

# A simulation study comparing advanced marker-assisted selection with genomic selection in tree breeding programs

Bernd Degen ,\* Niels A. Müller

Thünen Institute of Forest Genetics, Sieker Landstrasse 2, 22927, Grosshansdorf, Schleswig-Holstein, Germany

\*Corresponding author: Thünen Institute of Forest Genetics, Sieker Landstrasse 2, 22927, Grosshansdorf, Schleswig-Holstein, Germany. Email: bernd.degen@thuenen.de

## Abstract

Advances in DNA sequencing technologies allow the sequencing of whole genomes of thousands of individuals and provide several million single nucleotide polymorphisms (SNPs) per individual. These data combined with precise and high-throughput phenotyping enable genome-wide association studies (GWAS) and the identification of SNPs underlying traits with complex genetic architectures. The identified causal SNPs and estimated allelic effects could then be used for advanced marker-assisted selection (MAS) in breeding programs. But could such MAS compete with the broadly used genomic selection (GS)? This question is of particular interest for the lengthy tree breeding strategies. Here, with our new software “SNPscan breeder,” we simulated a simple tree breeding program and compared the impact of different selection criteria on genetic gain and inbreeding. Further, we assessed different genetic architectures and different levels of kinship among individuals of the breeding population. Interestingly, apart from progeny testing, GS using gBLUP performed best under almost all simulated scenarios. MAS based on GWAS results outperformed GS only if the allelic effects were estimated in large populations (ca. 10,000 individuals) of unrelated individuals. Notably, GWAS using 3,000 extreme phenotypes performed as good as the use of 10,000 phenotypes. GS increased inbreeding and thus reduced genetic diversity more strongly compared to progeny testing and GWAS-based selection. We discuss the practical implications for tree breeding programs. In conclusion, our analyses further support the potential of GS for forest tree breeding and improvement, although MAS may gain relevance with decreasing sequencing costs in the future.

**Keywords:** gBLUP, genetic gain, genomic selection, GWAS, marker-assisted selection, SNPscan breeder, stochastic simulation, Plant Genetics and Genomics

## Introduction

Due to the advance of next-generation sequencing technologies, the application of large marker sets with thousands of single nucleotide polymorphisms (SNPs) to estimate genomic breeding values [genomic prediction, genomic selection (GS)] has become a common practice in animal and plant breeding (Misztal *et al.* 2020; Sandhu *et al.* 2022). The effectiveness and accuracy of the genomic predictions have been studied with stochastic simulation models (Pérez-Enciso *et al.* 2017) and deterministic models (Grattapaglia and Resende 2011). In contrast to the success of GS, marker-assisted selection (MAS) could not fulfill the expectations in breeding programs (Grattapaglia and Kirst 2008; Kiszonas and Morris 2018) despite 4 decades of research and development (Nadeem *et al.* 2018). The main reason for this failure was the underestimated complexity of the genetic architecture for most traits of interest with many causal SNPs, most of them with small effects. In most studies, the number of causal gene markers was strongly underestimated and their effect sizes were overestimated (Grattapaglia 2022). Initially, the main focus of MAS was on the search for quantitative trait loci (QTLs), defined as regions of the genome associated with a particular phenotypic trait. Most studies used linkage mapping in biparental populations to identify QTLs (Würschum 2012). In the beginning, QTL studies were

performed with a relatively small number (<1000) of genetic markers (RAPDs, AFLPs, and SSRs). The total amount of phenotypic variation explained by the QTLs was small due to the polygenic nature of most studied traits and due to the low proportion of loci segregating for the causal alleles in biparental populations. Further, the application of the QTLs to other unrelated material was a challenge. With the advance of DNA sequencing techniques, many more markers, especially SNPs, have been used to detect QTLs, and the first genome-wide association studies (GWAS) in plants were performed to detect causal SNPs in populations of unrelated individuals (Atwell *et al.* 2010). GWAS identified many more causal loci compared to QTL mapping. However, these studies were mostly done with arrays of only a few thousand SNPs, so again, the SNPs identified in GWAS relied on linkage to the true causal variants and thus were difficult to use for unrelated individuals. Although the proportion of SNPs identified and the proportion of phenotypic variation explained increased a lot, it was still not enough to effectively estimate individual breeding values for a breeding program. As demonstrated in human genetics, the use of whole-genome sequencing data can have a large impact on GWAS results (Wainschtein *et al.* 2022). Causal SNPs could be identified directly without relying on linkage or imputation, and with increasing sample sizes, more and more of the “missing heritability” could be recovered (Yengo *et al.* 2022). These results

Received: May 31, 2023. Accepted: July 19, 2023

© Crown copyright 2023.

This Open Access article contains public sector information licensed under the Open Government Licence v3.0 (<http://www.nationalarchives.gov.uk/doc/open-government-licence/version/3/>).

raise the question of whether an advanced MAS approach, based on such whole-genome GWAS, would perform equal or even better than GS in a breeding program.

MAS and GS are of particular interest to tree breeding because of the difficulty of phenotyping large and often heterogeneous breeding populations and the long duration of breeding cycles. Even for the fast-growing *Eucalyptus* species, a breeding cycle takes at least 8 years (Mphahlele et al. 2020); for most other tree species, a breeding cycle takes decades. The limitation is caused by the late reproductive maturity and long phenotyping phase required for growth-related traits of at least one-third of the rotation period. This is the reason why GS is now “climbing the slope of enlightenment” in tree breeding programs (Grattapaglia 2022). In 2022, there were already 26 published GS studies for forest trees focusing on species of the genera *Eucalyptus*, *Picea*, and *Pinus* (Isik 2022). Usually, the breeding programs of these species include a few dozen to a few hundred individuals as the founder population, which were used to generate thousands of offspring for further selection in the following breeding cycles (Vidal et al. 2017; da Silva et al. 2019; Isik and McKeand 2019; Li et al. 2020). Modern tree breeding programs use SNP arrays with 5,000–150,000 SNPs as sources for genetic markers (Grattapaglia et al. 2011; Kastally et al. 2022; Nantongo et al. 2022). So far, whole-genome data were practically not used in forest tree breeding programs (Grattapaglia 2022). However, such data may enable the identification of causal SNPs also for complex traits as a basis for MAS.

To assess the potential of different breeding strategies, simulation studies can provide useful information (Liu et al. 2019). The simulation approaches in tree breeding used so far did not run on the scale of whole genomes with large numbers of markers (Resende et al. 2017; Grattapaglia 2022). Existing simulation studies explored the potentials of GS in tree breeding (Iwata et al. 2011; Resende et al. 2012; Li and Dungey 2018) but did not investigate the potential of whole-genome sequencing to identify causal SNPs and their use for advanced MAS. Here, we used our new simulation program *SNPscan breeder* to compare the expected performance of advanced MAS, GS, and the traditional selection criteria phenotype and progeny testing (backward selection) in a simple tree breeding program. Further, we investigated in different scenarios the impact of the number of genotyped and phenotyped individuals, the complexity of genetic architecture, the level of kinship among individuals, and the GWAS algorithm on the dynamics of genetic gain and inbreeding in a tree breeding program.

