

ARTICLE OPEN



Development and optimization of expected cross value for mate selection problems

Pouya Ahadi¹, Balabhaskar Balasundaram², Juan S. Borrero³ and Charles Chen³✉

© The Author(s) 2024

In this study, we address the mate selection problem in the hybridization stage of a breeding pipeline, which constitutes the multi-objective breeding goal key to the performance of a variety development program. The solution framework we formulate seeks to ensure that individuals with the most desirable genomic characteristics are selected to cross in order to maximize the likelihood of the inheritance of desirable genetic materials to the progeny. Unlike approaches that use phenotypic values for parental selection and evaluate individuals separately, we use a criterion that relies on the genetic architecture of traits and evaluates combinations of genomic information of the pairs of individuals. We introduce the *expected cross value* (ECV) criterion that measures the expected number of desirable alleles for gametes produced by pairs of individuals sampled from a population of potential parents. We use the ECV criterion to develop an integer linear programming formulation for the parental selection problem. The formulation is capable of controlling the inbreeding level between selected mates. We evaluate the approach on two applications: (i) improving multiple target traits simultaneously, and (ii) finding a multi-parental solution to design crossing blocks. We evaluate the performance of the ECV criterion using a simulation study. Finally, we discuss how the ECV criterion and the proposed integer linear programming techniques can be applied to improve breeding efficiency while maintaining genetic diversity in a breeding program.

Heredity (2024) 133:113–125; <https://doi.org/10.1038/s41437-024-00697-y>

INTRODUCTION

Plant and animal breeding consists of methodologies for the creation, selection, and fixation of superior phenotypes to fulfill the breeding goals of increasing productivity and financial returns, improving welfare, and reducing environmental impact Oldenbroek and van der Waaij (2015). Traditionally, breeders achieve these goals by identifying the individuals with desirable phenotypes and crossing them to produce the segregation of phenotypes in a new generation that allows further selection for advancement. This breeding strategy is perpetuated because high-volume crossing and evaluation led to the identification of the iconic *Green Revolution* varieties that successfully doubled rice and wheat yields from the 1960s to 1990s (Hesser 2006), despite the inevitable inefficiency of producing a high number of failed crosses. However, the future of food security and livestock will be driven not only by the demand but also by severe competition with other uses of land and water resource (Cassandro 2020). Therefore, more efficient breeding strategies ought to be considered because making many crosses with the knowledge that most fail is not justified either by theory or comparative experiments, and is also socially unacceptable.

Ultimately, the overall objective of a breeding program is to produce lines and varieties that are genetically homogeneous and perform at a high level, with end-use quality supportive of the intended market class. A wheat breeding pipeline, for instance, would begin with assembling parental stocks with a careful

examination of available germplasm and donor traits. In principle, this is to construct and partition parental stocks respective to a specific goal or goals, to create the genetic variability needed for producing an adapted, high-yielding pure-line variety with perceived quality demands in the future marketplace.

With the continued advancement of genomic technologies and steady decline in genotyping costs, breeders are now able to take full advantage of the availability of genetic information embedded in the genome (Hayes et al. 2009; Heffner et al. 2010). Nevertheless, except for the potential application of a higher selection intensity with GEBVs (genomic estimated breeding values) (Meuwissen et al. 2001), experimental data for the optimal number of crosses as well as the optimal numbers of progeny to sample from each cross required for selection as the initial investment to fulfill the breeding goal have not been reported in the literature (Donald 1968). This is unsurprising given that the number of individuals a breeding program can phenotypically evaluate is resource-limited (Rincent et al. 2017). For example, consider a single cross of two genetically distinct parental lines with 100 QTLs associated with variability among multiple desirable traits. Assuming independent assortment and co-dominance, the complete population from this single pair of founders will consist of $3^{100} \approx 5.1 \times 10^{47}$ genotypic combinations. Even considering a moderate number of 200 wheat lines in the parental stocks, the number of combinations to be evaluated in the field is astronomically high (Beans 2020). Consequently,

¹H. Milton Stewart School of Industrial and Systems Engineering, Georgia Institute of Technology, Atlanta, GA, USA. ²School of Industrial Engineering and Management, Oklahoma State University, Stillwater, OK, USA. ³Department of Biochemistry and Molecular Biology, Oklahoma State University, Stillwater, OK, USA. Associate editor: Armando Caballero ✉email: charles.chen@okstate.edu

Received: 5 February 2024 Revised: 11 June 2024 Accepted: 12 June 2024
Published online: 2 July 2024

analytical approaches based on operations research, mathematical optimization, and statistical learning to optimize breeding decision-making have gained prominence over the years (Byrum 2015; Byrum et al. 2017, 2016; Johnson et al. 1988; Kusmec et al. 2021).

There are two essential steps to addressing this problem using mathematical optimization. The first is to define a fitness criterion to evaluate individuals or crosses based on genetic information. The second step is to devise a mathematical optimization framework that incorporates the fitness criterion along with other essential requirements of the breeding program, and whose objective is to find the individuals or crosses that maximize the fitness criterion. The mathematical optimization framework, while faithfully capturing the various breeding requirements and objectives, must also be computationally viable in order for it to be useful in practice.

In contrast to addressing these breeding challenges in a traditional phenotype-centric paradigm, genetic improvement can also be more efficiently achieved by transferring desirable alleles from parents to progeny as a genetic process while avoiding alleles that show antagonistic pleiotropic effects. Therefore, our aim is to devise a multi-objective mathematical optimization framework that targets more than one phenotype and generates multiple crosses that identify a set of best parental pairs from populations to address multiple breeding goals simultaneously.

