in this species.

24 278

25 279 2. Pedigree correction and its effect on estimated breeding value (EBV)

26 280

27 281 Eleven pairs of genotypes were identified as synonymous and were renamed as 11 unique genotypes,  
 28 282 accounting for 5.3% of the genotyping set (2 in G0, 8 in G1, and 1 in G2). Thus, POP<sub>GS</sub> contained 401 unique  
 29 283 genotypes spread over three generations: 44 G0, 284 G1, and 73 G2. After the correction of synonymous  
 30 284 genotypes, the G matrix (dimension  $n=401$ ) derived from POP<sub>GS</sub> was compared to the corresponding subsample

31 61

32 62

33 63

34 64

35 65

285  $A_I$  matrix of the same dimensions (Fig. 2a). The  $A_I$  matrix coefficients were discrete variables with six relationship  
1 286 classes (0, 0.0625, 0.125, 0.25, 0.5 and 1), whereas the G matrix coefficients followed a continuous distribution  
2 287 ranging from -0.20 to 1.44. Negative coefficients in the G matrix suggested that some individuals had fewer  
3 288 markers in common than expected on the basis of allele frequencies. The diagonal elements of the  $A_I$  matrix  
4 289 indicate an absence of inbreeding (relationship of 1), whereas the diagonal elements of the G matrix ranged from  
5 290 0.82 to 1.44.

6 291 Overall, 238 of the 401 genotypes were involved in at least one relationship with  $A_I=0.25$ , and 264 were  
7 292 involved in at least one relationship with  $A_I=0.5$  (and 198 were involved in both types of relationships). For both  
8 293 the 0.25 and 0.5  $A_I$  classes, bimodal and asymmetric distributions were observed (Fig. 3a and Fig. 3b) with the  
9 294 largest peak close to the expected value, and a second peak close to zero, suggesting the existence of errors in the  
10 295 documented pedigree. P/P consistencies were first assessed in POP<sub>GS</sub> by evaluating the compatibility of SNP  
11 296 between parents and progenies. A gap between P/P associated with less than 115 non-concordant SNPs and P/P  
12 297 associated with more than 911 non-concordant SNPs was highlighted (Fig. 4) which justified the threshold  
13 298 considered to make the distinction between “true” and “false” P/P. Both parents were considered to have been  
14 299 correctly documented for 112 genotypes (“true” P/P) which highlighted between 46 and 113 non-concordant SNPs  
15 300 (i.e. between 0.31% and 0.77% of the available SNPs). The presence of non-concordant SNPs for “true” P/P was  
16 301 explained by the genotyping repeatability below 100%. In the present study, the repeatability (99.98%) was  
17 302 probably overestimated as based on only eight replicates. At least one parent was incorrectly documented for 37  
18 303 genotypes (“false” P/P) which were associated with 912 to 1405 non-concordant SNPs (6.2% to 9.5% of the  
19 304 available SNPs). As some parents were not genotyped, 226 individuals had at least one “undetermined” P/P. Their  
20 305 consistencies were then analysed by comparing  $A_I$  and G matrix coefficients for full-sib and half-sib relationships.  
21 306 The distribution of the percentage of relationships differing from more 0.2 was represented in Fig. 5a. Our strategy  
22 307 was to identify and correct pedigree for genotypes which clearly highlighted pedigree inconsistencies. In addition  
23 308 to the “false” P/P previously detected, we identified eight additional genotypes for which more than 40% of G  
24 309 coefficients were at least 0.2 lower than the  $A_I$  coefficient, suggesting an incorrect parent. In total, for 45 genotypes  
25 310 (41 G1 and 4 G2), corresponding to 11.2% of POP<sub>GS</sub>, at least one wrong parent was identified. Reassignment of  
26 311 the correct parents was possible for 14 genotypes (14 G1), and the wrong parents were replaced by unknown  
27 312 parents in the documented pedigree for the remaining 31 genotypes (27 G1 and 4 G2). Finally, an  $A_C$  matrix was  
28 313 built from the corrected pedigree and compared to the G matrix. As illustrated in Fig. 2b, the  $A_C$  matrix better  
29 314 matched the G matrix than did the  $A_I$  matrix. The G matrix coefficient statistics for each A matrix coefficient class  
30 315 were reported in Table 4. As expected, the  $A_C$  matrix coefficients fitted better the G matrix coefficients for the two  
31 316 classes considered for pedigree correction (0.25 and 0.5), bringing the mean G matrix coefficient closer to the A  
32 317 matrix coefficient, and reducing the corresponding standard deviations. The applied pedigree correction resolved  
33 318 the bimodal distribution, resulting in a single peak (Fig. 3c and Fig. 3d) and highly decreased the percentage of  
34 319 relationships with  $A_C$  and G differing from more than 0.2 (Fig. 5b). There was also a minor effect in other classes,  
35 320 due to the resulting changes in the number of relationships per A class (Fig. S1). After correction, G2 genotypes  
36 321 from the three main FS families accounted for 77% of G2 genotypes in POP<sub>GS</sub>: FS1 (14 genotypes), FS2 (27  
37 322 genotypes) and FS3 (15 genotypes).

38 323 The G matrix also uncovered relatedness between individuals, that was not expected based on the documented  
39 324 pedigree, i.e. hidden relationships. For example, for the 145,672 pairwise relationships from the  $A_C=0$  class of

unrelated individuals, a total of 1,668 G matrix coefficients (1.1%) were greater than 0.2, mostly in G0 and G1. For the  $A_C=0.25$  class, 52 of the 10,482 relationships (0.5%) had G matrix coefficients greater than 0.5, suggesting the existence of a few undocumented FS or P/P relationships. For hidden relationships not involved in pedigree error, the documented pedigree was not modified. For the three traits considered, Pearson correlation coefficients ( $\rho$ ) between  $EBV_I$  and  $EBV_C$  revealed a small but non-negligible effect of pedigree modification on EBV estimates (Fig. S2). At the  $POP_{GS}$  level,  $\rho$  was 0.99 for HT, 0.97 for DBH and 0.98 for SV. However, when we considered only the 45 genotypes for which parentage errors were highlighted, slightly higher deviations between  $EBV_I$  and  $EBV_C$  were observed, with  $\rho$  equal to 0.92 for HT, 0.81 for DBH and 0.73 for SV (red dots in Fig. S2).

**Table 4** G relationship coefficients according to the expected  $A_I$  and  $A_C$  matrix coefficients, for the 401 genotypes of  $POP_{GS}$

Relationship matrices	Expected coefficient	Number of relationships	Mean	Sd	Min	Max
$A_I$	0	144,000	-0.02	0.07	-0.20	0.76
	0.0625	58	0.12	0.03	0.00	0.18
	0.125	148	0.14	0.09	-0.05	0.30
	0.25	11,870	0.11	0.11	-0.16	0.72
	0.5	4,324	0.40	0.16	-0.15	0.77
	1	401	1.04	0.13	0.82	1.44
$A_C$	0	145,672	-0.03	0.07	-0.20	0.76
	0.0625	54	0.13	0.02	0.08	0.18
	0.125	140	0.15	0.09	-0.05	0.30
	0.25	10,482	0.15	0.07	-0.04	0.72
	0.5	4,052	0.44	0.10	0.10	0.77
	1	401	1.04	0.13	0.82	1.44

### 3. Genetic values and correlations over generations

We investigated the trends in genetic value over generations, and compared genetic values obtained independently from G0/G1 and G2 phenotypic data through a truncation process. Table 5 presents the descriptive statistics for  $EBV_C$  for each trait in each generation of  $POP_{GS}$ . For growth traits, the lowest  $EBV_C$  means were obtained for G1 (0.15 for HT and 0.01 for DBH), with higher values obtained in G2 (0.20 for HT and 0.12 for DBH). For SV, mean  $EBV_C$  decreased slightly over generations, with mean values of 0.23 for G0, 0.21 for G1 and 0.16 for G2. The mean accuracy of all  $EBV_C$  was high, regardless of the generation considered, ranging from 0.72 to 0.92 in  $POP_{GS}$ . This high accuracy was explained by the use of a meta-analysis based on a high number of clonal copies and a high level of connectivity in the pedigree.

**Table 5** Descriptive statistics (mean and standard deviation) for  $EBV_C$  and its mean accuracy ( $r$ ) for tree height (HT), diameter at breast height (DBH) and survival (SV) for each generation of  $POP_{GS}$ . The number of genotypes in each generation is given (Size).

Generations	Size	HT			DBH			SV		
		Mean	Sd	r	Mean	Sd	r	Mean	Sd	r
All	401	0.17	0.34	0.88	0.03	0.33	0.84	0.2	0.33	0.73
G0	44	0.20	0.35	0.92	0.04	0.31	0.89	0.23	0.38	0.75
G1	284	0.15	0.34	0.87	0.01	0.32	0.83	0.21	0.34	0.72
G2	73	0.20	0.35	0.89	0.12	0.35	0.85	0.16	0.30	0.76

Correlations between the EBV<sub>C</sub> for the three traits were assessed using POP<sub>GS</sub> and within each breeding generation (G0, G1 and G2) (Fig. 6). Strong, highly significant positive genetic correlations ( $p$ -value < 0.001) between growth traits were observed for the whole sample ( $\rho=0.85$ ), whatever the generation considered (0.86 for G0, 0.83 for G1, and 0.93 for G2). HT and SV were significantly but weakly correlated in POP<sub>GS</sub> (0.18), but strongly correlated in G2 (0.55). Correlation between DBH and SV in POP<sub>GS</sub> was not significant at the 5% level, even though the G2 correlation was similar to that between HT and SV.

These results highlighted that POP<sub>GS</sub> was not representative of POP<sub>TOT</sub> for which an increase of EBV<sub>C</sub> over the generation was observed whatever the trait (Table S2). Based on EBV<sub>C</sub> correlations, slight differences were also observed between POP<sub>TOT</sub> and the subsample POP<sub>GS</sub> (Fig. S3), although the same overall trends were found with the strongest correlations between growth traits (0.87) and weak correlations between growth traits and SV (0.25 for HT and 0.21 for DBH). The differences in EBV<sub>C</sub> and trait correlations between POP<sub>TOT</sub> and POP<sub>GS</sub> could be explained by a sampling effect, with a smaller number of genotypes from G0 (44), G1 (284), and G2 (73). G2 genotypes in POP<sub>GS</sub> consisted of 13 FS families, three of which accounted for 77% of G2 genotypes. By contrast, G2 in POP<sub>TOT</sub> was composed of 45,360 genotypes with 914 FS. The genotyped population was therefore poorly representative of the total diversity of POP<sub>TOT</sub>, particularly for the first breeding generation (G0) and the last one (G2).

Genetic values of G1 genotypes were estimated based on POP<sub>GS</sub> through two independent processes based on truncated phenotypic data either EBV<sub>C-T01</sub> or EBV<sub>C-T2</sub>, but keeping relatedness between genotypes. The correlations between G1 EBV<sub>C-T01</sub> and G1 EBV<sub>C-T2</sub> were weak (0.26 for HT, 0.41 for DBH, 0.14 for SV) (Fig. S4), suggesting that the G0/G1 and G2 generations made slightly different contributions to global EBV<sub>C</sub> estimates. This was confirmed by the strong correlations between G1 EBV<sub>C-T01</sub> and G1 EBV<sub>C</sub> (0.94 for HT, 0.83 for DBH, and 0.84 for SV) whereas the correlations between G1 EBV<sub>C-T2</sub> and G1 EBV<sub>C</sub> were only moderate (0.47 for HT, 0.59 for DBH and 0.56 for SV). This trend was also observed for genetic value in G0. In contrast, the correlation between G2 EBV<sub>C-T2</sub> and G2 EBV<sub>C</sub> (0.98 for HT, 0.99 for DBH and 0.96 for SV) was stronger than that between G2 EBV<sub>C-T01</sub> and G2 EBV<sub>C</sub> (0.46 for HT, 0.43 for DBH and -0.45 for SV). Thus, for G2, global genetic value (G2 EBV<sub>C</sub>) was determined principally from the data collected for the G2 generation.