## Methods

### Simulation scenarios

Here, we used our new stochastic simulation model *SNPscan breeder* (Degen and Müller 2023) to simulate a simple breeding program typical for a tree species with closed recurrent selection and separated generations of a monoecious or hermaphroditic, diploid species. The species had 10 chromosomes and a genome of 500 million base pairs (Mbp). In all scenarios, we simulated 1 million SNPs equally distributed over the genome, resulting in 100,000 SNPs per chromosome. For the proportions of bi-allelic, tri-allelic, and tetra-allelic SNPs as well as the distribution of minor allele frequencies (MAFs), we used the default values of *SNPscan breeder*, that is 95% bi-allelic, 4% tri-allelic, and 1% tetra-allelic as well as 40% MAF 0.01–0.05, 15% MAF 0.05–0.1, and 45% MAF 0.1–0.5. These values were based on population resequencing data of different tree species, specifically beech (*Fagus sylvatica*), oak (*Quercus robur*), and ash (*Fraxinus excelsior*). For details on these genomic data, see Sollars et al. (2017), Pfenninger et al. (2021), and

Plomion et al. (2016). The probability for a crossing-over was identical along the chromosomes, leading on average to 1.5 crossing-overs per chromosome.

We assumed a single target trait. The narrow-sense heritability ( $h^2$ ) was 0.5. The phenotype of each individual  $i$  ( $p_i$ ) was computed as the sum of the mean phenotype of the population at the beginning ( $\bar{p}$  = parameter of the model), the genetic value ( $g_i$ ), and an environmental value ( $e_i$ ):  $p_i = \bar{p} + g_i + e_i$ .

The genetic value of individual  $i$  ( $g_i$ ) was calculated as the sum of all additive effects at each causal locus  $j$  ( $a_{ij}$ ) + a correction  $m$  to centralize the mean of all genetic values of the initial population to 0 multiplied with the scale factor (sf):  $g_i = \sum_{j=1}^n a_{ij} + m \times \text{sf}$ .

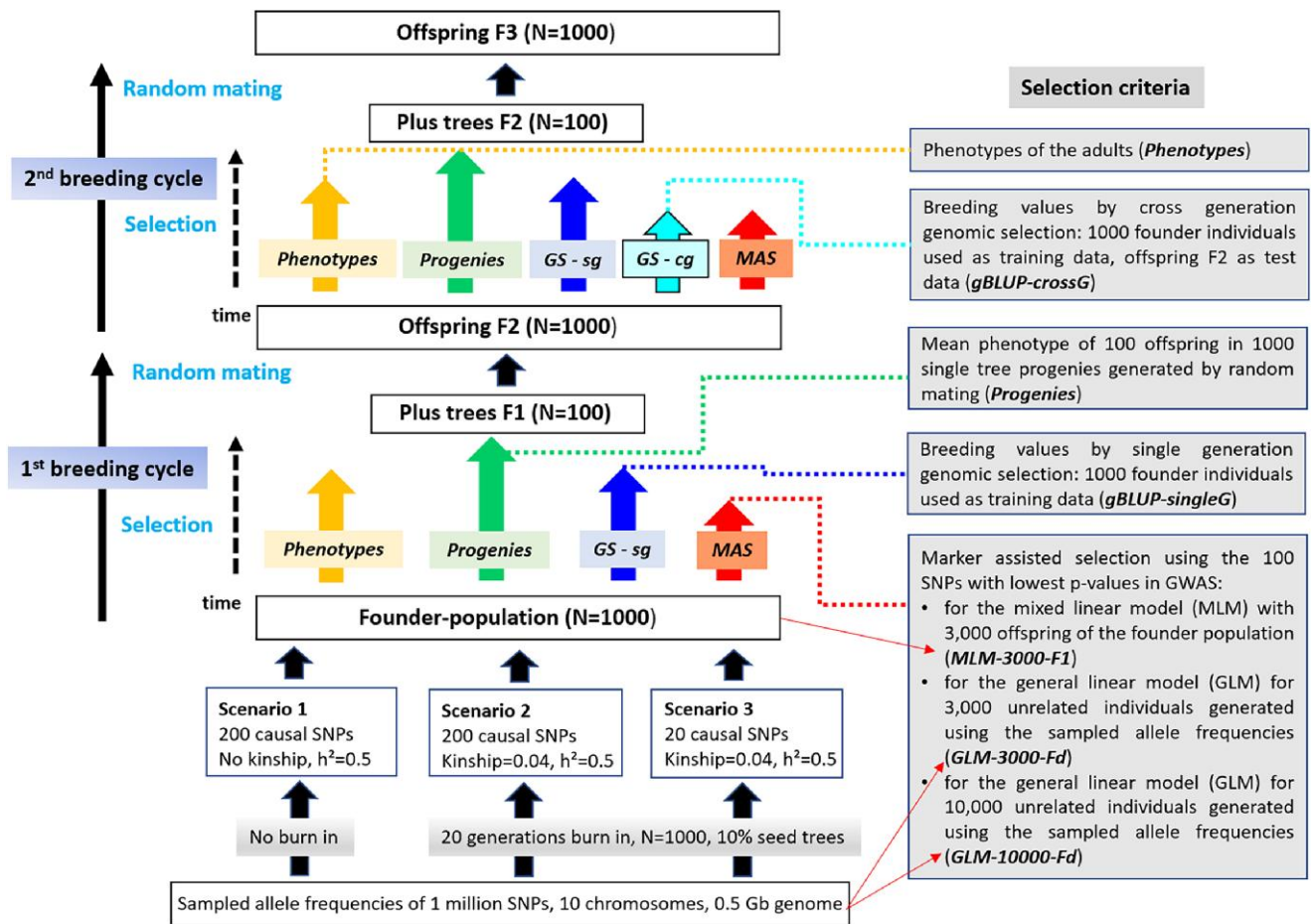
The sf was defined as  $\text{sf} = \sqrt{(\sigma_p^2 * h^2) / \sigma_{g\text{-values}}}$  [ $\sigma_p^2$  = variance of the phenotypes in the population (parameter of the model);  $h^2$  = heritability (0–1, parameter of the model);  $\sigma_{g\text{-values}}$  = standard deviation of the genetic values of the population in the beginning].

Environmental values ( $e_i$ ) were sampled from a normal distribution with a mean of 0 and the standard deviation of environmental effects:  $N(0, \sigma_{\text{env}})$ ,  $\sigma_{\text{env}} = \sqrt{(\sigma_p^2 - h^2 * \sigma_p^2)}$ .

In each generation, 1,000 trees were generated. In scenario 1, we had no burn-in. In scenarios 2 and 3, we ran 20 generations as a burn-in, with 100% of the males contributing to the random mating but only 10% randomly selected trees used for the harvest of seeds for the next generation. This burn-in created a gradually varying level of kinship among the individuals, with an average of 0.04.