As it is to be expected in any non-trivial multi-criteria decision-making setting, the criteria (breeding objectives) can be mutually conflicting, making it challenging to design an effective multi-objective optimization framework. For example, yield production in wheat has been found to be negatively correlated with grain protein content (Simmonds 1995), which is an essential factor for its commercial demand (Visscher et al. 1996). This makes concurrently fulfilling breeding goals of high yielding and protein content difficult. The negative correlation between the mass of beef cows and various measures of fertility and stayability could have attributed to the increasing concerns about compromised reproductive efficiency as a result of selection for growth (Berry and Evans 2014; Mwansa et al. 2002).

In this study, we propose a new fitness criterion called the *expected cross value* (ECV) that is inspired by a related fitness criterion called *predicted cross value* (PCV) introduced by Han et al. (2017). ECV returns a probabilistic measure of the fitness of the progeny of a specific pair of individuals based on the genetic architecture of trait variation. We consider the complexity of genetic architecture that underlies agronomic performance characteristics and develop an integer linear programming formulation of the parental pair selection problem that optimizes the ECV criterion. We further extend its capability to select multiple pairs of parents. Our optimization framework is based on the genetic transmission of all detectable genetic loci and can mitigate the potential impact of crossing within highly related individuals. Based on simulation studies, we demonstrate that using ECV as a fitness criterion would address the limitations of other related approaches for mate selection problems, and our multi-objective methodology can simultaneously improve a group of target phenotypic traits.

METHODS

We begin with some preliminaries needed to formally define the *expected cross value* (ECV) as a new fitness criterion for the mate selection problem. Considering all diploid and polyploid species that may behave as diploids cytologically, e.g., bread wheat (Riley and Chapman 1958), we assume that the variability of target traits is governed by segregating alleles at N different loci of all chromosomes in the genome. We use the index set notation $[a] = \{1, 2, \dots, a\}$ for a positive integer a , and define the genotype matrix next.

Definition 1. Given an individual k , we define its genotype matrix L^k as an $N \times 2$ binary matrix with the i, j -th entry for every $i \in [N]$ and $j \in [2]$ given by:

$$L_{ij}^k = \begin{cases} 1, & \text{if the allele in locus } i \text{ of gamete from parent } j \text{ is desirable,} \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

Genotype matrix information of all individuals is an input for the ECV. Hereafter, we refer to the allele in QTL i as the i -th allele for ease of discussion. In our simulations, alleles are desirable when they enhance the trait value, assuming larger the positive value, the better. Observe that each column of L^k represents a gamete from one of the parents of individual k .

We model how alleles transfer from parents to children, i.e., how a gamete inherits alleles from the parent, by using a random N -dimensional binary vector J_i with each component being a random variable J_i for each $i \in [N]$ defined as follows:

$$J_i = \begin{cases} 0, & \text{if the } i\text{-th allele is transferred from the first column of } L^k \text{ to the gamete,} \\ 1, & \text{otherwise.} \end{cases} \quad (2)$$

For a given individual k and QTL i , random variable J_i determines which of the gametes that comprise the genome of individual k will transfer the allele in the i -th locus.

Definition 2. (Han et al. (2017)). We say that the random vector $J \in \{0, 1\}^N$ follows an inheritance distribution with parameters $r \in [0, 0.5]^{N-1}$ and a_0 , denoted by $J \sim \mathcal{I}(r, a_0)$, if and only if

$$\Pr(J_1 = 0) = a_0, \Pr(J_1 = 1) = 1 - a_0, \quad (3)$$

$$\Pr(J_i = J_{i-1}) = 1 - r_{i-1}, \Pr(J_i = 1 - J_{i-1}) = r_{i-1}, \quad \forall i \in [N] \text{ and } i \geq 2. \quad (4)$$

In Definition 2, the $(N-1)$ -dimensional vector $r \in [0, 0.5]^{N-1}$ represents the recombination frequencies between the consecutive pairs of loci. The value of r_i is the probability that the i -th and $(i+1)$ -th alleles come from different gametes that comprise the genome of individual k . Note that if $r_i = 0$ for all $i \in [N-1]$, then the gamete produced by individual k is identical to exactly one of the parental gametes, while the maximum possible recombination between gametes is expected to be observed when $r_i = 0.5$ for all $i \in [N-1]$.

Deriving the closed-form marginal inheritance distributions

Given $J \sim \mathcal{I}(r, a_0)$, we now derive the marginal distribution of J_i for each $i \in [N]$. The closed-form expressions so obtained then allow us to compute the expectations required to obtain a general closed-form expression for the ECV. For each $i \in [N]$, define the recursive function $\phi_i : \mathbb{R}^{N-1} \rightarrow \mathbb{R}$ as follows:

$$\phi_1(r) = 0, \quad \phi_2(r) = r_1, \quad (5)$$

$$\phi_i(r) = r_{i-1} + (1 - 2r_{i-1})\phi_{i-1}(r), \quad \forall i \in \{3, 4, \dots, N\}. \quad (6)$$

Proposition 1. Suppose that $J \sim \mathcal{I}(r, a_0)$. Then, for each $i \in [N]$, the marginal distribution of J_i satisfies the following equations:

$$\Pr(J_i = 0) = a_0 + (1 - 2a_0)\phi_i(r), \quad (7a)$$

$$\Pr(J_i = 1) = 1 - a_0 + (2a_0 - 1)\phi_i(r). \quad (7b)$$

Proposition 1 (proved in the Supplementary Information) establishes the marginal distributions through a recursion, which can be used to obtain a closed-form expression. This result can be further simplified using the laws of inheritance that the allele pairs of a locus segregate randomly during meiosis, and each allele transmits to the gamete with equal probability. Specifically, Proposition 1 then implies the following corollary.