#### 4. Accuracy in cross-validation scenarios

The correlation between EBV and dEBV were very strong (>0.99 for all traits). We therefore used only dEBV as input variables for all GS scenarios. We first checked the effect of pedigree correction on GS accuracy. We plotted the correlation between GEBV and dEBV for the three traits according to the S0 scenario with and without pedigree correction (Fig. 7). For HT, DBH and SV, accuracy was slightly higher for dEBV<sub>C</sub> (0.46, 0.60 and 0.48, respectively) than for dEBV<sub>I</sub> (0.44, 0.55 and 0.48, respectively), with correction increasing accuracy by 5% for

389 HT and 9% for DBH (no change for SV). This benefit was also observed for the S1<sub>a</sub> and S2<sub>a</sub> scenarios (see Table  
1 390 S3), for which dEBV<sub>C</sub> gave slightly higher accuracies for HT (16% improvement for S1<sub>a</sub>, and 6% for S2<sub>a</sub>), DBH  
2 391 (39% for S1<sub>a</sub> and 1% for S2<sub>a</sub>), and SV (3% for S1<sub>a</sub> and 8% in S2<sub>a</sub>). Pedigree correction (11.2% of POP<sub>GS</sub>) had a  
3 392 significant effect on GS accuracy: i) mainly for the more precisely measured growth traits, and ii) for the S1<sub>a</sub>  
4 393 scenario generating estimates for G1, the generation for which the largest number of corrections were made (41 of  
5 394 the 45 genotypes corrected belonged to G1). The S2<sub>a</sub> scenario remained the most accurate for growth traits, with  
6 395 a smaller effect of correction, probably due to the smaller number of pedigree corrections for G2 genotypes.  
7 396 Whatever the scenario, our findings suggest that pedigree correction should be performed to correct pseudo-  
8 397 phenotypes before applying GS.

13 398 We then compared the nine scenarios (S0, S1<sub>a</sub>, S1<sub>b</sub>, S2<sub>a</sub>, S2<sub>b</sub>, S2<sub>c</sub>, S2<sub>d</sub>, S3<sub>a</sub> and S3<sub>b</sub>) considering only pseudo-  
14 399 phenotypes after pedigree correction, as shown in Fig. 8 (and Table S4). In the S0 scenario, genotypes were  
15 400 randomly assigned to either the CP (321 genotypes) or VP (80 genotypes), and the corresponding accuracies were  
16 401 higher for DBH (0.60), than for HT (0.46) or SV (0.48). This higher accuracy for DBH was observed in most of  
17 402 the scenarios tested. In the S1 scenarios, the addition of 20% of the G1 genotypes to the G0 genotypes used in the  
18 403 CP (S1<sub>b</sub>) improved accuracy to 55% for HT, 16% for DBH and 17% for SV, for the prediction of G1 genotypes.  
19 404 Overall, S1 accuracies were lower than the accuracies achieved for scenarios S0, S2<sub>a</sub> and S2<sub>b</sub>, which also used  
20 405 dEBV<sub>C</sub> as an input variable. In the S2 scenarios, dEBV<sub>C</sub> or dEBV<sub>C-T01</sub> for all G0 (44) and G1 (284) genotypes were  
21 406 chosen for the prediction of G2 genetic values in the VP (73). Accuracy was higher for dEBV<sub>C</sub> (S2<sub>a</sub> and S2<sub>b</sub>) than  
22 407 for dEBV<sub>C-T01</sub> (S2<sub>c</sub> and S2<sub>d</sub>), for which accuracy was non-significant for HT and, surprisingly, negative for SV. S2  
23 408 scenarios using dEBV<sub>C-T2</sub> in the VP (S2<sub>b</sub> and S2<sub>d</sub>) were slightly more accurate than those using dEBV<sub>C</sub> (e.g. the  
24 409 S2<sub>b</sub> scenario gave increases in accuracy of 4% for HT, 6% for DBH and 23% for SV relative to the S2<sub>a</sub> scenario).

31 410 The poor accuracy of scenarios S2<sub>c</sub> and S2<sub>d</sub> in comparison to scenarios S2<sub>a</sub> and S2<sub>b</sub> can be explained by  
32 411 the pseudo-phenotypes considered in CP. Indeed, even if dEBV<sub>C</sub> and dEBV<sub>C-T01</sub> were highly correlated for  
33 412 generation G1, this was not the case when considering specifically the parents of the three main G2 families (which  
34 413 represents 77% of the VP). Scenarios S2<sub>a</sub> and S2<sub>b</sub> calibrated with dEBV<sub>C</sub> (including information from G2  
35 414 phenotypes) were more efficient to predict G2 genotypes than scenarios S2<sub>c</sub> and S2<sub>d</sub> based on dEBV<sub>C-T01</sub> (including  
36 415 only G0 and G1 phenotypes). In the two S3 scenarios, following the addition of G2 genotypes to the CP (14 G2),  
37 416 accuracy was highest for the three traits when dEBV<sub>C-T2</sub> was used: 0.65 for HT, 0.78 for DBH, and 0.59 for SV  
38 417 (Fig. 8). This suggests that the addition of G2 genotypes to the CP greatly influences the quality of prediction for  
39 418 the remaining G2 genotypes in the VP due to increased relatedness between CP and VP. In scenarios S3<sub>a</sub> and S3<sub>b</sub>,  
40 419 the accuracies for the three traits were more dispersed than in S0 and similar to those in S1<sub>b</sub>, scenarios for which  
41 420 iterations were also performed. This suggests that composition of CP was affecting prediction accuracy and thus  
42 421 its optimisation could maximise the accuracy of GS predictions.

50 422

## 423 Discussion

1 424

2  
3 425 GS implementation in an advanced forest tree breeding programme requires a better knowledge of the change  
4 426 in genomic prediction accuracy for quantitative traits over breeding generations. Such investigations have been  
5  
6 427 conducted in conifers (Bartholomé et al., 2016; Isik et al., 2016; Thistlethwaite et al., 2019), but not in eucalypts.  
7 428 The advanced *E. globulus* breeding programme of Altri Florestal focuses on growth and survival, providing a great  
8  
9 429 opportunity to evaluate GS accuracy over three breeding generations.

10 430 The quantitative genetics of growth traits have been described in detail for *E. globulus*, with low to medium  
11  
12 431 heritabilities, suggesting polygenic determinism for both primary and secondary growth traits (Lopez et al. 2002;  
13 432 Raymond 2002; Potts et al. 2004). The genetic architecture of this species has been studied and a large number of  
14  
15 433 quantitative trait loci (QTLs) have been localised to different linkage groups, varying over time and/or  
16 434 environments (Freeman et al. 2013; Bartholomé et al. 2013, 2020). Genetic control has been shown to be mostly  
17  
18 435 additive, although non-negligible dominance effects have also been identified (Denis and Bouvet 2013; Tan et al.  
19 436 2018; Thavamanikumar et al. 2020). In *E. globulus*, the reported coefficients for genetic correlations between  
20  
21 437 height and diameter range from 0.55 to 0.93 (Volker et al. 1998; Hamilton et al. 2010; Rojas 2017). We observed  
22 438 strong additive genetic correlations between height and diameter for each generation (G0, G1, G2). Fewer data  
23  
24 439 have been published on tree survival, even though this has become a key breeding objective for companies wishing  
25 440 to expand their planting areas to less optimal climatic conditions (Costa e Silva et al. 2008). Survival can be defined  
26  
27 441 as the ability of a genotype to cope with a set of undefined environmental factors (both biotic and abiotic  
28 442 constraints), which is particularly important in the early stages of the tree's life (Lopez et al. 2002). Reported  
29  
30 443 heritabilities for survival in *E. globulus* are low to moderate, ranging from 0.02 to 0.38 (Chambers et al. 1996;  
31 444 Lopez et al. 2002; Hamilton et al. 2015; Mora and Serra 2014). The genetic determinism of survival may varies  
32  
33 445 over time, with the age of the tree (Dutkowski and Potts 1999) and the breeding generation considered. For long-  
34 446 lived species, such as forest trees, the survival recorded at the start of a breeding programme may depend on  
35  
36 447 genetic drivers different from those in contemporary measurements, due to changes in climate and the emergence  
37  
38 448 of new pests. Conflicting results concerning the correlation between survival and growth have been published for  
39 449 *E. globulus*, from weak negative genetic correlations (Lopez et al. 2002; Mora and Serra 2014) to highly positive  
40  
41 450 genetic correlations (Hamilton et al. 2010), and from non-significant phenotypic correlations (Lopez et al. 2002)  
42 451 to significant phenotypic correlations in a wide-ranging collection of open-pollinated *E. globulus* seeds from parent  
43  
44 452 trees growing in native stands in Australia (Dutkowski and Potts 1999). We found weak genetic correlations  
45 453 between growth traits and survival in the studied breeding zone characterised by no major constraints related to  
46  
47 454 coldness and drought. These correlations suggest the opportunity to select these two traits without trade-offs  
48 455 through this specific breeding zone, highlighting the importance of evaluating both of them in GS-based breeding  
49  
50 456 strategies.

51 457 Increasing numbers of studies are investigating the presence of pedigree errors in breeding programmes for  
52  
53 458 forest trees (Isik 2014). Two principal types of error have been highlighted: i) identity errors (synonymous labels  
54 459 or genetically different clonal replicates), and ii) parentage errors, when either one or both documented parents are  
55  
56 460 incorrect. In eucalypts, mislabelled ramets were found for four of 10 commercial clones from several organisations  
57 461 (Keil and Griffin 1994). Reported rates of parentage errors are highly variable, ranging from 2.8% in *Picea rubens*  
58  
59 462 (Doerksen and Herlinger 2008), to 30.2% in *Pseudotsuga menziesii* and 33.3% in *Pinus taeda* (Adams et al. 1988).

463 In *Pinus sylvestris* L. seed orchards, ramet assignment error rates range from 5.8% to 37.7% (Przybylski et al.  
1 464 2019). In a breeding population of *Pinus radiata*, 10% of documented relationships were found to be incorrect  
2 465 after pedigree verification (Kumar and Richardson 2005). Some progenies from open-pollinated families of *Picea*  
3 466 *glauca* were found to have a relationship coefficient of zero, suggesting pedigree errors (Gamal El-Dien et al.  
4 467 2016). Most of these pedigree verifications compared allelic consistencies between parents and progenies. We  
5 468 propose here an additional method based on comparison between A and G matrices applied to siblings for detecting  
6 469 pedigree errors. Our original approach identified 11.2% parentage errors in addition to the 5.3% identity errors  
7 470 detected by genomic fingerprinting. Our results confirmed the efficacy of SNP analyses for revealing incorrectly  
8 471 inferred relationships between individuals and for identifying previously unknown relationships as shown  
9 472 previously (Munoz et al. 2014; Tan et al. 2017; Lenz et al. 2020; Thumma et al. 2022). It remains tricky to define  
10 473 suitable thresholds for parentage errors in cases of large variances of pairwise relatedness estimators (Blouin 2003),  
11 474 as observed here within FS families. However, individuals for which more than 40% of G matrix coefficients  
12 475 deviated by more than 0.20 from the expected A coefficient were considered to have parentage errors. Our  
13 476 conservative approach allowed the most glaring pedigree errors to be corrected, thereby keeping the risk of false  
14 477 correction low.