In scenarios 1 and 2, the trait was encoded by 200 causal SNPs and in scenario 3 by 20 SNPs. The effects at these SNPs followed a standard normal distribution with a mean of 0 and a standard deviation of 1. There was random mating among all selected parents. In the forward simulations, the top 10% of parents for the next breeding cycle were selected according to the following selection criteria (Fig. 1):

- a) Phenotype of adults => Phenotypes
- b) Breeding values estimated by progeny tests (100 offspring per tree generated by random mating among all 1,000 trees). The mean phenotype of the offspring was used to rank the adults (estimate the genomic breeding values in a backward selection) => Progenies
- c) Single-generation GS with gBLUP using 10,000 randomly selected SNPs (excluding causal SNPs) recalculated in each generation. Here, all individuals provided genotypes and phenotypes to the analysis, and thus, there was no separation between training and test data sets. For this, *SNPscan breeder* used the function “kin.blup” integrated in the R package “rrBLUP” (Endelman 2011; R Core-Team 2022). During the simulations, *SNPscan breeder* created a subdirectory “R” in the project folder in order to store the input files, r-script for the calculation and the results. The genotypes at the 10,000 SNPs were transformed in a  $-1, 0, 1$  matrix for bi-allelic markers. The Euclidian distance among the so-transformed genotypes served as the estimator of the genomic matrix. The algorithm computes “Predicted Genomic Breeding Values (PGBV)” for all individuals in the actual generation  $F$  with the help of their phenotypes and the kinship matrix of all individuals => gBLUP-singleG
- d) Cross-generation GS with gBLUP using 10,000 randomly selected SNPs (excluding causal SNPs) recalculated for each generation with phenotypes of  $F_{n-1}$  and genotypes of  $F_{n-1}$  and  $F_n$ . Here, the kinship matrix and the phenotypes of  $F_{n-1}$  serve as the training set and the next-generation  $F_n$  is the test data set => gBLUP-crossG. The cross-generation



**Fig. 1.** Schematic presentation of the simulated tree breeding program for the first 2 breeding cycles. The breeding cycles 3–5 followed the scheme of the second breeding cycle.

GS started in the second breeding cycle. In the first breeding cycle, it was a single-generation GS since no training data of a former generation existed.

- e) MAS using estimated allelic effects of the 100 SNPs with the lowest  $P$ -values in a GWAS with the mixed linear model (MLM) using the 3,000 phenotypes and genomes of an F1. The individuals of the F1 were generated with random mating of all 1,000 individuals of F0. The simulated genomes of the offspring were stored as a HapMap file, and their phenotypes were stored in a text file. Then, the GWAS were computed with Tassel version 5.0 (Bradbury et al. 2007). During the simulations, the breeding values of individuals are estimated as the sum of the allelic effects at all identified (true + false positive) associated SNPs => MLM-3000-F1
- f) MAS using estimated allelic effects of the 100 SNPs with the lowest  $P$ -values in a GWAS with the general linear model (GLM) using 3,000 unrelated individuals of the founder population => GLM-3000-Fd
- g) MAS using estimated allelic effects of the 100 SNPs with the lowest  $P$ -values in a GWAS (GLM) using 10,000 unrelated individuals of the founder population => GLM-10000-Fd

All scenarios were repeated 10 times.

It should be noted that in practical breeding programs, the application of the different selection criteria results in huge differences in costs, labor requirements, and time needed to complete the selection (see Discussion).

## Simulation outputs

### Genetic gain

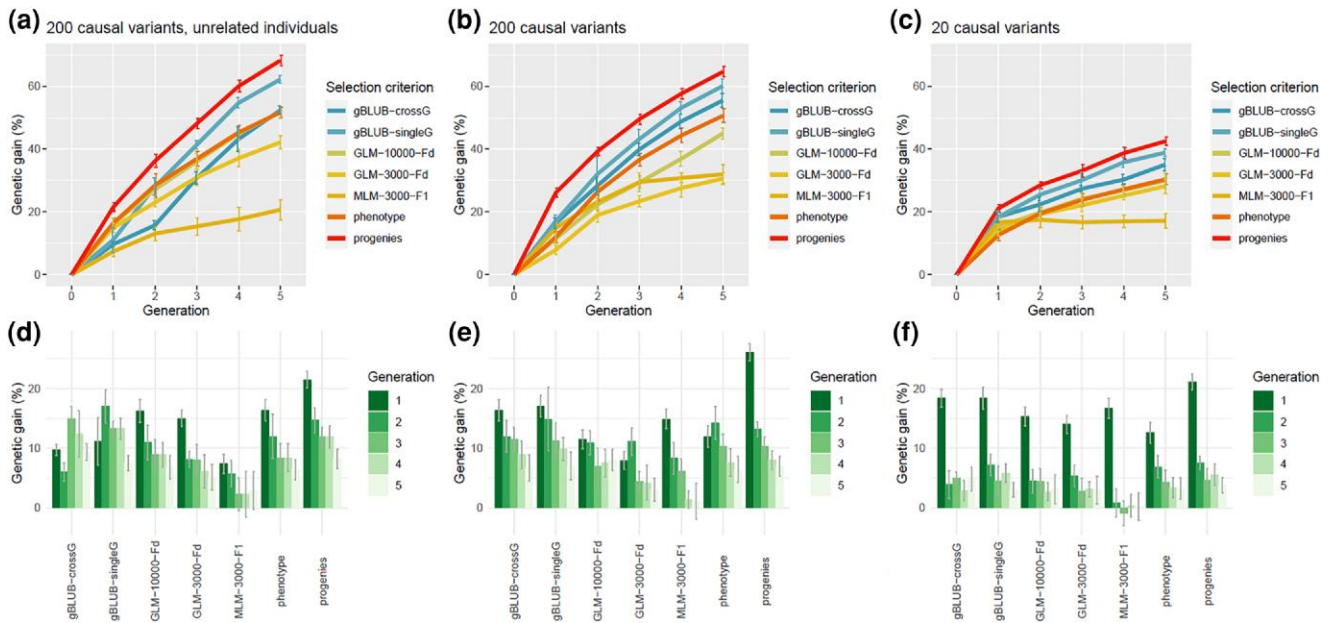
The genetic gain was calculated as the relative difference of the population mean at the target trait between 2 consecutive generations.

### Inbreeding $F$ (0–1)

The inbreeding coefficient  $F$  was computed as the average probability that the 2 alleles of a homozygous genotype are identical by descent. It was realized in the simulations by using 100 “dummy” loci with unique alleles for all individuals in the initial founder population. In consecutive generations, homozygotes at these loci represent alleles that were identical by descent.

### Post hoc analysis of GWAS accuracy

The Tassel (Bradbury et al. 2007) results of GWAS using the GLM and MLM algorithms have been loaded into *SNPscan breeder* for post hoc analysis. For user-defined thresholds of the association probabilities, the proportion of true- and false-positive associated SNPs was analyzed. Here, we selected the threshold so that 100 SNPs are identified as significant. The estimated allelic effects were compared to the true effects using the Pearson correlation coefficient. Further, correlation coefficients between true and estimated individual breeding values were computed considering the identified SNPs and their allelic effects.