Corollary 1. Assume that Mendel's second law holds and $\alpha_0 = 0.5$. Then,
 $\Pr(J_i = 0) = \Pr(J_i = 1) = 0.5, \forall i \in [N].$ (8)

Furthermore,

$$\mathbb{E}(J_i) = 0 \times \Pr(J_i = 0) + 1 \times \Pr(J_i = 1) = 0.5, \forall i \in [N], \quad (9)$$

where $\mathbb{E}(\cdot)$ represents the expectation operator.

The gamete and loss functions

The inheritance distribution characterizes the source of alleles transmitted from a parent to its gametes. Therefore, we can define a so-called *gamete function* to specify the alleles in the gamete according to the inheritance distribution. Given this gamete function, we derive a closed-form expression for the ECV of a pair of individuals.

Definition 3. (Han et al. (2017)). Given an individual with genotype matrix L and a vector $J \sim \mathcal{I}(r, \alpha_0)$, the vector-valued gamete function $\text{gam}: (L, J) \mapsto g$ outputs the binary gamete vector g defined as follows for each $i \in [N]$:

$$g_i = \begin{cases} L_{i,1}, & \text{if } J_i = 0, \\ L_{i,2}, & \text{if } J_i = 1. \end{cases} \quad (10)$$

Equivalently, $g_i = L_{i,1}(1 - J_i) + L_{i,2}J_i$.

Suppose we have two individuals with genotype matrices L^1 and L^2 , and two independent random vectors $J^1, J^2 \sim \mathcal{I}(r, \alpha_0)$. By crossing these two individuals, the genotype matrix for a child in the progeny is given by matrix $[g^1, g^2]$ where $g^1 = \text{gam}(L^1, J^1)$ and $g^2 = \text{gam}(L^2, J^2)$. Then, the gamete that is produced by a child of this progeny for the next generation is given by:

$$g^3 = \text{gam}([g^1, g^2], J^3), \quad (11)$$

where $J^3 \sim \mathcal{I}(r, \alpha_0)$ is independent of J^1 and J^2 . Below, we define a *loss function* in terms of the g^3 gamete vector that will lead us to the ECV criterion.

Definition 4. Suppose L^1 and L^2 are the genotype matrices of two individuals and let $J^k, k = 1, 2, 3$, be independent random vectors following the distribution $\mathcal{I}(r, \alpha_0)$ for some given r and α_0 . Let $g^k = \text{gam}(L^k, J^k)$ for $k = 1, 2$ and $g^3 = \text{gam}([g^1, g^2], J^3)$. We define the loss function associated with L^1, L^2, r , and α_0 as the following random variable:

$$\text{loss}(L^1, L^2, r, \alpha_0) = \sum_{i=1}^N (1 - g_i^3) = N - \sum_{i=1}^N g_i^3. \quad (12)$$

The loss function counts the number of undesirable alleles in the gamete g^3 . If the loss function is equal to 0, then all alleles in g^3 are desirable, while the opposite is true if it is equal to N . Before deriving our ECV criterion, we introduce the related PCV criterion of Han et al. (2017).

Definition 5. (Han et al. (2017)). Let L^1 and L^2 be the genotype matrices of two individuals, and let r and α_0 be given. Define the gamete g^3 using Eq. (11). Then, the PCV associated with L^1, L^2, r , and α_0 is the probability that the gamete g^3 contains only desirable alleles. That is,

$$\text{PCV}(L^1, L^2, r, \alpha_0) = \Pr(\text{loss}(L^1, L^2, r, \alpha_0) = 0). \quad (13)$$

The expected cross value criterion

Next, we use the loss function in Definition 4 to define the ECV, an alternative criterion to PCV, based on allelic information of individuals. The measure depends on the gamete g^3 defined in Eq. (11) and can evaluate a pair of individuals that could be mated.

Definition 6. For a selected pair of individuals with genotype matrices L^1 and L^2 , the ECV is the expected number of desirable alleles in gamete g^3 defined in Equation (11). As the loss function represents the number of

undesirable alleles in g^3 , the ECV can be computed as:

$$\text{ECV}(L^1, L^2, r, \alpha_0) = N - \mathbb{E}(\text{loss}(L^1, L^2, r, \alpha_0)) = \mathbb{E}\left(\sum_{i=1}^N g_i^3\right). \quad (14)$$

A pair of individuals with the highest ECV value could be selected as parents for crossing. Theorem 1 (proved in the Supplementary Information) constitutes our main result that provides a closed-form expression for calculating ECV for a pair of parents.