15 478 In most forest tree breeding programmes, genetic evaluations are performed with a BLUP analysis based  
16 479 on the mixed model methodology (Henderson 1975), and genetic covariances are expressed in the pedigree-based  
17 480 relationship matrix (the A matrix). EBV accuracies are, therefore, highly dependent on the correctness of the  
18 481 documented pedigree. In practice, the A matrix does not provide information about all existing relationships,  
19 482 whereas the G matrix can reveal undocumented relationships, pedigree errors, and capture the variation arising  
20 483 from Mendelian sampling (Powell et al. 2010). Both hidden and incorrectly documented relationships may affect  
21 484 the accuracy of genetic parameters, biasing EBV and, by extension, decreasing GS accuracy if EBV (or dEBV)  
22 485 are used as pseudo-phenotypes in GS methodology. We show here that, in Altri's multigenerational breeding  
23 486 programme, the 11.2% parentage errors revealed by G matrix information had a significant impact on EBV for the  
24 487 three traits studied, highlighting the importance of pedigree checking before running GS models with pseudo-  
25 488 phenotypes. In *Pinus pinaster* polycross trials, EBV was shown to be improved by the use of a pedigree  
26 489 reconstructed from parentage analysis and allowing paternal identification (Vidal et al. 2015). In *Populus nigra*,  
27 490 pedigree-based BLUP models based on a corrected A matrix were found to be more accurate than models based  
28 491 on an uncorrected matrix (Pégard et al. 2020). The removal of pedigree errors had a negligible effect on additive  
29 492 genetic variance structure in *Picea rubens*, but the authors suggested that the population studied may have been  
30 493 too small for the assessment of variance components or that the magnitude of error was too small (Doerksen and  
31 494 Herbinger 2010). As Altri's breeding population remains largely ungenotyped, we can hypothesise that additional  
32 495 yet to be identified parentage errors may still bias *E. globulus* genetic estimates.

33 496 We took the incompleteness of genotypic data and the heterogeneity of phenotyping due to the use of data  
34 497 from many trials of different ages into account by using GS models based on pseudo-phenotypes as done in several  
35 498 reports on forest tree species (Resende et al. 2012; Bartholomé et al. 2016; Isik et al. 2016; Thistlethwaite et al.  
36 499 2019). The pseudo-phenotypes came here from a meta-analysis of 31 trials located in the same breeding zone  
37 500 explaining that GxE interaction was not included in the model. When breeding value is used as a pseudo-  
38 501 phenotype, deregression can improve GS accuracy by reducing estimate shrinkage toward parental means, and by  
39 502 taking into account the heterogeneity of EBV reliability (Garrick et al. 2009). This method generates contrasting



503 results, according to the deregression process used. In maritime pine, the use of either EBV or dEBV had no effect  
1 504 on GS accuracy (Isik et al. 2016), whereas, in Douglas fir, the use of dEBV taking mean parental effect into account  
2  
3 505 resulted in a much lower accuracy (Thistlethwaite et al. 2019). Despite these conflicting results, many GS studies  
4 506 in forest trees have used dEBV rather than EBV as the pseudo-phenotype in order to take into account phenotypic  
5  
6 507 information from ungenotyped individuals, and is suitable in case of unbalanced data. For the three traits studied  
7 508 here, GS accuracy was estimated with dEBV, which was strongly correlated with EBV ( $>0.99$ ). This strong  
8  
9 509 correlation may be due to the use of a deregression process without the removal of parental average effect (as  
10 510 many crosses were of unknown paternity), as well as, to the high degree of relatedness in the  $POP_{GS}$  and the large  
11 511 number of clonal copies, both of which contributed to the high accuracy of EBV regardless of the generation  
12 512 considered. For  $dEBV_I$  and  $dEBV_C$ , pedigree correction, principally applied to G1 genotypes, did not change GS  
13 513 accuracy (SV predictions in the S0 scenario) or increased it up to 39% for DBH in the  $S1_a$  scenario. Munoz et al.  
14  
15 514 (2014) reported higher predictive abilities (from 2% to 5%) for the use of  $dEBV_C$  for various traits related to  
16 515 growth and tree architecture in loblolly pine. GS accuracy is commonly evaluated through the random allocation  
17  
18 516 of individuals to either the validation or calibration population. With the S0 scenario encompassing the three  
19 517 generations indifferently, HT, DBH and SV GS accuracies (0.46, 0.60 and 0.48, respectively) were all consistent  
20  
21 518 with published values in *E. globulus* (Durán et al. 2017; Ballesta et al. 2018). Even if the different scenarios implied  
22 519 a low number of genotypes (from 44 in the  $S1_a$  scenario to 342 in the S3 scenario), the effective size of  $POP_{GS}$   
23 520 estimated from status number  $N_S$  (Lindgren et al. 1996) was 34. This limited effective size as well as the high  
24 521 marker density (13.3 SNPs / cM in average) were both parameters impacting favourably accuracies of GS  
25 522 (Grattapaglia and Resende 2011). As suggested a study in maritime pine with  $N_S=25$  (Bartholomé et al. 2016), GS  
26 523 accuracy resulted more from the high relatedness between the CP and the VP than from historical LD associations  
27 524 between markers and QTLs.

525 GS accuracy must be evaluated over generations to determine the value of GS for advanced breeding  
34 526 programmes. Through successive generations of a breeding process, genetic recombination between haplotypes  
35 527 may change the extent of LD, thereby limiting the efficacy of GS over several generations, as LD must be  
36 528 conserved between the CP and VP (Hayes et al. 2009a). The impact of such changes in LD over generations has  
37  
38 529 mostly been investigated in simulation studies. For *Eucalyptus*, Denis and Bouvet (2013) found that GS accuracy  
39 530 decreased over successive breeding cycles for a breeding population with an effective size of  $N_e=100$ . As a means  
40 531 of coping with the loss of GS predictive ability across generations, predictive models should be refreshed by  
41 532 aggregating data from the two most recent breeding cycles, as in the simulations for oil palm performed by Cros  
42 533 et al. (2018). We investigated GS over generations with different scenarios (S1 and S2) in which the CP consisted  
43 534 of individuals from previous generations (G0, G1), used to make predictions for the most recent generations (G1  
44 535 and G2). We obtained moderate-to-high prediction accuracies, with a value of 0.68 for DBH in the  $S2_a$  scenario,  
45 536 indicating that it was possible to predict G2 genotypes from data for the parents (G1) and grandparents (G0) in the  
46 537 CP. Moreover growth trait prediction remained similar to that for conventional cross-validation (S0), suggesting  
47 538 that the predictive model was not altered over generations. This conclusion is consistent with other reports for a  
48 539 five consecutive progeny set in *Hordeum vulgare* (Sallam et al. 2015) and in *Avena sativa* L. (Asoro et al. 2011).  
49 540 Similarly, Bartholomé et al. (2016) showed, in a study on *Pinus pinaster*, that high accuracies (0.70 for height and  
50 541 0.79 for circumference) could be obtained with only G0 and G1 genotypes for the CP, and G2 genotypes for the  
51 542 VP. In an F1 progeny test on Douglas fir, the accuracy of GS for juvenile height was evaluated at 0.9 in the F2  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

543 generation (Thistlethwaite et al. 2019). Here, we found that accuracy for the prediction of survival was slightly  
1 544 lower for the per-generation scenario (S1 and S2) than for S0. Moreover, scenarios S2<sub>c</sub> and S2<sub>a</sub>, calibrated with  
2  
3 545 dEBV<sub>C-T01</sub>, also provided strongly negative mean accuracies for SV prediction in G2 (-0.47 and -0.56,  
4 546 respectively). We can hypothesise that abiotic and/or biotic constraints affected the first two generations (G0 and  
5  
6 547 G1) differently from the last generation (G2), making it difficult to predict complex traits, such as survival, with a  
7 548 CP and VP encompassing different environmental conditions. Accuracy decrease was also observed for growth  
8  
9 549 traits when pseudo-phenotypes included in CP were poorly estimated (scenarios S2<sub>c</sub> and S2<sub>d</sub> vs. scenarios S2<sub>a</sub> and  
10 550 S2<sub>b</sub>). Interestingly, the addition of progeny genotypes (scenarios S1<sub>b</sub>, S3<sub>a</sub> and S3<sub>b</sub>) improved the prediction of both  
11 551 G1 and G2 genotypes for all the traits studied, to 0.65 for HT and 0.78 for DBH in the S3<sub>b</sub> scenario, and for SV,  
12 552 with an optimum of 0.59. This result was expected, as higher levels of relatedness between calibration and  
13 553 validation populations has been shown to improve GS accuracy in other species, such as *Picea glauca* (Beaulieu  
14  
15 554 et al. 2014).  
16  
17  
18 555

## 19 556 **Conclusion**

20 557  
21 558 In conclusion, we report here encouraging results for applied GS in *Eucalyptus globulus*. Given the relatively  
22 559 small population size, we were able to predict the breeding value of the most recent generation reasonably  
23 560 accurately, by aggregating data from the first two generations. In addition, pedigree correction for identity and  
24 561 parentage errors increased the accuracy of GS for all traits. Including a few relatives from targeted families in GS  
25 562 models also improved accuracy for all traits. Further investigations in *E. globulus* are required, particularly as  
26 563 concerning optimisation of the calibration and validation populations, as proposed in GS approaches for other  
27 564 species (Ahmadi and Bartholomé 2022). As POP<sub>GS</sub> was not representative from the breeding population of Altri  
28 565 Florestal (POP<sub>TOT</sub>), this study must be considered as a proof-of-concept. Genotyping efforts in this breeding  
29 566 population will need to continue before implementing concretely GS. In addition, the genotyping of the base  
30 567 population could be compared with the 13 races and eight genetic groups defined in previous studies (Dutkowski  
31 568 and Potts 1999; Costa et al. 2017) to define meta-founders usable for GS (Legarra et al. 2015). An alternative to  
32 569 deregressed EBV for taking the performance of non-genotyped individuals into account would be so-called  
33 570 “single-step genomic BLUP” (Legarra et al. 2014). This methodology has recently been successfully tested in *E.*  
34 571 *globulus* (Callister et al. 2021; Quezada et al. 2022) as the training population size can be increased by including  
35 572 both genotyped and not genotyped individuals.  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

574 **Funding :** This work was funded by CIRAD (Centre de coopération internationale en recherche agronomique pour  
1 575 le développement) and by the European Union’s Horizon 2020 research and innovation programme under grant  
2 576 agreement n°773383 (B4EST). In this context, G. Haristoy received a doctoral fellowship.  
3 577