**Fig. 2.** Genetic gains achieved by different selection criteria over 5 generations in 3 simulated breeding populations. a–c) Average cumulative genetic gains ( $n = 10$ ) over 5 breeding cycles are shown for selection based on gBLUP, i.e. GS (blue), GWAS, i.e. MAS (green and yellow), phenotypes (orange), and progeny tests (red). d–f) Average genetic gains ( $n = 10$ ) per breeding cycle (generation) are shown for the different selection criteria. Error bars represent the standard deviation of the mean. The 3 simulated breeding populations exhibit different genetic architectures, with 200 a, b, d, e) vs 20 c, f) causal variants, and different levels of kinship, with unrelated individuals a, d) vs average kinship values of approximately 0.04 b, c, e, f).

## Results

### Genetic gain

The selection based on progeny tests performed best in all tested scenarios, irrespective of the number of causal variants (20 vs 200) or the level of kinship (0 vs 0.04). The progeny tests delivered high genetic gains of more than 20% in the first breeding cycle and cumulative genetic gains between 40 and 70% at the end of 5 cycles (Fig. 2a–f). In both scenarios with a burn-in phase and a resulting kinship structure, GS (gBLUP-singleG and gBLUP-crossG) outperformed the selection by phenotypes. As expected, the single-generation GS (gBLUP-singleG) delivered slightly better results compared to the cross-generation GS (gBLUP-crossG). Overall, MAS based on different GWAS analyses resulted in the lowest genetic gains. Irrespective of the genetic architecture, MAS could not compete even against simple selection by phenotypes and was clearly outperformed by GS. The only exception where MAS achieved higher genetic gains than GS was the first breeding cycle in a population of unrelated individuals; however, here, the selection based on phenotypes was on a similar level. When the target trait is controlled by a small number of causal variants (such as 20 SNPs in scenario 3), MAS can lead to rapid fixation of the identified alleles in the breeding population and prevent further genetic gains. This is illustrated by using allelic effects estimated by a GWAS with 3,000 individuals in an F1 (MLM-3000-F1) where no additional genetic gain was realized after breeding cycle 1.

### Precision of GWAS

In each scenario, we used MAS with allelic effects estimated in 3 different GWAS subscenarios. First, we ran a MLM using the genomic data and kinship information for 3,000 individuals created in an F1 (MLM-3000-F1). In addition, we estimated allelic effects using 2 GLM analyses with 3,000 (GLM-3000-Fd) and 10,000 (GLM-10000-Fd) unrelated individuals from the founder population before the 20 generations of burn-in. In order to estimate

the potential for a GWAS using extreme phenotypes, we selected from 10,000 simulated phenotypes the largest 1,500 and smallest 1,500 and repeated the GLM (GLM-3000-extreme). The probability thresholds were set so that for each of the 9 GWAS, the top 100 SNPs were selected (Table 1). Generally, the use of unrelated individuals in the GWAS outperformed the use of related individuals in terms of the proportion of true-positive SNPs and cumulative genetic gains after the 5 breeding cycles (Table 1). For the scenarios with 200 causal SNPs, the increase in the number of individuals in the GWAS led to a higher number of identified SNPs, more precise estimates of the allelic effects, and higher genetic gains. Interestingly, nearly the same proportions of identified true-positive SNPs and genetic gains were realized with a GWAS using only the 3,000 extreme phenotypes (GLM-3000-extreme). The correlation between the estimated and true allelic effects varied between 0.60 and 0.96 and exceeded 0.9 in all subscenarios using 10,000 individuals or extreme phenotypes (“r allele effects” in Table 1).

### Inbreeding

In all tested scenarios, the inbreeding increased in each breeding cycle (Fig. 3a–c). However, the increase was different depending on the selection criteria. In scenario 1 with 200 causal variants and no burn-in, inbreeding values between 0.051 and 0.117 were observed in the fifth breeding cycle with an average  $\Delta F$  per breeding cycle of 0.010–0.023 (Fig. 3a). In scenarios 2 and 3 with 200 and 20 causal variants and a burn-in, inbreeding reached values between 0.072 and 0.122 ( $\Delta F$  per breeding cycle of 0.007–0.017) in scenario 2 and values between 0.065 and 0.117 in scenario 3 ( $\Delta F$  per breeding cycle of 0.006–0.017). The large genetic gains of the GS are linked to a stronger increase of inbreeding (Fig. 4). The selection based on progeny tests had the best combination of high genetic gains and low levels of inbreeding. The selection based on GWAS resulted in lower levels of inbreeding compared to GS

**Table 1.** Results of a post hoc analysis of the number and proportion of true-positive SNPs ( $N$  true SNPs, % true SNPs) in the simulated GWAS, the Pearson correlation coefficient between estimated and true allelic effects ( $r$  allele effects), and the mean cumulative genetic gain after 5 breeding cycles [cumulative genetic gain F5 (%)].

No.	Scenario	Name GWAS	P threshold ( $-\log_{10}$ ) for the top 100 SNPs	N true SNPs	% true SNPs	r allele effects	Cumulative genetic gain F5 (%)
1	1	MLM-3000-F1	4.0017	11	5.5	0.88	20.7
2	2	MLM-3000-F1	4.2420	5	2.5	0.60	32.0
3	3	MLM-3000-F1	7.9000	4	20.0	0.96	17.2
4	1	GLM-3000-Fd	4.3800	28	14.0	0.87	42.2
5	2	GLM-3000-Fd	4.5410	23	12.5	0.79	30.6
6	3	GLM-3000-Fd	4.8717	12	60.0	0.84	28.2
7	1	GLM-3000-extreme	5.7800	44	22.0	0.94	50.5
8	2	GLM-3000-extreme	6.0700	43	21.5	0.96	44.3
9	3	GLM-3000-extreme	9.9100	11	55.0	0.96	25.8
10	1	GLM-10000-Fd	7.1800	45	22.5	0.91	51.8
11	2	GLM-10000-Fd	7.6800	45	22.5	0.96	45.0
12	3	GLM-10000-Fd	11.9000	11	55.0	0.95	30.0

but was in most cases not as effective in terms of genetic gains, as detailed above.

## Discussion

### Good performance of GS and moderate gains from MAS based on GWAS

In our simulation study, the GS performed better in almost all cases compared to MAS scenarios using allelic effects estimated with GWAS. The cross-generation GS (gBLUP-crossG) provided nearly the same genetic gains as the single-generation GS (gBLUP-singleG) in most cases. This is of high relevance for practical breeding programs with trees because it suggests that an effective selection of future parents can be done without phenotypes and thus several years earlier. The MAS based on GWAS with large sample sizes of unrelated individuals outperformed GS (first breeding cycle in scenario 1) only when there was no or only a weak kinship structure. Our results on GWAS using extreme phenotypes indicate that a strong reduction of sample size is possible without losing much performance. It should be noted that the GWAS results not only can be used for MAS but also provide information on the genetic architecture of traits, which can be implemented in GS models to improve prediction accuracy (Morgante *et al.* 2018), as well as underlying candidate genes and biological processes, which can serve as a starting point for trait improvement via gene editing or smart breeding (Wei *et al.* 2021).