Theorem 1. Assume Mendel's second law holds true and let L^1 and L^2 be the genotype matrices of two individuals. The ECV corresponding to the desired phenotypic trait can be computed using the following equation:

$$\text{ECV}(L^1, L^2, r, 0.5) = 0.25 \sum_{i=1}^N (L_{i,1}^1 + L_{i,2}^1 + L_{i,1}^2 + L_{i,2}^2). \quad (15)$$

Remark 1. Without relying on Mendel's second law, the ECV can still be computed in closed-form more generally as:

$$\begin{aligned} \text{ECV}(L^1, L^2, r, \alpha_0) &= \sum_{i=1}^N \left(L_{i,1}^1 + [1 - \alpha_0 + (2\alpha_0 - 1)\phi_i(r)](L_{i,2}^1 - 2L_{i,1}^1 + L_{i,1}^2) \right. \\ &\quad \left. + [1 - \alpha_0 + (2\alpha_0 - 1)\phi_i(r)]^2(L_{i,2}^2 + L_{i,1}^1 - L_{i,2}^1 - L_{i,1}^2) \right). \end{aligned}$$

Theorem 1 provides a closed-form expression for the ECV criterion that enables us to formulate the parental selection problem as an integer linear program.

Single-trait parental selection problem

We develop an integer programming (IP) formulation for the parental selection problem using the ECV criterion as the single optimization objective (see Supplementary Formulation (27)) and the constraint system (and decision variables) from the mixed-integer programming formulation for the PCV introduced by Han et al. (2017). The formulation finds the best pair of individuals maximizing the ECV criterion based on a desired phenotypic trait. In addition, we restrict the inbreeding between selected individuals by preventing pairs of individuals with a sufficiently large inbreeding value from being selected as parents. By using the marker genotype information we can construct the genomic matrix G that quantifies the genomic relationship between any pair of individuals in the population (VanRaden 2008). Any pair in the population that has a genomic relationship (i.e., inbreeding value) higher than a pre-determined parameter ϵ , will be excluded from the set of feasible pairs using a family of constraints we include in the formulation.

In a breeding program, we may also seek to find multiple pairs for crossing, rather than just a single pair. In order to do so, we introduce Supplementary Algorithm 1 that iteratively solves our IP formulation for the single-trait parental pair selection problem. Note that solving the Supplementary Formulation (27) will identify a pair of individuals as the optimal solution for the problem. By adding "conflict constraints" corresponding to this optimal pair to the formulation, we can exclude just this optimal solution from the set of feasible solutions and reoptimize to find the next optimal pair. We can repeat this process until the required number of pairs have been chosen for the crossing program (assuming that many solutions exist).

The flowchart in Fig. 1 illustrates the workflow of the proposed ECV approach for mate selection problems for a single trait. The process begins with an initial population where genetic marker and QTL information are available for the selection of parental lines to assemble the crossing block to advance specific breeding targets (Velu and Singh 2013). The ECV criteria can be optimized over several generations (denoted by T in Fig. 1). In each generation, genomic information related to QTLs and genetic markers, along with a genetic relationship matrix (G) is used for constructing the optimization model detailed in Supplementary Formulation (27), and solving it to find an optimal set of mating pairs for crossing. The workflow for solving the multi-trait parental selection by optimizing the ECV metric mirrors the process in Fig. 1 for single-trait ECV