4 577  
5 578 **Acknowledgements:** Thanks to J. Ventura (Altri Florestal) who carried out work related to phenotyping, leaves  
6 579 sampling and handling, as well as, sending all samples to France. DNA extractions were performed at the  
7 580 BIOGECO research unit (Cestas, France) with the valuable advice of C. Lalanne (INRAE) and C. Larue (INRAE).  
8 581 We also thank H. San Clemente (CNRS), for the additional bioinformatic analysis during the reviewing process,  
9 582 revealing putative genes with polymorphic SNP. DNA quantification was performed at the Genome Transcriptome  
10 583 Facility of Bordeaux (grants from the Conseil Regional d’Aquitaine: n°20030304002FA and 20040305003FA, the  
11 584 European Union: FEDER n°2003227 and ANR: n°ANR-10-EQPX-16 Xyloforest). Genotyping was performed by  
12 585 Thermo Fisher Scientific (Santa Clara, USA). We are also grateful to J. Bartholomé (CIRAD), D. Cros (CIRAD),  
13 586 C. Dubos (INRAE), D. Lopez (CIRAD), P. Rozenberg (INRAE), and D. This (INRAE, L’Institut Agro) for helpful  
14 587 discussions.  
15 588

16 588  
17 589 **Conflicts of interest:** The authors have no conflict of interest to declare.  
18 590  
19 591

20 591 **Authors’ contributions:** Conceptualization and PhD supervision: LB and JMG. Management of the breeding  
21 592 programme and curation of field data: LF, LL, JPP. Laboratory work: GH. Data analysis: LB, JMG, GH. Drafting  
22 593 the manuscript: LB, JMG, GH. Critical revision of the manuscript: LF, LL, JAPP, JPP. All authors have approved  
23 594 the publication of the final version of the manuscript.  
24 595

25 595  
26 596 **Data archiving statement:** Genomic, pedigree, and pseudo-phenotypic data will be deposited to the Data INRAE  
27 597 portal: <https://data.inrae.fr/>. The supplemental data “Marker\_data.csv” contains the molecular profiles of the 401  
28 598 genotypes of POP<sub>GS</sub> (individuals in row, SNP in columns). The pedigree of POP<sub>GS</sub> genotypes is available in the  
29 599 supplemental data file “Pedigree\_data.csv”. The first column contains the identifier of each individual, the second  
30 600 and third columns refer to the mother and father documented in the initial pedigree, and the fourth and fifth columns  
31 601 contain the mother and father after pedigree corrections. The pseudo-phenotypic data are available in 12 different  
32 602 “.csv” files (4 files for each trait), containing either dEBV<sub>I</sub>, dEBV<sub>C</sub>, dEBV<sub>C-T01</sub>, dEBV<sub>C-T2</sub> (refers to the file name).  
33 603 The first column is the identifier of the individuals, and the second column contains the deregressed estimated  
34 604 breeding value.  
35 605

606 **References**

- 1 607 Adams WT, Neale DB, Loopstra CA (1988) Verifying controlled crosses in conifer tree-improvement programs.  
2 608 *Silvae Genet* 37:147–152  
3
- 4 609 Ahmadi N, Bartholomé J (eds) (2022) *Complex Trait Prediction: Methods and Protocols*. Springer US, New York  
5
- 6 610 Ahmar S, Ballesta P, Ali M, Mora-Poblete F (2021) Achievements and Challenges of Genomics-Assisted Breeding  
7 611 in Forest Trees: From Marker-Assisted Selection to Genome Editing. *Int J Mol Sci* 22:10583.  
8 612 <https://doi.org/10.3390/ijms221910583>  
9
- 10 613 Amadeu RR, Cellon C, Olmstead JW, et al (2016) AGHmatrix: R Package to Construct Relationship Matrices for  
11 614 Autotetraploid and Diploid Species: A Blueberry Example. *Plant Genome* 9:1.  
12 615 <https://doi.org/10.3835/plantgenome2016.01.0009>  
13
- 14 616 Asoro FG, Newell MA, Beavis WD, et al (2011) Accuracy and Training Population Design for Genomic Selection  
15 617 on Quantitative Traits in Elite North American Oats. *Plant Genome* 4:132.  
16 618 <https://doi.org/10.3835/plantgenome2011.02.0007>  
17
- 18 619 Ballesta P, Serra N, Guerra F, et al (2018) Genomic Prediction of Growth and Stem Quality Traits in *Eucalyptus*  
20 620 *globulus* Labill. at Its Southernmost Distribution Limit in Chile. *Forests* 9:779.  
21 621 <https://doi.org/10.3390/f9120779>  
22
- 23 622 Bartholomé J, Mabiala A, Burlett R, et al (2020) The pulse of the tree is under genetic control: eucalyptus as a  
24 623 case study. *Plant J* 103:338–356. <https://doi.org/10.1111/tjp.14734>  
25
- 26 624 Bartholomé J, Mandrou E, Mabiala A, et al (2015) High- resolution genetic maps of *Eucalyptus* improve  
27 625 *Eucalyptus grandis* genome assembly. *New Phytol* 206:1283–1296. <https://doi.org/10.1111/nph.13150>  
28
- 29 626 Bartholomé J, Salmon F, Vigneron P, et al (2013) Plasticity of primary and secondary growth dynamics in  
30 627 *Eucalyptus* hybrids: a quantitative genetics and QTL mapping perspective. *BMC Plant Biol* 13:120.  
31 628 <https://doi.org/10.1186/1471-2229-13-120>  
32
- 33 629 Bartholomé J, Van Heerwaarden J, Isik F, et al (2016) Performance of genomic prediction within and across  
34 630 generations in maritime pine. *BMC Genom* 17:604. <https://doi.org/10.1186/s12864-016-2879-8>  
35
- 36 631 Beaulieu J, Doerksen T, Clément S, et al (2014) Accuracy of genomic selection models in a large population of  
37 632 open-pollinated families in white spruce. *Heredity* 113:343–352. <https://doi.org/10.1038/hdy.2014.36>  
38
- 39 633 Blouin MS (2003) DNA-based methods for pedigree reconstruction and kinship analysis in natural populations.  
40 634 *Trends Ecol Evol* 18:503–511. [https://doi.org/10.1016/S0169-5347\(03\)00225-8](https://doi.org/10.1016/S0169-5347(03)00225-8)  
41
- 42 635 Borralho NMG (1995) The impact of individual tree mixed models (BLUP) in tree breeding. In: BM Potts, NMG  
43 636 Borralho, JB Reid, RN Cromer, WN Tibbits and CA Raymond (eds) *Proc CRCTHF-IUFRO Conference*  
44 637 *Eucalypt Plantation: Improving Fibre Yield And Quality*, Hobart, Tasmania, pp 141-145  
45
- 46 638 Borralho NMG, Almeida MH, Potts BM (2007) O melhoramento do eucalipto em Portugal. In: AM Alves and JS  
47 639 Pereira and JMN Silva (eds), *Eucaliptal em Portugal: Impactes Ambientais e Investigação Científica*,  
48 640 ISAPress, Lisbon, Portugal, pp 61-110  
49
- 50 641 Borralho NMG, Pina JP, Leal L, Araujo J (2018) The gain achieved from *Eucalyptus globulus* tree improvement  
51 642 programs in Portugal, a joint analysis of RAIZ and ALTRI trials. In: *Proc Tecnicelpa XXIV Internacional*  
52 643 *Forest, Pulp and Paper Conference*, Porto, Portugal  
53
- 54 644 Butler DG, Cullis BR, Gilmour AR, et al (2017) ASReml estimates variance components under a general linear.  
55 645 VSN International Ltd  
56
- 57 646 Callister AN, Bradshaw BP, Elms S, et al (2021) Single-step genomic BLUP enables joint analysis of disconnected  
58 647 breeding programs: an example with *Eucalyptus globulus* Labill. *G3: Genes, Genomes, Genet*  
59 648 11:jkab253. <https://doi.org/10.1093/g3journal/jkab253>  
60

- 649 Chambers PGS, Borralho NMG, Potts BM (1996) Genetic Analysis of Survival in *Eucalyptus globulus* ssp.  
650 *globulus*. *Silvae Genet* 45:107–112
- 651 Costa e Silva F, Shvaleva A, Broetto F, et al (2008) Acclimation to short-term low temperatures in two *Eucalyptus*  
652 *globulus* clones with contrasting drought resistance. *Tree Physiol* 29:77–86.  
653 <https://doi.org/10.1093/treephys/tpn002>
- 654 Costa J, Vaillancourt RE, Steane DA, et al (2017) Microsatellite analysis of population structure in *Eucalyptus*  
655 *globulus*. *Genome* 60:770–777. <https://doi.org/10.1139/gen-2016-0218>
- 656 Cros D, Tchounke B, Nkague-Nkamba L (2018) Training genomic selection models across several breeding cycles  
657 increases genetic gain in oil palm in silico study. *Molecular Breeding* 38:89.  
658 <https://doi.org/10.1007/s11032-018-0850-x>
- 659 Denis M, Bouvet J-M (2013) Efficiency of genomic selection with models including dominance effect in the  
660 context of *Eucalyptus* breeding. *Tree Genet Genomes* 9:37–51. [https://doi.org/10.1007/s11295-012-](https://doi.org/10.1007/s11295-012-0528-1)  
661 0528-1
- 662 Doerksen TK, Herbinger CM (2008) Male reproductive success and pedigree error in red spruce open-pollinated  
663 and polycross mating systems. *Can J For Res* 38:1742–1749. <https://doi.org/10.1139/X08-025>
- 664 Doerksen TK, Herbinger CM (2010) Impact of reconstructed pedigrees on progeny-test breeding values in red  
665 spruce. *Tree Genet Genomes* 6:591–600. <https://doi.org/10.1007/s11295-010-0274-1>
- 666 Doyle J (1991) DNA Protocols for Plants. In: Hewitt GM, Johnston AWB, Young JPW (eds) *Molecular*  
667 *Techniques in Taxonomy*. Springer Berlin Heidelberg, Berlin, Heidelberg, pp 283–293
- 668 Durán R, Isik F, Zapata-Valenzuela J, et al (2017) Genomic predictions of breeding values in a cloned *Eucalyptus*  
669 *globulus* population in Chile. *Tree Genet Genomes* 13:74. <https://doi.org/10.1007/s11295-017-1158-4>
- 670 Dutkowski GW, Potts BM (1999) Geographic Patterns of Genetic Variation in *Eucalyptus globulus* ssp. *globulus*  
671 and a Revised Racial Classification. *Aust J Bot* 47:237. <https://doi.org/10.1071/BT97114>
- 672 Ericsson (1999) The effect of pedigree error by misidentification of individual trees on genetic evaluation of a full-  
673 sib experiment. *Silvae Genet* 48:239-242.
- 674 Freeman JS, Potts BM, Downes GM, et al (2013) Stability of quantitative trait loci for growth and wood properties  
675 across multiple pedigrees and environments in *Eucalyptus globulus*. *New Phytol* 198:1121–1134.  
676 <https://doi.org/10.1111/nph.12237>
- 677 Gamal El-Dien O, Ratcliffe B, Klápště J, et al (2016) Implementation of the Realized Genomic Relationship Matrix  
678 to Open-Pollinated White Spruce Family Testing for Disentangling Additive from Nonadditive Genetic  
679 Effects. *G3: Genes, Genomes, Genet* 6:743–753. <https://doi.org/10.1534/g3.115.025957>
- 680 Garrick DJ, Taylor JF, Fernando RL (2009) Deregressing estimated breeding values and weighting information  
681 for genomic regression analyses. *Genet, Sel, Evol* 41:55. <https://doi.org/10.1186/1297-9686-41-55>
- 682 Grattapaglia D (2017) Status and Perspectives of Genomic Selection in Forest Tree Breeding. In: Varshney RK,  
683 Roorkiwal M, Sorrells ME (eds) *Genomic Selection for Crop Improvement*. Springer International  
684 Publishing, Cham, pp 199–249
- 685 Grattapaglia D, Resende MDV (2011) Genomic selection in forest tree breeding. *Tree Genet Genomes* 7:241–255.  
686 <https://doi.org/10.1007/s11295-010-0328-4>
- 687 Gu Z, Eils R, Schlesner M (2016) Complex heatmaps reveal patterns and correlations in multidimensional genomic  
688 data. *Bioinformatics* 32:2847–2849. <https://doi.org/10.1093/bioinformatics/btw313>
- 689 Hamilton MG, Acuna M, Wiedemann JC, et al (2015) Genetic control of *Eucalyptus globulus* harvest traits. *Can*  
690 *J For Res* 45:615–624. <https://doi.org/10.1139/cjfr-2014-0428>