Our findings are in accordance with many other publications showing the usefulness of GS in breeding programs. This has been demonstrated with simulation studies on trees several times (Grattapaglia and Resende 2011; Iwata *et al.* 2011; Li and Dungey 2018) and successfully implemented in practical breeding operations first in dairy cattle breeding programs (Su *et al.* 2010), then in agricultural crops (Robertson *et al.* 2019), and little later in forestry species (El-Kassaby *et al.* 2012; Grattapaglia *et al.* 2018; Grattapaglia 2022; Isik 2022).

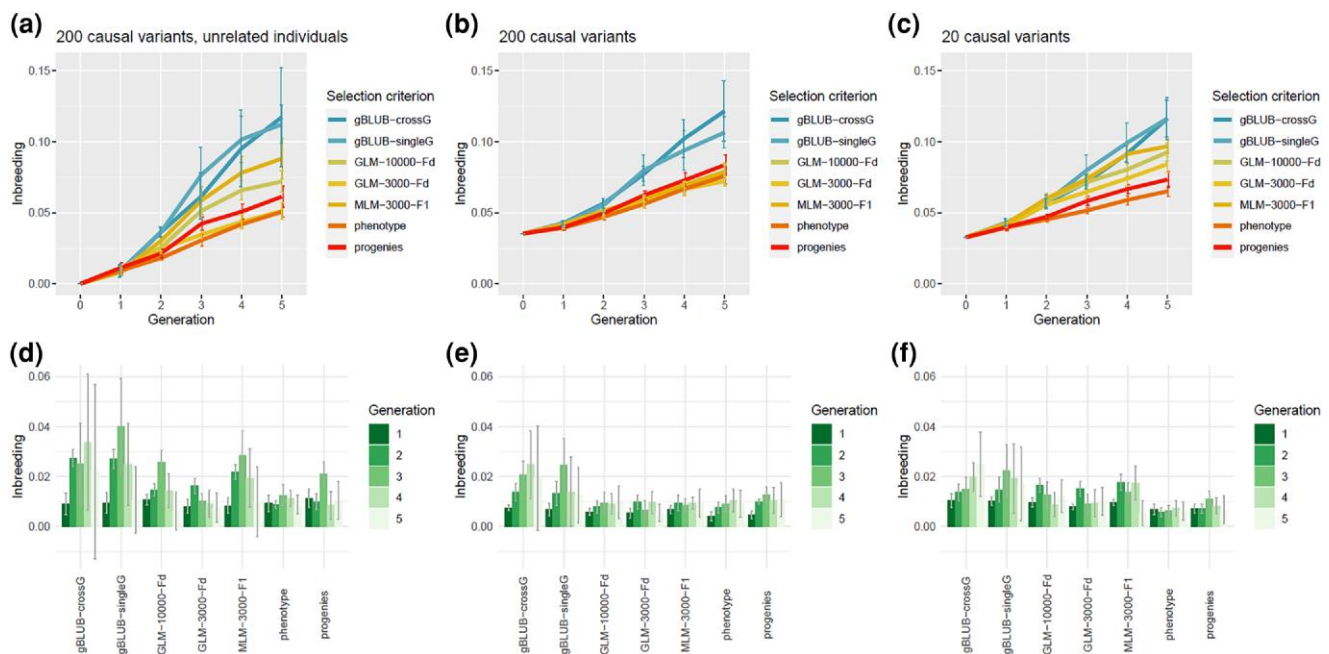
So far, the application of MAS in forest breeding programs could not be successfully established. Grattapaglia *et al.* (2018) stated that while MAS was mostly fruitless, GS has been proven to be very successful. More positive conclusions on the application of GWAS results in breeding have been drawn from crop (Cortes *et al.* 2021; Saini *et al.* 2022) and animal breeding (Gutierrez-Reinoso *et al.* 2021). The main explanation for the difficulties in the application of MAS has been the undiscovered complexity of the genetic architectures of the target traits with many causal loci with

small effects and low frequencies. SNPs with larger effects and SNPs with moderate effects but higher allele frequencies are picked up in our simulations by GWAS. These alleles become fixed in the first breeding cycles and thereby limit additional genetic gains in the following cycles. With high numbers of unrelated individuals ( $>10,000$ ), more causal SNPs can be identified, but because of their small allelic effects and low frequencies, their impact on the genetic gain is small. In most practical breeding programs, phenotyping is much cheaper than whole-genome resequencing. We did some first promising simulations on GWAS with extreme phenotypes (e.g. the 15% edges of the phenotypic distribution). Here, we observe higher genetic gains even with relatively small sample sizes of 3,000 individuals; however, these scenarios need to be studied in more detail.

### Critical increase of inbreeding

In the simulated breeding program, the increase of inbreeding per breeding cycle ( $\Delta F$ ) varied in the different scenarios between 0.006 and 0.023. The highest  $\Delta F$  values were observed for the scenarios with GS. It should be noted that we did not include the negative effects of inbreeding on the target trait and the tree fitness in our simulations and thus overestimated the growth of inbred individuals and potentially the expected genetic gains. This is also true for all other simulations on GS in trees that we are aware of. It is known that inbreeding in trees leads to a higher proportion of homozygotes and thus to inbreeding depression for many traits (Durel *et al.* 1996; Sorensen 1999). Higher mortality and lower growth performance of individuals with higher levels of inbreeding would keep the level of inbreeding lower in real tree breeding programs. Nevertheless, lower growth performance of highly inbred individuals would also be picked up by GS. Thus, we probably overestimated in our simulations the level of inbreeding more than the expected genetic gains.

The optimization of a maximum genetic gain and minimal inbreeding and thus low loss of genetic diversity has been a challenge for most breeding programs. Wu *et al.* (2016) studied the impact of inbreeding depression in a detailed simulation study for various tree breeding strategies. Although they did not directly cover the impact of GS, they found for all simulated breeding strategies a “considerable fixation of unfavorable alleles rendered the purging performance of selfing...” Generally, changes in the selected individuals for the mating and particular crossing schemes using mathematical algorithms are applied to achieve an optimization (Woolliams *et al.* 2015). Usually, breeding programs, especially animal breeding programs, mitigate inbreeding



**Fig. 3.** Inbreeding dynamics under different selection criteria over 5 generations in 3 simulated breeding populations. a–c) Average cumulative inbreeding values ( $n = 10$ ) over 5 breeding cycles are shown for selection based on gBLUP, i.e. GS (blue), GWAS, i.e. MAS (green and yellow), phenotypes (orange), and progeny tests (red). d–f) Average inbreeding values ( $n = 10$ ) per breeding cycle (generation) are shown for the different selection criteria. Error bars represent the standard deviation of the mean. The 3 simulated breeding populations exhibit different genetic architectures, with 200 a, b, d, e) vs 20 c, f) causal variants, and different levels of kinship, with unrelated individuals a, d) vs average kinship values of approximately 0.04 b, c, e, f).

by optimum contribution selection (OCS). This method aims to keep the average coancestry of the selected parents at a certain level and thus controls the short-term and long-term inbreeding. The OCS approach has been improved in several steps, and the most recent methods also consider different levels of genetic introgression (Kohl et al. 2020).

Our simulated  $\Delta F$  values are similar to real tree breeding programs with comparable census numbers and selection intensity in the breeding population. For example, the breeding program of *Pinus taeda* in North Carolina started with 935 selected trees and controlled the inbreeding not to exceed 0.0625 (Isik and McKeand 2019). The optimal balance between the benefits of increased genetic gains and the drawbacks of elevated levels of inbreeding needs to be carefully considered for each specific breeding program.