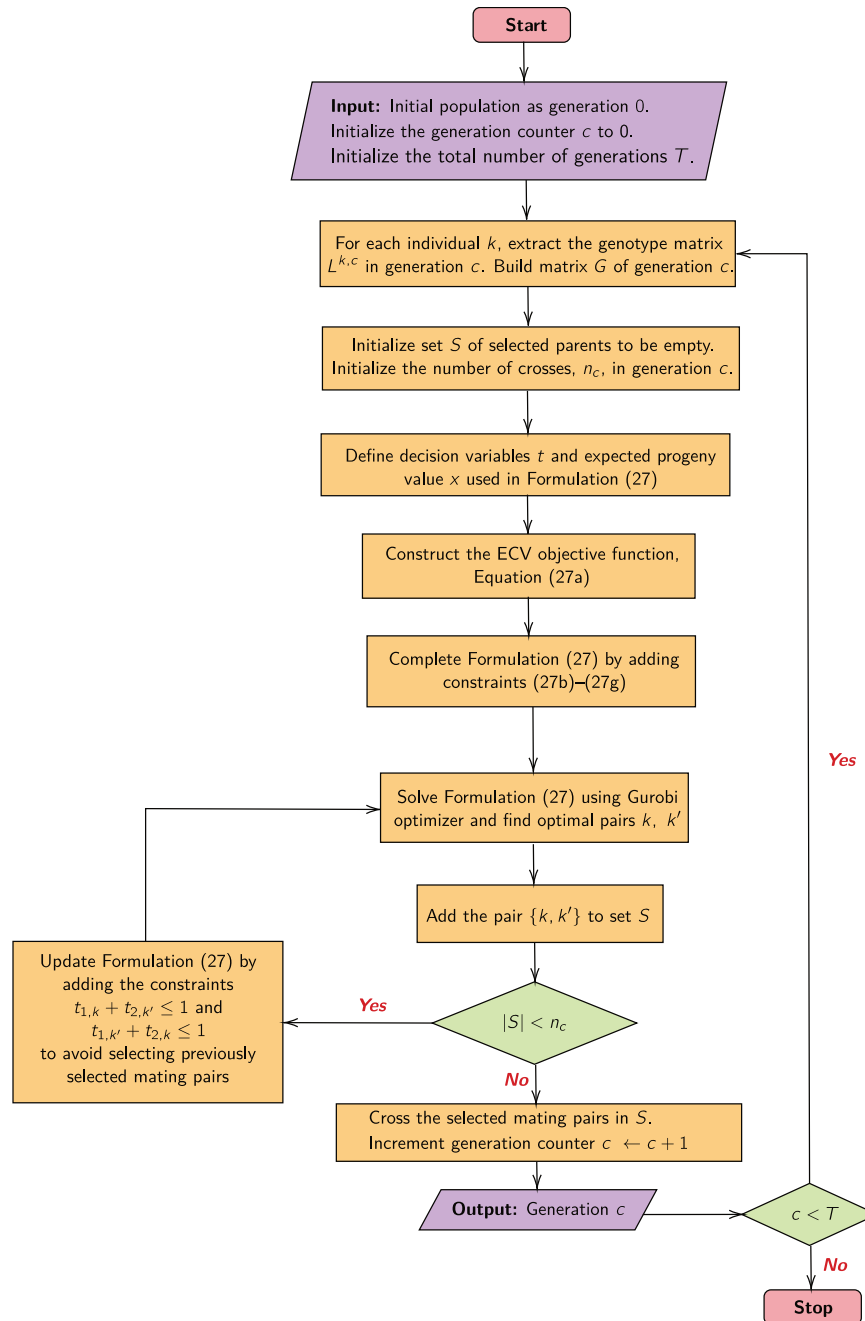


Fig. 1 A flowchart describing the overall process for single-trait ECV optimization producing n_c mating pairs in each generation $c = 1, 2, \dots, T$. The optimization formulation and equations referred in the flowchart are included in the Supplementary Information.

optimization. The key difference is that we solve the Supplementary Formulation (29) via lexicographic optimization with user-specified degradation tolerances as described in detail in the Supplementary Information.

Multi-trait parental selection problem

In general, breeders may be interested in improving several phenotypic traits simultaneously. In this case, we need to extend the ECV criterion to account for multiple traits. We assume there are M target traits in the breeding program and that the ℓ -th desired trait for every $\ell \in [M]$ is affected by N_ℓ different QTL in the genome. For each individual we define M genotype matrices, one for each trait. Each such matrix is an $N_\ell \times 2$ binary matrix in which each row represents the pair of alleles in the corresponding genetic locus. Thus, we extend the previous definitions as follows.

Definition 7. For $k \in [K]$ and $\ell \in [M]$, the genotype matrix $L^{k,\ell}$ associated with the k -th individual and the ℓ -th trait is defined as:

$$L_{ij}^{k,\ell} = \begin{cases} 1 & \text{if } i\text{-th allele of gamete from parent } j \text{ is desirable for trait } \ell, \\ 0 & \text{otherwise.} \end{cases} \quad \forall i \in [N_\ell], j \in [2]. \quad (16)$$

Consider two individuals with genotype matrices $L^{1,\ell}$ and $L^{2,\ell}$ for target trait $\ell \in [M]$, and suppose that we have three independent random vectors J^1 , J^2 and J^3 following an inheritance distribution $\mathcal{I}(r, a_0)$. Using the definition of gamete function (10), the genotype matrix corresponding to the ℓ -th trait for a child in the progeny is represented by matrix $[g^{1,\ell} g^{2,\ell}]$ where $g^{1,\ell} = \text{gam}(L^{1,\ell}, J^1)$ and $g^{2,\ell} = \text{gam}(L^{2,\ell}, J^2)$. The gamete that is

produced by this progeny for the next generation is then given by:

$$g^{3,\ell} = \text{gam}([g^{1,\ell}, g^{2,\ell}], J^3). \quad (17)$$

Definition 8. For the ℓ -th target trait and a selected pair of individuals with genotype matrices $L^{1,\ell}$ and $L^{2,\ell}$, the ECV^ℓ is the expected number of desirable alleles of trait ℓ in gamete $g^{3,\ell}$. Following Equation (17), ECV^ℓ , for each $\ell \in [M]$ is defined as:

$$ECV^\ell(L^{1,\ell}, L^{2,\ell}, r, a_0) = N_\ell - \mathbb{E}(\text{loss}(L^{1,\ell}, L^{2,\ell}, r, a_0)) = \mathbb{E}\left(\sum_{i=1}^{N_\ell} g_i^{3,\ell}\right). \quad (18)$$

Following Theorem 1, we can obtain a closed-form expression for ECV^ℓ function.

Theorem 2. Assume Mendel's second law holds true and let $\ell \in [M]$. Then, for a selected pair of individuals with genotype matrices $L^{1,\ell}$ and $L^{2,\ell}$, the ECV corresponding to the ℓ -th target phenotypic trait can be computed as:

$$ECV^\ell(L^{1,\ell}, L^{2,\ell}, r, 0.5) = 0.25 \sum_{i=1}^{N_\ell} (L_{i,1}^{1,\ell} + L_{i,2}^{1,\ell} + L_{i,1}^{2,\ell} + L_{i,2}^{2,\ell}). \quad (19)$$

Ideally, a breeding program would like to select parental pair(s) that simultaneously optimize all the ECV^ℓ functions. Such an optimum is not likely to exist in practice because some phenotypic traits are negatively correlated. Therefore, improving one trait might worsen others. In order to achieve a reasonable trade-off, one turns to the theory of multi-objective optimization.