- 691 Hamilton MG, Potts BM, Greaves BL, Dutkowski GW (2010) Genetic correlations between pulpwood and solid-  
692 wood selection and objective traits in *Eucalyptus globulus*. *Ann For Sci* 67:511–511.  
693 <https://doi.org/10.1051/forest/2010013>
- 694 Hayes BJ, Bowman PJ, Chamberlain AJ, Goddard ME (2009a) Invited review: Genomic selection in dairy cattle:  
695 Progress and challenges. *J Dairy Sci* 92:433–443. <https://doi.org/10.3168/jds.2008-1646>
- 696 Hayes BJ, Visscher PM, Goddard ME (2009b) Increased accuracy of artificial selection by using the realized  
697 relationship matrix. *Genet Res* 91:47–60. <https://doi.org/10.1017/S0016672308009981>
- 698 Henderson C (1950) Estimation of Genetic Parameters. *Ann Math Stat* 309–310
- 699 Henderson CR (1975) Best Linear Unbiased Estimation and Prediction under a Selection Model. *Biometrics*  
700 31:423. <https://doi.org/10.2307/2529430>
- 701 Hill WG, Weir BS (2011) Variation in actual relationship as a consequence of Mendelian sampling and linkage.  
702 *Genet Res* 93:47–64. <https://doi.org/10.1017/S0016672310000480>
- 703 Hudson CJ, Freeman JS, Kullán AR, et al (2012) A reference linkage map for *Eucalyptus*. *BMC Genom* 13:240.  
704 <https://doi.org/10.1186/1471-2164-13-240>
- 705 ICNF (2019) 6º Inventário Florestal Nacional - Relatório Final. Instituto da Conservação da Natureza e das  
706 Florestas Lisboa, Portugal
- 707 Isik F (2014) Genomic selection in forest tree breeding: the concept and an outlook to the future. *New For* 45:379–  
708 401. <https://doi.org/10.1007/s11056-014-9422-z>
- 709 Isik F, Bartholomé J, Farjat A, et al (2016) Genomic selection in maritime pine. *Plant Sci* 242:108–119.  
710 <https://doi.org/10.1016/j.plantsci.2015.08.006>
- 711 Isik F, Holland J, Maltecca C (2017) Genetic Data Analysis for Plant and Animal Breeding. Springer International  
712 Publishing, Cham
- 713 Jarvis SF, Borralho NMG, Potts BM (1995) Implementation of a Multivariate BLUP Model for Genetic Evaluation  
714 of *Eucalyptus globulus* in Australia. In: BM Potts, NMG Borralho, JB Reid, RN Cromer, WN Tibbits,  
715 and CA Raymond (eds) Proc CRCTHF–IUFRO Conference Eucalypt Plantation: Improving Fibre Yield  
716 And Quality, Hobart, Tasmania, pp 212-216
- 717 Jones TH, Steane DA, Jones RC, et al (2006) Effects of domestication on genetic diversity in *Eucalyptus globulus*.  
718 *For Ecol Manage* 234:78–84. <https://doi.org/10.1016/j.foreco.2006.06.021>
- 719 Keil M, Griffin AR (1994) Use of random amplified polymorphic DNA (RAPD) markers in the discrimination  
720 and verification of genotypes in *Eucalyptus*. *Theor Appl Genet* 89:442–450.  
721 <https://doi.org/10.1007/BF00225379>
- 722 Klápště J, Lstibůrek M, El-Kassaby YA (2014) Estimates of genetic parameters and breeding values from western  
723 larch open-pollinated families using marker-based relationship. *Tree Genet Genomes* 10:241–249.  
724 <https://doi.org/10.1007/s11295-013-0673-1>
- 725 Kumar S, Richardson TE (2005) Inferring relatedness and heritability using molecular markers in radiata pine.  
726 *Mol Breed* 15:55–64. <https://doi.org/10.1007/s11032-004-2059-4>
- 727 Lebedev VG, Lebedeva TN, Chernodubov AI, Shestibratov KA (2020) Genomic Selection for Forest Tree  
728 Improvement: Methods, Achievements and Perspectives. *Forests* 11:1190.  
729 <https://doi.org/10.3390/f11111190>
- 730 Legarra A, Christensen OF, Aguilar I, Misztal I (2014) Single Step, a general approach for genomic selection.  
731 *Livest Sci* 166:54–65. <https://doi.org/10.1016/j.livsci.2014.04.029>

- 732 Legarra A, Christensen OF, Vitezica ZG, et al (2015) Ancestral Relationships Using Metafounders: Finite  
1 733 Ancestral Populations and Across Population Relationships. *Genetics* 200:455–468.  
2 734 <https://doi.org/10.1534/genetics.115.177014>  
3
- 4 735 Lenz PRN, Nadeau S, Azaiez A, et al (2020) Genomic prediction for hastening and improving efficiency of  
5 736 forward selection in conifer polycross mating designs: an example from white spruce. *Heredity* 124:562–  
6 737 578. <https://doi.org/10.1038/s41437-019-0290-3>  
7
- 8 738 Li Y, Dungey HS (2018) Expected benefit of genomic selection over forward selection in conifer breeding and  
9 739 deployment. *PLoS ONE* 13:208–232. <https://doi.org/10.1371/journal.pone.0208232>  
10
- 11 740 Lindgren D, Gea LD, Jefferson PA (1996) Loss of genetic diversity monitored by status number. *Silvae Genet*  
12 741 45:52–59  
13
- 14 742 Lopez GA, Potts BM, Dutkowski G, et al (2002) Genetic variation and inter-trait correlations in *Eucalyptus*  
15 743 *globulus* base population trials in Argentina. *For Genet* 9:217–231  
16
- 17 744 MAPA (2019) Anuario de Estadística Forestal 2019. Ministerio de Agricultura, Pesca y Alimentación, Madrid,  
18 745 Spain  
19
- 20 746 Meuwissen THE, Hayes BJ, Goddard ME (2001) Prediction of Total Genetic Value Using Genome-Wide Dense  
21 747 Marker Maps. *Genetics* 157:1819–1829. <https://doi.org/10.1093/genetics/157.4.1819>  
22
- 23 748 Mora F, Serra N (2014) Bayesian estimation of genetic parameters for growth, stem straightness, and survival in  
24 749 *Eucalyptus globulus* on an Andean Foothill site. *Tree Genet Genomes* 10:711–719.  
25 750 <https://doi.org/10.1007/s11295-014-0716-2>  
26
- 27 751 Mrode RA (2013) Linear models for the prediction of animal breeding values, 3rd ed. CABI, Boston  
28
- 29 752 Munoz F, Rodriguez LS (2020) breedR: Statistical Methods for Forest Genetic Resources Analysts. R package  
30 753 version 0.12-5  
31
- 32 754 Munoz PR, Resende MFR, Huber DA, et al (2014) Genomic Relationship Matrix for Correcting Pedigree Errors  
33 755 in Breeding Populations: Impact on Genetic Parameters and Genomic Selection Accuracy. *Crop Sci*  
34 756 54:1115–1123. <https://doi.org/10.2135/cropsci2012.12.0673>  
35
- 36 757 Myburg AA, Grattapaglia D, Tuskan GA, et al (2014) The genome of *Eucalyptus grandis*. *Nature* 510:356–362.  
37 758 <https://doi.org/10.1038/nature13308>  
38
- 39 759 Pégard M, Segura V, Muñoz F, et al (2020) Favorable Conditions for Genomic Evaluation to Outperform Classical  
40 760 Pedigree Evaluation Highlighted by a Proof-of-Concept Study in Poplar. *Front Plant Sci* 11:581954.  
41 761 <https://doi.org/10.3389/fpls.2020.581954>  
42 762
- 43 762 Potts BM, Vaillancourt RE, Jordan G, et al (2004) Exploration of the *Eucalyptus globulus* gene pool. In: Borralho  
44 763 NMG, Pereira JS, Marques CMP, Coutinho J, Madeira M, Tomé M (eds) *Eucalyptus* in a changing world  
45 764 Proc IUFRO Conference, Aveiro, Portugal, pp 46–61  
46 764  
47
- 48 765 Powell JE, Visscher PM, Goddard ME (2010) Reconciling the analysis of IBD and IBS in complex trait studies.  
49 766 *Nat Rev Genet* 11:800–805. <https://doi.org/10.1038/nrg2865>  
50
- 51 767 Przybylski P, Kowalczyk J, Odrzykoski I, Matras J (2019) Identifying alien genotypes and their consequences for  
52 768 genetic variation in clonal seed orchards of *Pinus sylvestris* L. *Dendrobiology* 81:40–46.  
53 769 <https://doi.org/10.12657/denbio.081.005>  
54
- 55 770 Quezada M, Aguilar I, Balmelli G (2022) Genomic breeding values’ prediction including populational selfing rate  
56 771 in an open-pollinated *Eucalyptus globulus* breeding population. *Tree Genet Genomes* 18:10.  
57 772 <https://doi.org/10.1007/s11295-021-01534-7>  
58
- 59 773 Raymond CA (2002) Genetics of *Eucalyptus* wood properties. *Ann For Sci* 59:525–531.  
60 774 <https://doi.org/10.1051/forest:2002037>  
61  
62  
63  
64  
65