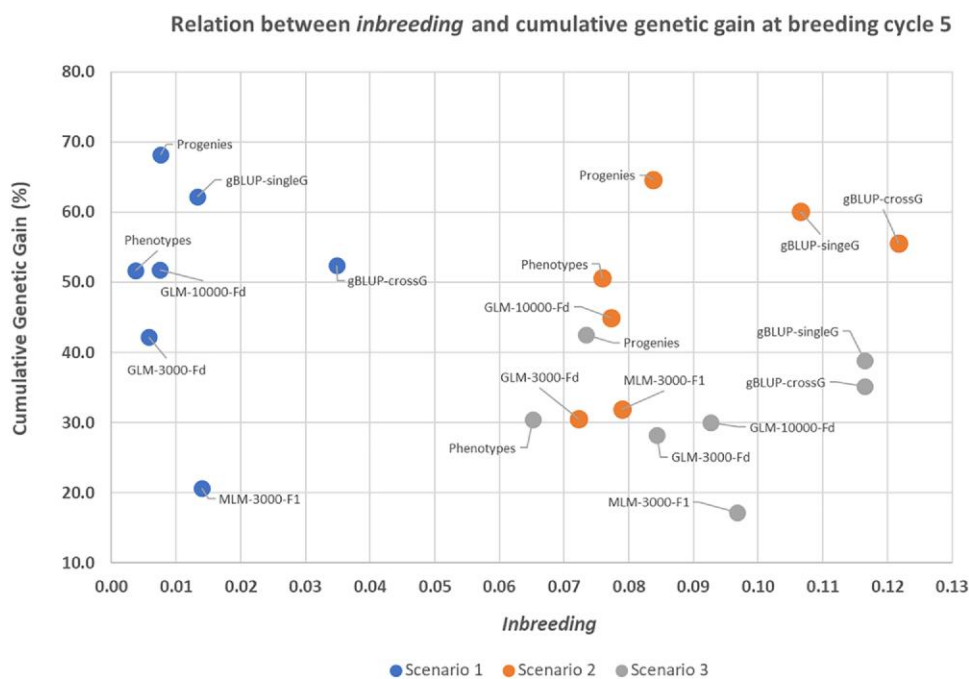
### Implications for practical tree breeding programs

We found in our simulations that the traditional concept of progeny testing (backward selection) would lead to higher genetic gains compared to pure phenotypic selection, GS, and MAS. The main drawbacks of this approach are the enormous workload, time, and costs to establish and phenotype the many different progenies. In our simulations for each of the 1,000 potential parents, a progeny test was performed. For that, the trees first need to reach reproductive maturity, and second, the phenotypes must exhibit a sufficiently stable age–age correlation in order to estimate the breeding values of the mother trees. For most traits, this is the case at one-third of the rotation period (Xie and Ying 1996; Hanaoka and Kato 2022). The simple forward selection based on phenotypes does not come with additional costs for genotyping or the extra time and workload of progeny testing (backward selection), but still, at least one-third of the rotation period is needed to make an accurate measure of the phenotypes,

and for traits with low heritability, the selection by phenotypes is imprecise.

With GS and MAS, the selection of individuals for the next breeding cycle can be done in the seedling stage without phenotyping. In order to make use of the saved years, 2 strategies are possible: (a) the selected individuals could be treated with methods to accelerate the reproductive maturity (phytohormones, top grafting, and in-house seed orchards), and (b) the top selected individuals could be vegetatively propagated and used as deploying clones for plantations after field testing (Li and Dungey 2018). Both the GS and the phenotypic selection rely on repeated phenotyping of the breeding population. In our simulations using MAS, we computed a GWAS only once for a diverse founder population of unrelated individuals. The results from this single GWAS were then used during the entire simulated breeding program. Depending on the target species, this methodological difference can have a profound impact on the feasibility of the whole program. Thus, the repeated phenotyping needed in GS to keep the training population accurate will delay the breeding program by several years compared to selection by MAS. Moreover, there is the advantage of MAS to be more suitable for the integration of new unrelated material into a breeding program and whenever phenotyping is limiting.

Besides the technical considerations, cost–benefit analyses are important for a successful tree breeding strategy (Chamberland et al. 2020). GS and potentially MAS save years of a breeding program and thus enlarge the genetic gain per unit time. They further save costs for phenotyping. But on the other side, the genotyping of SNP arrays comes with costs of about 20–30 US dollars (USD) per individual (Grattapaglia 2022) and whole-genome resequencing of species with the reference genome and small genome size with costs of 100–150 USD per individual, although prices may considerably decrease in the next years. For GS, the optimal training population that needs to be phenotyped and genotyped has a



**Fig. 4.** Relation between inbreeding and the cumulative genetic gain (%) at the end of breeding cycle 5 in 3 different breeding scenarios. Scenario 1 (blue) with 200 causal SNPs and no burn-in, scenario 2 (orange) with 200 causal SNPs and burn-in, and scenario 3 (gray) with 20 causal SNPs and burn-in.

size of a few thousand individuals and should be updated every breeding cycle (Isik 2022). Thus, successful identification of a sufficient proportion of causal SNPs in a GWAS for MAS requires a large investment at the beginning of a breeding program but no further cost later, while GS would require regular investments for the updating of the training data set. With decreasing costs for whole-genome resequencing, the calculation could be in favor of GWAS and MAS. The calculation of genotype–environment associations based on whole-genome data may be another strategy to identify causal SNPs relevant for MAS (Sang et al. 2022; Mueller et al. 2023). With this approach, breeding for traits that are difficult to phenotype such as drought tolerance may be possible. This should be further explored in simulation and experimental studies.

## Conclusions

Using simulations of different breeding populations and strategies, our results further support the potential of GS for forest tree breeding and improvement. Nevertheless, using whole-genome data of large sample sizes or extreme phenotypes for GWAS may provide advantages over GS under certain conditions and could revive efforts for MAS, especially when phenotyping represents a bottleneck. We will study in more detail the possibilities of MAS based on GWAS with extreme phenotypes. Considering the implementation of GS methods in the field of forest tree breeding in the last decade, it will be exciting to follow the impacts on the actual forest ecosystems and further develop strategies to adapt forest tree species to the rapidly changing environmental conditions. In forestry, the genomic revolution has only just begun.

## Data availability

SNPscan breeder has been programmed with Visual Studio 2019 as a .NET application (Framework 4.7.2) and compiled as 64-bit versions for the operating system Microsoft Windows (Windows

11). The program, user manual, and different videos that explain the program are available on our website: <https://www.thuenen.de/en/institutes/forest-genetics/software/SNPscan>.

## Acknowledgments

We are thankful to members of the Center for Integrated Breeding Research (CiBreed) at the University of Göttingen for helpful discussions on SNPscan breeder. We would like to thank Malte Mader for critical testing of the program and for helpful suggestions for its improvement. We thank 3 anonymous reviewers for their helpful comments on a former version of the manuscript.