Consider a vector of objective functions $F(t, x) = (f_1(t, x^1), f_2(t, x^2), \dots, f_M(t, x^M))$ where $f_\ell(t, x^\ell)$ denotes the ECV function (19) corresponding to ℓ -th trait. Supplementary Formulation (29) for the multi-trait parental selection problem seeks to find a pair of individuals that will “maximize” the vector-valued objective function. Similar to the single-trait optimization model, this formulation also excludes pairs of individuals with genomic relationship exceeding the tolerance threshold from the set of feasible solutions. Furthermore, as explained in the previous section, this approach can also be extended to select multiple parental pairs for the breeding program by iteratively adding conflict constraints. The differences lie in the handling of multiple traits, especially in the vector objective function.

Multi-objective or vector optimization problems are commonly handled by scalarization—converting the vector optimization problem into one or more scalar optimization problems (Miettinen 2012; Sawaragi et al. 1985); see survey by Miettinen et al. (2016) for interactive and other methods. One approach is to use a weighted combination of the individual objective functions to produce an optimization problem with a scalar objective. The weights, which are predetermined by the user, need to be carefully chosen to ensure they reflect the relative importance of the individual objectives and also scale them appropriately as necessary. Another approach, *lexicographic optimization*, prioritizes the objective functions based on their importance and optimizes them sequentially, starting with the most important. While optimizing lower priority objectives, we restrict the feasible region to only those solutions that will not degrade the higher priority objectives, or limit their degradation by user-specified tolerances.

The weighted sum approach, where we aggregate the individual objectives into a single objective using user-defined weights, requires a vector of weights that capture the importance of each phenotypic trait in the breeding program. In practice, it is difficult to identify a precise and meaningful weight for each trait as there are many factors of the breeding program (some of them potentially unknown) that might play a role in defining it. By contrast, it might be simpler for a breeding program to order the traits based on their importance.

The lexicographic optimization approach is not without drawbacks, as it could degenerate into single-objective optimization with the highest priority objective if we subsequently allow no degradation of higher priority objectives. In the worst case, if the first objective has a unique optimal solution and we tolerate no degradation on the first objective, the subsequent objectives are irrelevant. The use of tolerance is therefore important as it allows limited degradation of a higher priority objective when optimizing a lower priority objective, but allows for a larger feasible

solution space for the lower priority objective (when compared against using zero tolerance). Hence, we will be using lexicographic maximization with positive tolerances in solving Supplementary Formulation (29).

Assume without loss of generality that the vector of objective functions, $F(t, x) = (f_\ell(t, x^\ell))_{\ell=1}^M$, is already in decreasing order of importance. Thus trait ℓ is more important than trait $\ell + 1$, for each $\ell \in [M - 1]$. The solver we use in our computational studies is capable of lexicographic optimization with degradation tolerances for objectives specified by the decision maker. Let us denote these tolerances by $\tau = (\tau_1, \tau_2, \dots, \tau_M)$, where $\tau_\ell \in [0, 1]$ for each trait $\ell \in [M]$. The solver optimizes the first objective function $f_1(t, x)$ and then, among those feasible solutions within a factor $(1 - \tau_1)$ of the optimal objective value of the first objective function, optimizes the second objective function. This process is repeated until the last objective is optimized. In particular, this method assures that the optimal solution for the ℓ -th objective, for $\ell = 2, \dots, M$, is within a factor $(1 - \tau_\ell)$ of the optimal value of the i -th objective, for every $i \in [\ell - 1]$. As $f_M(t, x^M)$ is the last objective function to optimize, there is no need for a tolerance τ_M , and hence we set $\tau_M = 0$ for all our experiments.

Simulation study

Simulations were conducted to evaluate the performance of ECV compared to other parent selection approaches using phenotypes and breeding values (GEBV). Two simulation experiments were considered in this study. First, we considered a single-trait optimization problem to solely improve Trait 1, simulated as a mixture of traits with oligogenic and polygenic genetic architectures. Next, we examined a multiple-trait parent selection problem where the breeding program was tasked to simultaneously improve all traits of interest. In this experiment, we simulated a polygenic architecture for Trait 3, representing a trait such as yielding capacity that is usually governed by a large number of loci where each allele has a small impact on the expression of the trait and in a negative genetic correlation with Trait 1, in addition to an oligogenic phenotype (Trait 2) that may imitate the genetic architecture underlying disease resistance.

For all experiments, two metrics were reported from the simulations, average desirable allele frequency and average phenotypic trait values of the progeny, to compare the performance of the methods in each generation. We also recorded the average genomic relationship for the selected individuals for all three approaches. In the case of multi-parental pair selection, we sorted pairs of individuals based on the summation of their trait values or GEBVs and made selection decisions based on the summations of trait values. Moreover, by default, there was no control over the genomic relationship between selected parent pairs for the phenotypic selection and GEBV selection approaches; however, we assumed that self-crossing is not a feasible choice in these approaches.

The QU-GENE engine and QuLinePlus proposed by Ali et al. (2020) were used to simulate initial populations and the progeny in the subsequent generations. The QU-GENE engine establishes the initial population with inputs of genetic effects for segregating alleles, recombination frequencies and the number of desired individuals. We considered an initial population such that the allele frequency at all loci was set at 0.5. In our experiments, QuLinePlus took the genotypic information of a population and a list of selected pairs, simulated the progeny by crossing the selected parental pairs, and output genotypic and phenotypic information for all individuals in the subsequent generation. The GEBVs were calculated using the “rrBLUP” package (Endelman 2011). The Gurobi Optimization Solver (Gurobi Optimization, LLC 2024) was used to solve the integer linear programming formulations that were implemented in the Python programming language.