- 775 Resende MDV, Resende MFR, Sansaloni CP, et al (2012) Genomic selection for growth and wood quality in  
1 776 *Eucalyptus*: capturing the missing heritability and accelerating breeding for complex traits in forest trees.  
2 777 New Phytol 194:116–128. <https://doi.org/10.1111/j.1469-8137.2011.04038.x>  
3  
4 778 Rezende GDSP, de Resende MDV, de Assis TF (2014) *Eucalyptus* Breeding for Clonal Forestry. In: Fenning T  
5 779 (ed) Challenges and Opportunities for the World's Forests in the 21st Century. Springer Netherlands,  
6 780 Dordrecht, pp 393–424  
7  
8 781 Rojas PV (2017) Breeding *Eucalyptus globulus* for lower rainfall sites in the Bío-Bío Region of Chile. Australian  
9 782 Forestry. <https://doi.org/10.1080/00049158.2017.1319260>  
10  
11 783 Rstudio Team (2021) RStudio: Integrated Development for R. RStudio, Inc., Boston  
12  
13 784 Salas M, Nieto V, Perafán L, et al (2014) Genetic parameters and comparison between native and local landraces  
14 785 of *Eucalyptus globulus* Labill. ssp. *globulus* growing in the central highlands of Colombia. Ann For Sci  
15 786 71:405–414. <https://doi.org/10.1007/s13595-013-0342-4>  
16  
17 787 Sallam AH, Endelman JB, Jannink J-L, Smith KP (2015) Assessing Genomic Selection Prediction Accuracy in a  
18 788 Dynamic Barley Breeding Population. Plant Genome 8:plantgenome2014.05.0020.  
19 789 <https://doi.org/10.3835/plantgenome2014.05.0020>  
20  
21 790 Sinnwell JP, Therneau TM, Schaid DJ (2014) The kinship2 R Package for Pedigree Data. Human Heredity 78:91–  
22 791 93. <https://doi.org/10.1159/000363105>  
23  
24 792 Tan B, Grattapaglia D, Martins GS, et al (2017) Evaluating the accuracy of genomic prediction of growth and  
25 793 wood traits in two *Eucalyptus* species and their F1 hybrids. BMC Plant Biol 17:110.  
26 794 <https://doi.org/10.1186/s12870-017-1059-6>  
27  
28 795 Tan B, Grattapaglia D, Wu HX, Ingvarsson PK (2018) Genomic relationships reveal significant dominance effects  
29 796 for growth in hybrid *Eucalyptus*. Plant Science 267:84–93. <https://doi.org/10.1016/j.plantsci.2017.11.011>  
30  
31 797 Thavamanikumar S, Arnold RJ, Luo J, Thumma BR (2020) Genomic Studies Reveal Substantial Dominant Effects  
32 798 and Improved Genomic Predictions in an Open-Pollinated Breeding Population of *Eucalyptus pellita*. G3:  
33 799 Genes, Genomes, Genetics 10:3751–3763. <https://doi.org/10.1534/g3.120.401601>  
34  
35 800 Thistlethwaite FR, Ratcliffe B, Klápště J, et al (2019) Genomic selection of juvenile height across a single-  
36 801 generational gap in Douglas-fir. Heredity 122:848–863. <https://doi.org/10.1038/s41437-018-0172-0>  
37  
38 802 Thumma BR, Joyce KR, Jacobs A (2022) Genomic studies with preselected markers reveal dominance effects  
39 803 influencing growth traits in *Eucalyptus nitens*. G3: Genes, Genomes, Genetics 12:jkab363.  
40 804 <https://doi.org/10.1093/g3journal/jkab363>  
41  
42 805 VanRaden PM (2008) Efficient Methods to Compute Genomic Predictions. Journal of Dairy Science 91:4414–  
43 806 4423. <https://doi.org/10.3168/jds.2007-0980>  
44  
45 807 Vidal M, Plomion C, Harvengt L, et al (2015) Paternity recovery in two maritime pine polycross mating designs  
46 808 and consequences for breeding. Tree Genet Genomes 11:105. [https://doi.org/10.1007/s11295-015-0932-](https://doi.org/10.1007/s11295-015-0932-4)  
47 809 4  
48  
49 810 Volker PW, Dean CA, Tibbits WN, Ravenwook IC (1998) Genetic parameters and gains expected from selection  
50 811 in *Eucalyptus globulus* in Tasmania. Silvae Genet 18–21  
51  
52 812 Wickham H (2016) ggplot2. Springer International Publishing, Cham  
53  
54 813  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



814 **Figure captions:**

1 815

2  
3 816 **Fig. 1** Initial pedigree of the 412 POP<sub>TOT</sub> genotypes selected for genotyping (blue dots) (grey dots represent  
4 817 ancestors not available from clonal archives and lines indicate parent-progeny relationships documented in the  
5  
6 818 initial pedigree)

7 819

8  
9 820 **Fig. 2** Heatmap of the A coefficients (under the red diagonal) vs. G coefficients (above the diagonal) of 401  
10 821 individuals from POP<sub>GS</sub>. Two A matrix coefficients are shown: A<sub>I</sub> coefficients without the correction of pedigree  
11 822 errors (a), and A<sub>C</sub> coefficients with pedigree correction (b). In both cases, the 401 genotypes were ordered by  
12 823 generation, from top to bottom, and left to right (G0, G1, G2)

13 824

14  
15 825 **Fig. 3** Distribution of G coefficients within two A classes: 0.25 (a and c) and 0.5 (b and d). The upper two  
16 826 histograms (a and b) show G distributions based on the initial pedigree (A<sub>I</sub> matrix), and the lower two histograms  
17 827 (c and d) show G distributions after pedigree correction (A<sub>C</sub> matrix)

18 828

19 829 **Fig. 4** Distribution of non-concordant SNPs number in the parent-progeny relationships (P/P).

20 830 The threshold (115 non-concordant SNPs) below which the parent of an individual is considered “true” is  
21 831 represented by the red dotted line

22 832

23 833 **Fig. 5** Distribution of the percentage of half-sib and full-sib relationships differing from more 0.2 between A and  
24 834 G matrices for genotypes involved in “undetermined” P/P. The threshold (40%) above which the pedigree was  
25 835 considered inconsistent is represented by the dotted line.

26 836 a) A<sub>I</sub> was considered i.e. initial documented pedigree b) A<sub>C</sub> was considered i.e. pedigree after corrections

27 837

28 838 **Fig. 6** Genetic correlations between EBV<sub>C</sub> for tree height (HT-EBV<sub>C</sub>), diameter at breast height (DBH-EBV<sub>C</sub>) and  
29 839 survival (SV-EBV<sub>C</sub>) across the whole POP<sub>GS</sub> sample (in black) and in the three generations (G0 in red, G1 in  
30 840 green, and G2 in blue) of POP<sub>GS</sub>

31 841

32 842 **Fig. 7** GS accuracy of the S0 scenario for height (HT), diameter at breast height (DBH) and survival (SV), with  
33 843 either dEBV<sub>I</sub> (in grey) or dEBV<sub>C</sub> (in red) used as a pseudo-phenotype (the means are indicated by coloured dots)

34 844

35 845 **Fig. 8** Accuracy of GS models for height (HT), diameter at breast height (DBH) and survival (SV) according to  
36 846 the nine scenarios tested: S0 (in red), S1<sub>a</sub> and S1<sub>b</sub> (green), S2<sub>a</sub>, S2<sub>b</sub>, S2<sub>c</sub>, and S2<sub>d</sub> (blue), and S3<sub>a</sub> and S3<sub>b</sub> (purple).  
37 847 In all scenarios, deregressed and corrected EBV (dEBV<sub>C</sub>) were used as pseudo-phenotypes. Significance is shown  
38 848 for each assessment of accuracy in Table S4)

39 849

40

41

42

43

44

45

46

47

48

850 **Supplementary material captions:**

1 851

2  
3 852 **Table S1** Settings for SNP quality control analysis in Axiom Suite Analysis software

4 853

5  
6 854 **Table S2** Descriptive statistics for  $EBV_C$ ,  $EBV_{C-T01}$ ,  $EBV_{C-T2}$  and their mean accuracy ( $r$ ) for tree height (HT),  
7 855 diameter at breast height (DBH) and survival (SV) for each generation of  $POP_{TOT}$

8  
9 856

10 857 **Table S3** Mean accuracy with  $dEBV_I$  or  $dEBV_C$ , for height (HT), diameter at breast height (DBH), and survival  
11 858 (SV), for scenarios  $S_0$ ,  $S_{1a}$  and  $S_{2a}$ . (1) In  $S_0$ , accuracy is the mean of 100 per-iteration accuracies, and the  
12 859 corresponding significance threshold is that for at least 95% of the 100 iterations

13  
14 860

15 861 **Table S4** Significance of GS accuracies by scenario. In cases of iteration (\*), the accuracy given is the mean of  
16 862 100 per-iteration accuracies, and the corresponding significance threshold is the threshold for at least 95% of the  
17 863 100 iterations. The most globally significant accuracies are shown in bold

18  
19 864

20 865 **Fig. S1** Distribution of G coefficients according to  $A_I$  or  $A_C$  coefficients in  $POP_{GS}$

21  
22 866

23 867 **Fig. S2** Regression of  $EBV_I$  on  $EBV_C$  for tree height (A), diameter at breast height (B), and survival (C) for the  
24 868 401 genotypes of  $POP_{GS}$ . The black dots referred to the 356 genotypes without pedigree correction, and the red  
25 869 dots, the 45 genotypes for which the pedigree was corrected. Pearson correlation of  $EBV_I$  with  $EBV_C$  for the 45  
26 870 genotypes is indicated in red

27  
28 871

29 872 **Fig. S3** Genetic correlations between  $EBV_C$  for tree height (HT- $EBV_C$ ), diameter at breast height (DBH- $EBV_C$ )  
30 873 and survival (SV- $EBV_C$ ) across the whole  $POP_{TOT}$  sample (in black) and in the three generations (G0 in red, G1  
31 874 in green, and G2 in blue) of  $POP_{TOT}$

32  
33 875

34 876 **Fig. S4** Correlation matrix for the  $EBV_C$ ,  $EBV_{C-T01}$ , and  $EBV_{C-T2}$  estimates for the three traits (HT, DBH and SV)  
35 877 based on the  $POP_{GS}$  sample. Correlation coefficients are indicated for the whole  $POP_{GS}$  population (in grey) and  
36 878 for each generation (G0 in red, G1 in green and G2 in blue). The significance threshold is indicated as follows:  
37 879 5% (\*), 1% (\*\*), and 0.1% (\*\*\*)

38  
39

40  
41

42  
43

44  
45

46  
47

48  
49

50  
51

52  
53

54  
55

56  
57

58  
59

60  
61

62  
63

64  
65

Figure 1

[Click here to access/download;Figure;Fig1.jpg](#)