## Conflicts of interest statement

The authors declare no conflict of interest.

## Literature cited

- Atwell S, Huang YS, Vilhjalmsón BJ, Willems G, Horton M, Li Y, Meng DZ, Platt A, Tarone AM, Hu TT, et al. Genome-wide association study of 107 phenotypes in *Arabidopsis thaliana* inbred lines. *Nature*. 2010;465(7298):627–631. doi:10.1038/nature08800.
- Bradbury PJ, Zhang Z, Kroon DE, Casstevens TM, Ramdoss Y, Buckler ES. Tassel: software for association mapping of complex traits in diverse samples. *Bioinformatics*. 2007;23(19):2633–2635. doi:10.1093/bioinformatics/btm308.
- Chamberland V, Robichaud F, Perron M, Gelinat N, Bousquet J, Beaulieu J. Conventional versus genomic selection for white spruce improvement: a comparison of costs and benefits of plantations on Quebec public lands. *Tree Genet Genomes*. 2020;16(1):17. doi:10.1007/s11295-019-1409-7.
- Cortes LT, Zhang ZW, Yu JM. Status and prospects of genome-wide association studies in plants. *Plant Genome*. 2021;14(1):e20077. doi:10.1002/tpg2.20077.

- Da Silva PHM, Brune A, Alvares CA, do Amaral W, de Moraes MLT, Grattapaglia D, de Paula RC. Selecting for stable and productive families of *Eucalyptus urophylla* across a country-wide range of climates in Brazil. *Can J For Res.* 2019;49(1):87–95. doi:10.1139/cjfr-2018-0052.
- Degen B, Müller NA. SNPscan breeder—a computer program to test genomic tools in breeding programs. *Silvae Genet.* 2023;72(1):126–131. doi:10.2478/sg-2023-0013.
- Durel CE, Bertin P, Kremer A. Relationship between inbreeding depression and inbreeding coefficient in maritime pine (*Pinus pinaster*). *Theor Appl Genet.* 1996;92(3–4):347–356. doi:10.1007/BF00223678.
- El-Kassaby YA, Klápště J, Guy RD. Breeding without breeding: selection using the genomic best linear unbiased predictor method (GBLUP). *New Forests.* 2012;43(5–6):631–637. doi:10.1007/s11056-012-9338-4.
- Endelman JB. Ridge regression and other kernels for genomic selection with R package rrBLUP. *Plant Genome.* 2011;4(3):250–255. doi:10.3835/plantgenome2011.08.0024.
- Grattapaglia D. Twelve years into genomic selection in forest trees: climbing the slope of enlightenment of marker assisted tree breeding. *Forests.* 2022;13(10):1554. doi:10.3390/f13101554.
- Grattapaglia D, Kirst M. *Eucalyptus* applied genomics: from gene sequences to breeding tools. *New Phytol.* 2008;179(4):911–929. doi:10.1111/j.1469-8137.2008.02503.x.
- Grattapaglia D, Resende MDV. Genomic selection in forest tree breeding. *Tree Genet Genomes.* 2011;7(2):241–255. doi:10.1007/s11295-010-0328-4.
- Grattapaglia D, Silva-Junior OB, Kirst M, de Lima BM, Faria DA, Pappas GJ. High-throughput SNP genotyping in the highly heterozygous genome of *Eucalyptus*: assay success, polymorphism and transferability across species. *BMC Plant Biol.* 2011;11(1):65. doi:10.1186/1471-2229-11-65.
- Grattapaglia D, Silva-Junior OB, Resende RT, Cappa EP, Müller BSF, Tan BY, Isik F, Ratcliffe B, El-Kassaby YA. Quantitative genetics and genomics converge to accelerate forest tree breeding. *Front Plant Sci.* 2018;9:1693. doi:10.3389/fpls.2018.01693.
- Gutierrez-Reinoso MA, Aponte PM, Garcia-Herreros M. Genomic analysis, progress and future perspectives in dairy cattle selection: a review. *Animals (Basel).* 2021;11(3):599. doi:10.3390/ani11030599.
- Hanaoka S, Kato K. Estimation of optimal timing of early selection based on time trends of genetic parameters in *Abies sachalinensis*. *Silvae Genet.* 2022;71(1):31–38. doi:10.2478/sg-2022-0004.
- Isik F. Genomic prediction of complex traits in perennial plants: a case for forest trees. In: Ahmadi N, Bartholomé J, editors. *Genomic Prediction of Complex Traits: Methods and Protocols*. New York: Springer; 2022. p. 493–520.
- Isik F, McKeand SE. Fourth cycle breeding and testing strategy for *Pinus taeda* in the NC State University Cooperative Tree Improvement Program. *Tree Genet Genomes.* 2019;15(5):12. doi:10.1007/s11295-019-1377-y.
- Iwata H, Hayashi T, Tsumura Y. Prospects for genomic selection in conifer breeding: a simulation study of *Cryptomeria japonica*. *Tree Genet Genomes.* 2011;7(4):747–758. doi:10.1007/s11295-011-0371-9.
- Kastally C, Niskanen AK, Perry A, Kujala ST, Avia K, Cervantes S, Haapanen M, Kesalahti R, Kumpula TA, Mattila TM, et al. Taming the massive genome of Scots pine with PiSy50k, a new genotyping array for conifer research. *Plant J.* 2022;109(5):1337–1350. doi:10.1111/tbj.15628.
- Kiszonas AM, Morris CF. Wheat breeding for quality: a historical review. *Cereal Chem.* 2018;95(1):17–34. doi:10.1094/CCHEM-05-17-0103-FI.
- Kohl S, Wellmann R, Herold P. Advanced optimum contribution selection as a tool to improve regional cattle breeds: a feasibility study for Vorderwald cattle. *Animal.* 2020;14(1):1–12. doi:10.1017/S1751731119001484.
- Li YJ, Dungey HS. Expected benefit of genomic selection over forward selection in conifer breeding and deployment. *PLoS One.* 2018;13(12):e0208232. doi:10.1371/journal.pone.0208232.
- Li X, Liu XT, Wei JT, Li Y, Tigabu M, Zhao XY. Genetic improvement of *Pinus koraiensis* in China: current situation and future prospects. *Forests.* 2020;11(2):13. doi:10.3390/f11020148.
- Liu HM, Tessema BB, Jensen J, Cericola F, Andersen JR, Sorensen AC. ADAM-Plant: a software for stochastic simulations of plant breeding from molecular to phenotypic level and from simple selection to complex speed breeding programs. *Front Plant Sci.* 2019;9:15. doi:10.3389/fpls.2018.01926.
- Misztal I, Lourenco D, Legarra A. Current status of genomic evaluation. *J Anim Sci.* 2020;98(4):14. doi:10.1093/jas/skaa101.
- Morgante F, Huang W, Maltecca C, Mackay TFC. Effect of genetic architecture on the prediction accuracy of quantitative traits in samples of unrelated individuals. *Heredity (Edinb).* 2018;120(6):500–514. doi:10.1038/s41437-017-0043-0.
- Mphahlele MM, Isik F, Mostert-O'Neill MM, Reynolds SM, Hodge GR, Myburg AA. Expected benefits of genomic selection for growth and wood quality traits in *Eucalyptus grandis*. *Tree Genet Genomes.* 2020;16(4):12. doi:10.1007/s11295-020-01443-1.
- Mueller NA, Gessner C, Mader M, Blanc-Jolivet C, Fladung M, Degen B. Genomic variation of a keystone forest tree species reveals patterns of local adaptation and future maladaptation. *bioRxiv.* 2023.1105.1111.540382. <https://doi.org/10.1101/2023.1105.1111.540382>.
- Nadeem MA, Nawaz MA, Shahid MQ, Dogan Y, Comertpay G, Yildiz M, Hatipoglu R, Ahmad F, Alsaleh A, Labhane N, et al. DNA molecular markers in plant breeding: current status and recent advancements in genomic selection and genome editing. *Biotechnol Biotechnol Equip.* 2018;32(2):261–285. doi:10.1080/13102818.2017.1400401.
- Nantongo JS, Potts BM, Klápště J, Graham NJ, Dungey HS, Fitzgerald H, O'Reilly-Wapstra JM. Genomic selection for resistance to mammalian bark stripping and associated chemical compounds in radiata pine. *G3 (Bethesda).* 2022;12(11):jkac245. doi:10.1093/g3journal/jkac245.
- Pérez-Enciso M, Forneris N, de los Campos G, Legarra A. Evaluating sequence-based genomic prediction with an efficient new simulator. *Genetics.* 2017;205(2):939–953. doi:10.1534/genetics.116.194878.
- Pfenninger M, Reuss F, Kiebler A, Schonnenbeck P, Caliendo C, Gerber S, Cocchiarraro B, Reuter S, Bluthgen N, Mody K, et al. Genomic basis for drought resistance in European beech forests threatened by climate change. *eLife.* 2021;10:e65532. doi:10.7554/eLife.65532.
- Plomion C, Aury JM, Amselem J, Alaeitabar T, Barbe V, Belsler C, Berges H, Bodenes C, Boudet N, Boury C, et al. Decoding the oak genome: public release of sequence data, assembly, annotation and publication strategies. *Mol Ecol Resour.* 2016;16(1):254–265. doi:10.1111/1755-0998.12425.
- R Core Team. R: a language and environment for statistical computing. R foundation for statistical computing 2022. [accessed]. <https://www.R-project.org>
- Resende MDV, Resende MFR, Sansaloni CP, Petroli CD, Missiaggia AA, Aguiar AM, Abad JM, Takahashi EK, Rosado AM, Faria DA, et al. Genomic selection for growth and wood quality in *Eucalyptus*: capturing the missing heritability and accelerating breeding for complex traits in forest trees. *New Phytol.* 2012;194(1):116–128. doi:10.1111/j.1469-8137.2011.04038.x.