The initial population consisted of 10,000 individuals, with 200 biallelic genetic loci and 100 markers. Of these, 40 genetic loci had effects on Trait 1, 10 on Trait 2, and 70 on Trait 3. The markers had no genetic effects on any of the traits. Trait 3 and Trait 1 share 20 common loci with pleiotropic effects, which resulted in a negative correlation between those phenotypic traits. We conducted all of the experiments for four generations and for each cross we simulated 100 progeny for the next generation. We performed two sets of simulation studies, assuming a consistent growing environment across generations. In the first simulation study, the number of parental pairs that we chose from the initial population, generations one, two, and three, was 50, 10, 5, and 5, respectively. Thus, the population size in the simulation studies for generations one, two, three and four were, respectively, 5000, 1000, 500, and 500, respectively.

To further investigate the effectiveness of our methodology, we explored the impact of selection intensity in our second simulation study.

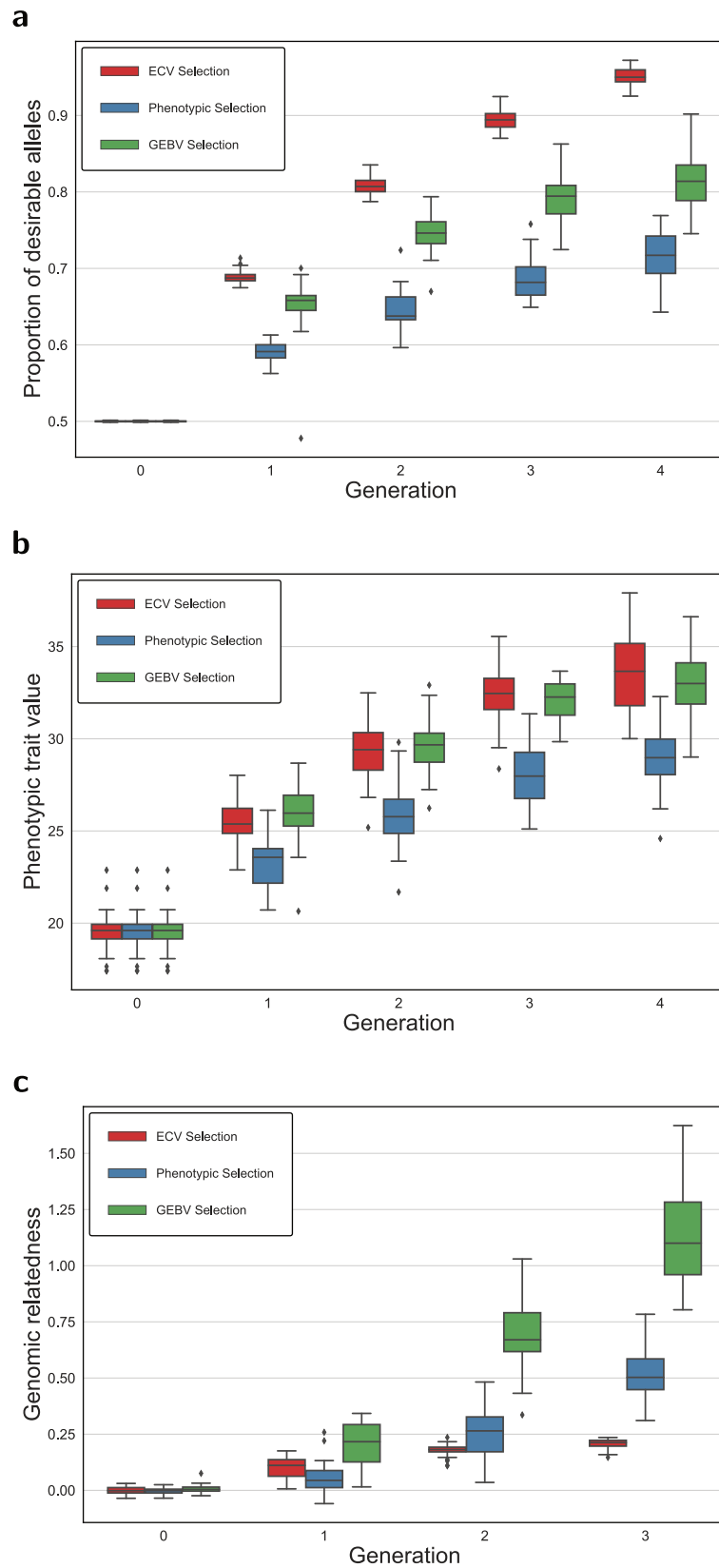


Fig. 2 Performance of ECV, phenotypic, and GEBV selection methods in 30 simulation runs for single trait improvement. **a** Proportion of desirable alleles for Trait 1. **b** Phenotypic value for Trait 1. **c** Genomic relatedness of selected parents.