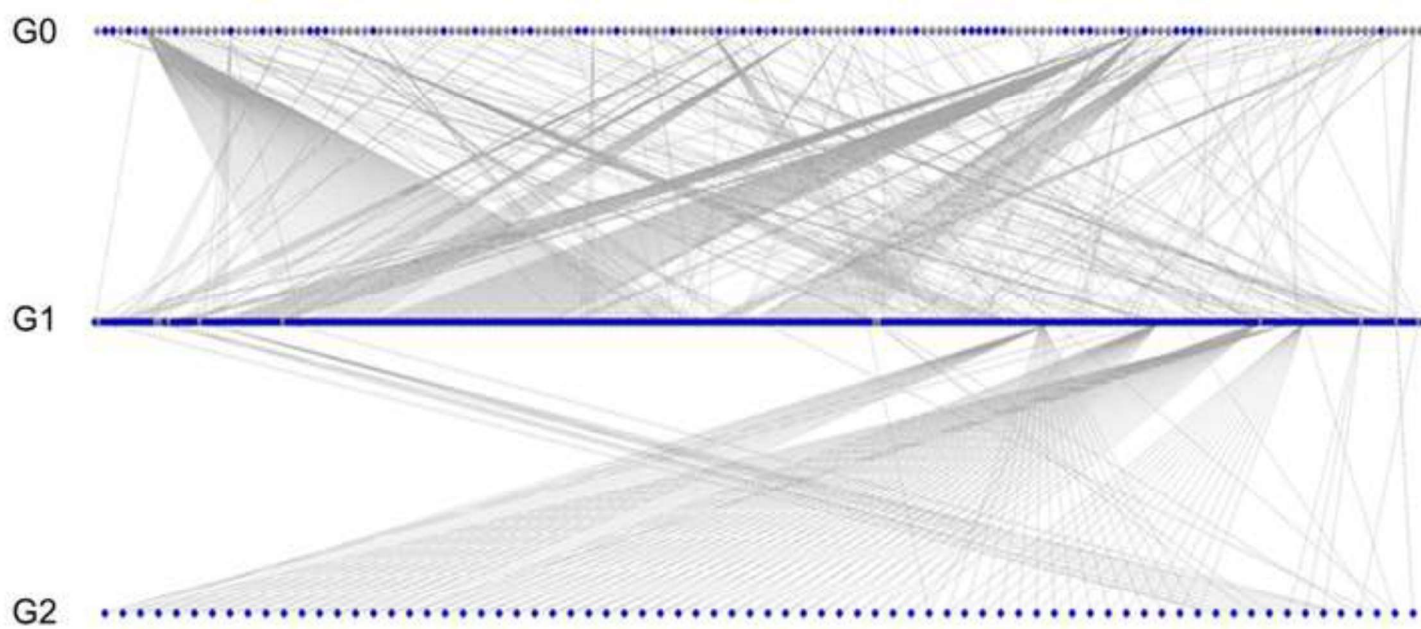
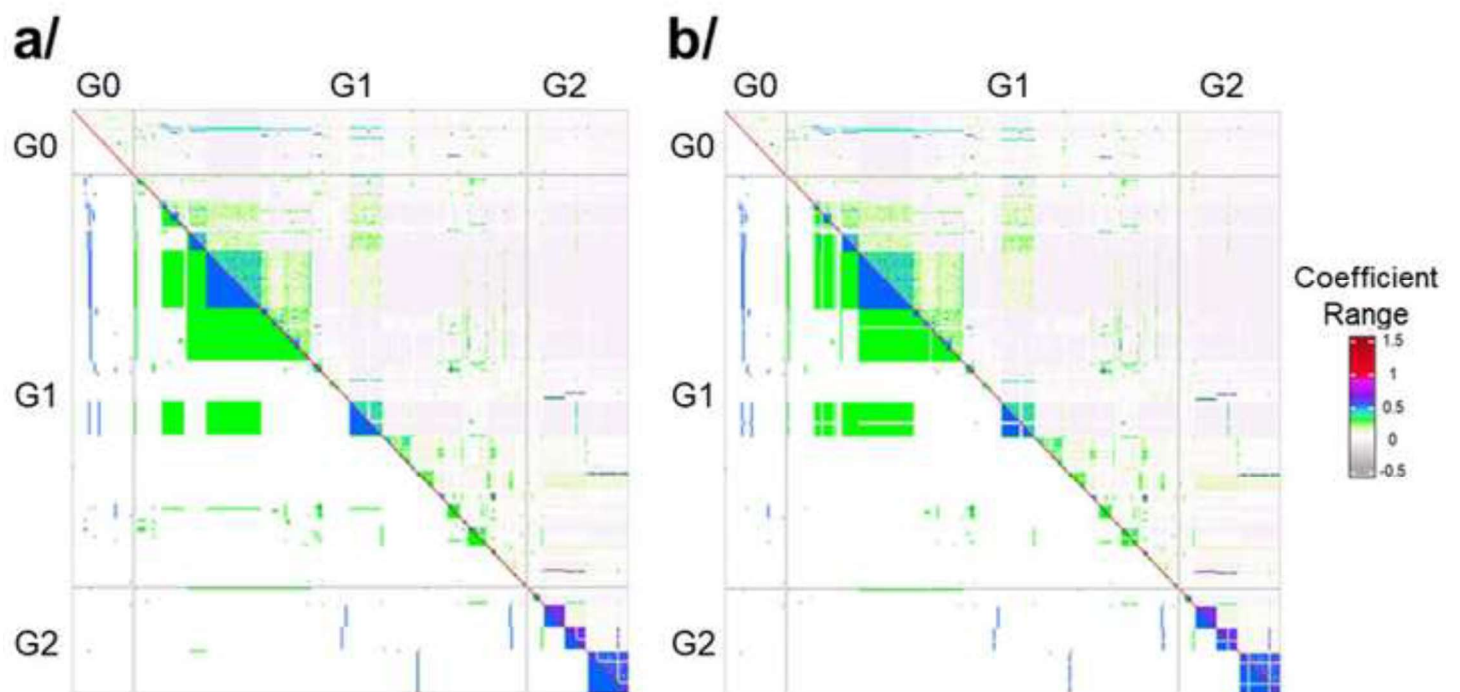


Figure 2

[Click here to access/download;Figure;Fig2.jpg](#)



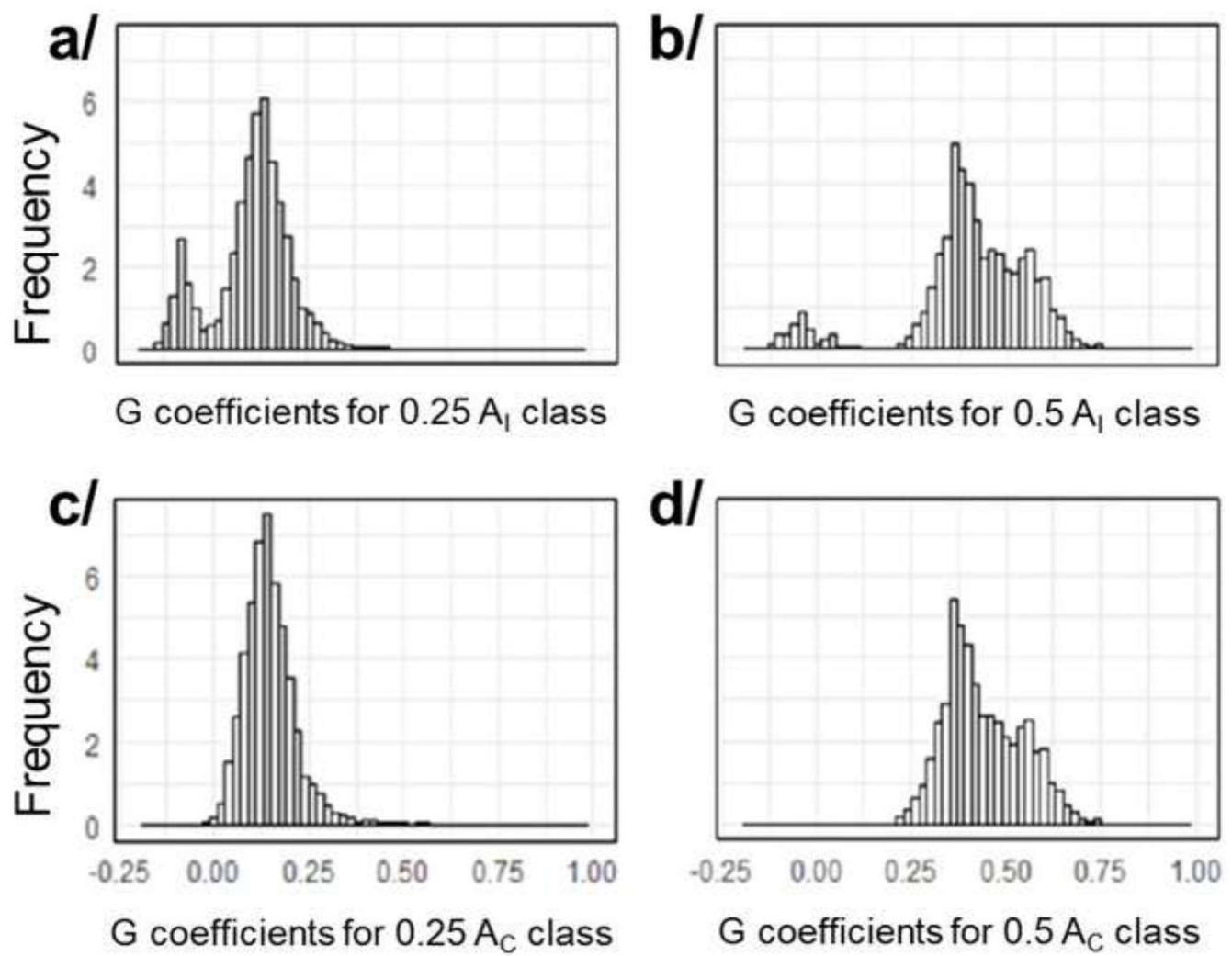


Figure 4

[Click here to access/download;Figure;Fig4.jpg](#)

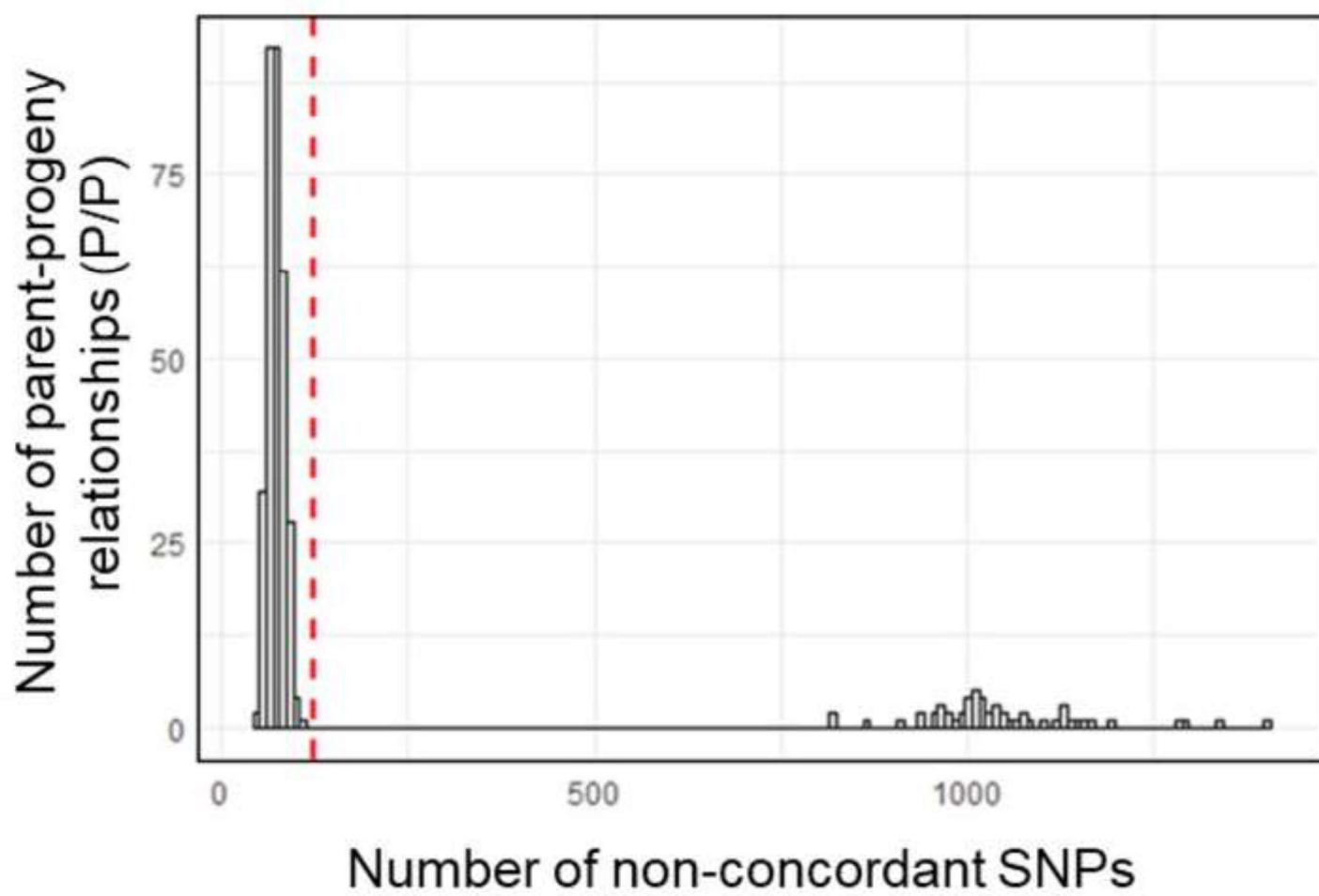


Figure 5

[Click here to access/download;Figure;Fig5.jpg](#)