- Resende RT, Resende MDV, Silva FF, Azevedo CF, Takahashi EK, Silva-Junior OB, Grattapaglia D. Assessing the expected response to genomic selection of individuals and families in *Eucalyptus* breeding with an additive-dominant model. *Heredity* (Edinb). 2017;119(4):245–255. doi:10.1038/hdy.2017.37.
- Robertson CD, Hjortshoj RL, Janss LL. Genomic selection in cereal breeding. *Agronomy-Basel*. 2019;9(2):95. doi:10.3390/agronomy9020095.
- Saini DK, Chopra Y, Singh J, Sandhu KS, Kumar A, Bazzar S, Srivastava P. Comprehensive evaluation of mapping complex traits in wheat using genome-wide association studies. *Mol Breed*. 2022;42(1):1. doi:10.1007/s11032-021-01272-7.
- Sandhu KS, Merrick LF, Sankaran S, Zhang ZW, Carter AH. Prospectus of genomic selection and phenomics in cereal, legume and oilseed breeding programs. *Front Genet*. 2022;12:829131. doi:10.3389/fgene.2021.829131.
- Sang YP, Long ZQ, Dan XM, Feng JJ, Shi TT, Jia CF, Zhang XX, Lai Q, Yang GL, Zhang HY, et al. Genomic insights into local adaptation and future climate-induced vulnerability of a keystone forest tree in East Asia. *Nat Commun*. 2022;13(1):6541. doi:10.1038/s41467-022-34206-8.
- Sollars ESA, Harper AL, Kelly LJ, Sambles CM, Ramirez-Gonzalez RH, Swarbreck D, Kaithakottil G, Cooper ED, Uauy C, Havlickova L, et al. Genome sequence and genetic diversity of European ash trees. *Nature*. 2017;541(7636):212–216. doi:10.1038/nature20786.
- Sorensen FC. Relationship between self-fertility, allocation of growth, and inbreeding depression in three coniferous species. *Evolution*. 1999;53(2):417–425. doi:10.2307/2640778.
- Su G, Guldbbrandtsen B, Gregersen VR, Lund MS. Preliminary investigation on reliability of genomic estimated breeding values in the Danish Holstein population. *J Dairy Sci*. 2010;93(3):1175–1183. doi:10.3168/jds.2009-2192.
- Vidal M, Plomion C, Raffin A, Harvengt L, Bouffier L. Forward selection in a maritime pine polycross progeny trial using pedigree reconstruction. *Ann For Sci*. 2017;74(1):21. doi:10.1007/s13595-016-0596-8.
- Wainschtein P, Jain D, Zheng ZL, Cupples LA, Shadyab AH, McKnight B, Shoemaker BM, Mitchell BD, Psaty BM, Kooperberg C, et al. Assessing the contribution of rare variants to complex trait heritability from whole-genome sequence data. *Nature Genet*. 2022;54(3):263–273. doi:10.1038/s41588-021-00997-7.
- Wei X, Qiu J, Yong KC, Fan JJ, Zhang Q, Hua H, Liu J, Wang Q, Olsen KM, Han B, et al. A quantitative genomics map of rice provides genetic insights and guides breeding. *Nature Genet*. 2021;53(2):243–253. doi:10.1038/s41588-020-00769-9.
- Woolliams JA, Berg P, Dagnachew BS, Meuwissen THE. Genetic contributions and their optimization. *J Anim Breed Genet*. 2015;132(2):89–99. doi:10.1111/jbg.12148.
- Wu HX, Hallingback HR, Sanchez L. Performance of seven tree breeding strategies under conditions of inbreeding depression. *G3* (Bethesda). 2016;6(3):529–540. doi:10.1534/g3.115.025767.
- Würschum T. Mapping QTL for agronomic traits in breeding populations. *Theor Appl Genet*. 2012;125(2):201–210. doi:10.1007/s00122-012-1887-6.
- Xie CY, Ying CC. Heritabilities, age-age correlations, and early selection in lodgepole pine (*Pinus contorta* ssp *latifolia*). *Silvae Genet*. 1996;45(2–3):101–107.
- Yengo L, Vedantam S, Marouli E, Sidorenko J, Bartell E, Sakaue S, Graff M, Eliassen AU, Jiang YX, Raghavan S, et al. A saturated map of common genetic variants associated with human height. *Nature*. 2022;610(7933):704–712. doi:10.1038/s41586-022-05275-y.

Editor: E. Huang