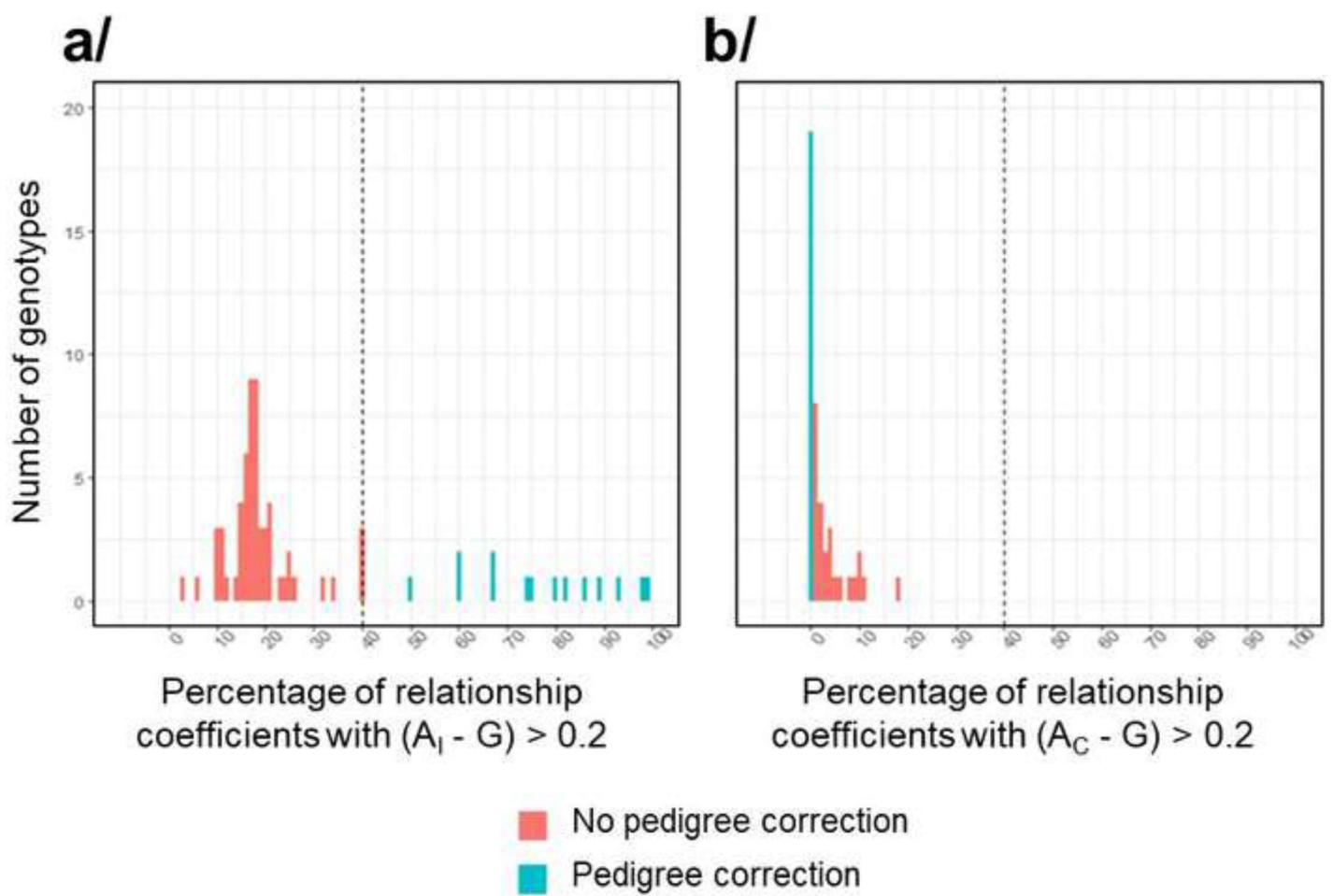


Figure 6

[Click here to access/download;Figure;Fig6.jpg](#)

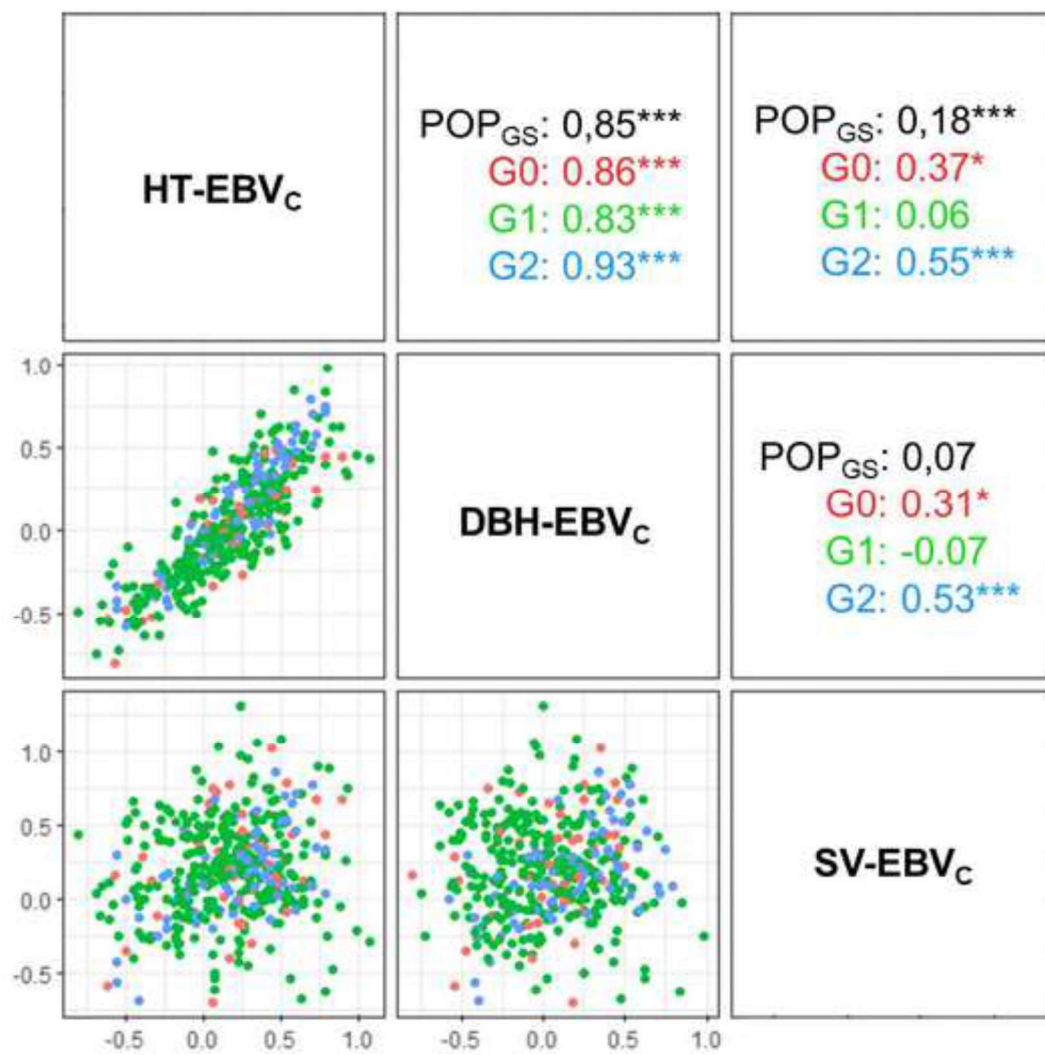




Figure 7

[Click here to access/download;Figure;Fig7.jpg](#)

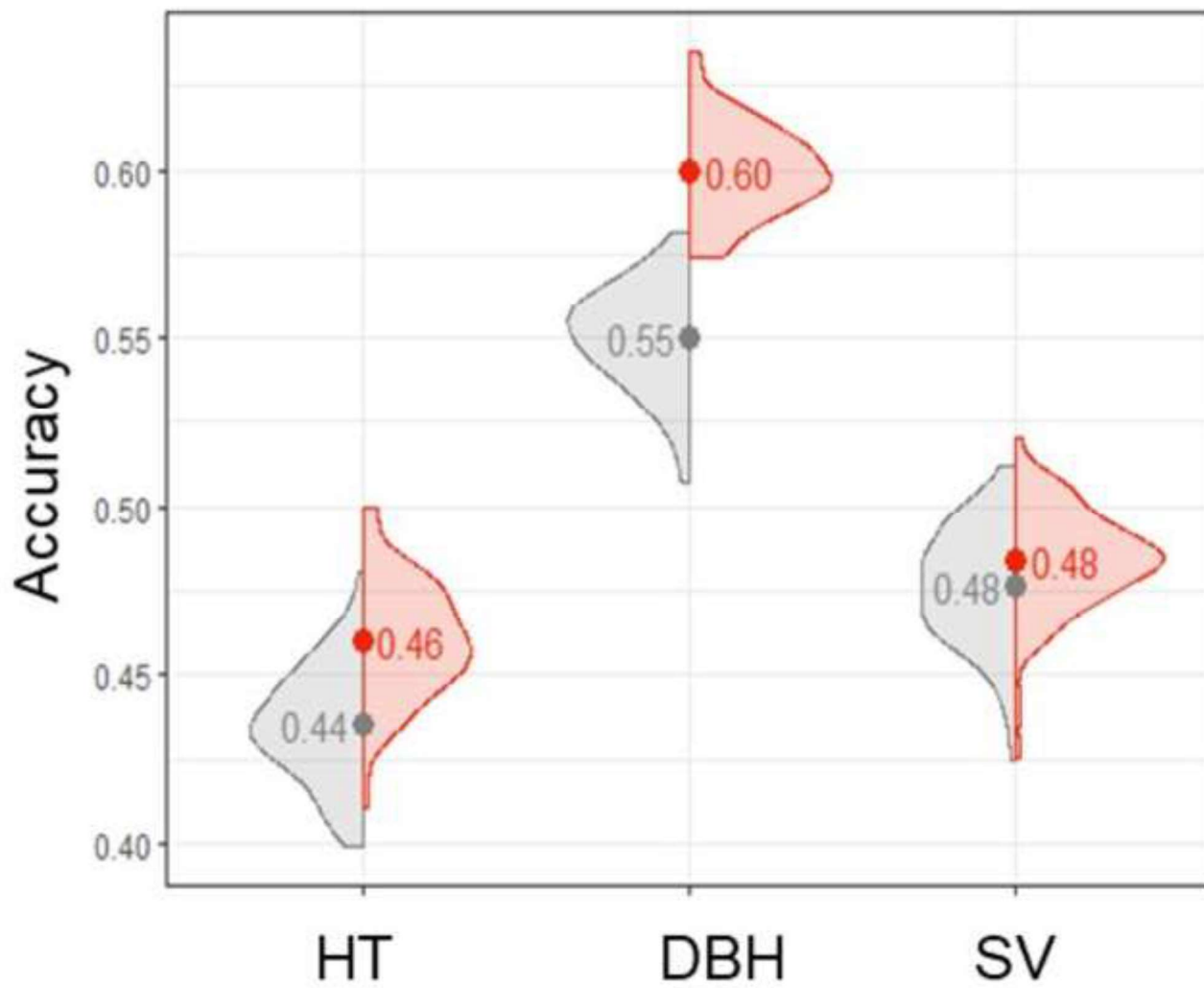


Figure 8

[Click here to access/download;Figure;Fig8.jpg](#)