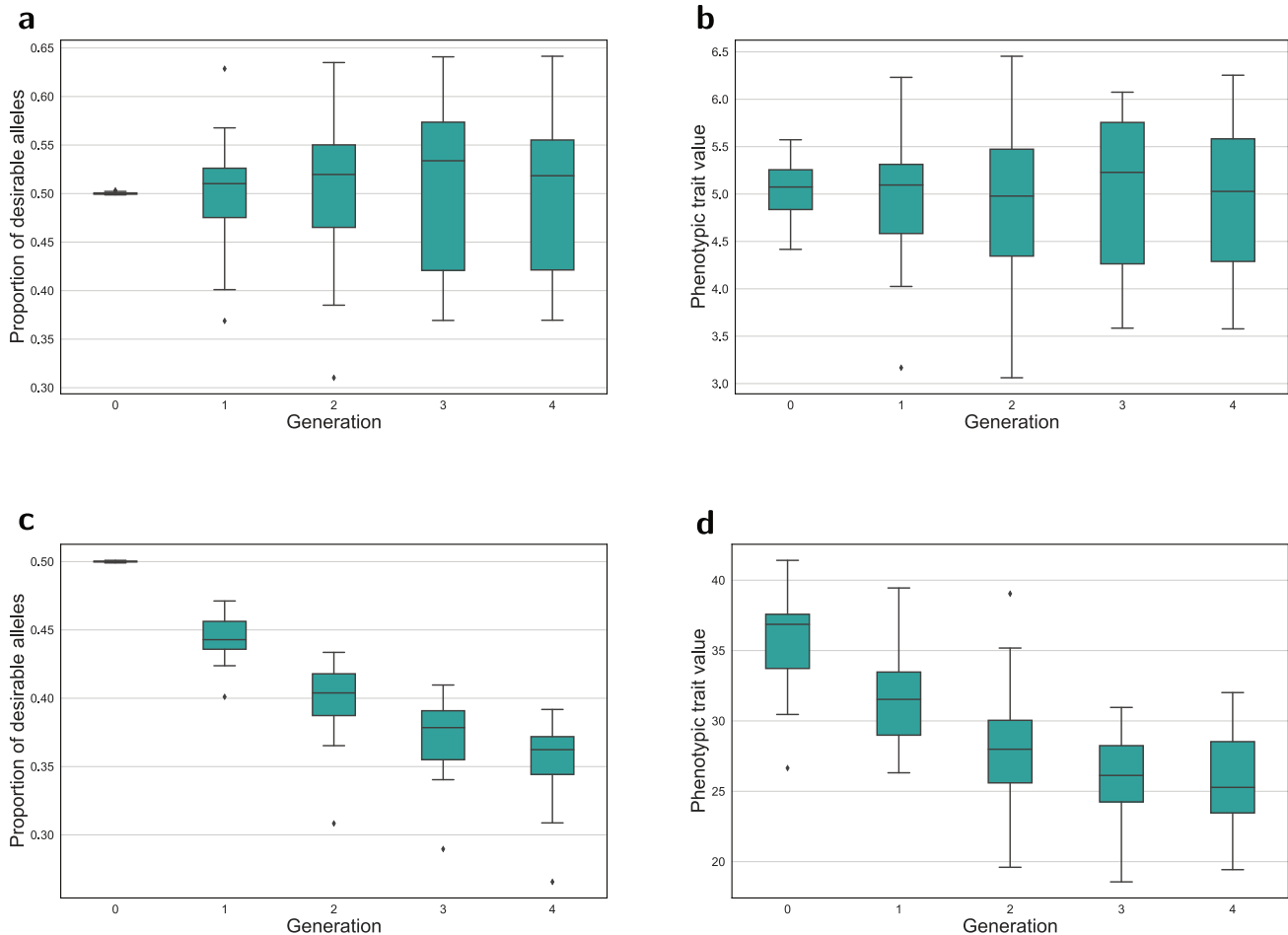


Fig. 3 Effect of improving Trait 1 using ECV optimization on the performance of Trait 2 and Trait 3 based on proportion of desirable alleles and phenotypic trait value. **a** Proportion of desirable alleles for Trait 2. **b** Phenotypic values for Trait 2. **c** Proportion of desirable alleles for Trait 3. **d** Phenotypic values for Trait 3.

Scenario A imposed a higher selection intensity with 50 crosses made from the initial population (generation 0), and 10, 3, and 3, respectively, for generations one, two, and three. For intermediate selection intensity (Scenario B, same as the first simulation study), from the initial population, generations one, two, and three, we chose 50, 10, 5, and 5 parental pairs, respectively. Finally, in Scenario C, representing a case of reduced selection intensity, the numbers of parental pairs selected in generation one was 25, and 5 parental pairs for the generations two and three.

RESULTS

The single-trait simulation results over five generations are summarized in Fig. 2. For all traits considered, ECV significantly increased the proportion of desirable alleles (see Fig. 2a) while showing the capacity to regulate the relatedness within the breeding population by avoiding crossing closely related individuals (see Fig. 2c). Further, although statistically insignificant, genetic crosses done by phenotypically superior individuals returned the lowest means of the progeny in all traits, compared to genetics-informed approaches, like GEBVs and ECV (see Fig. 2b). However, populations generated by ECV provided a greater potential for advancing individuals with larger phenotypic values.

Single-trait optimization does not guarantee improvement for phenotypes other than the target trait. Figure 3 shows boxplots for Trait 2 and Trait 3 when we optimize Trait 1 in a single-trait ECV optimization framework. The frequency for the desirable allele (Fig. 3a) as well as the phenotypic values (Fig. 3b) remained unimproved for Trait 2. The scenario could be worse if target traits

are determined by QTLs with antagonistic pleiotropic effects. This can be seen in Fig. 3c, d, which depict a significant decrease in the proportion of desirable alleles and phenotypic values of Trait 3 as a result of optimizing for Trait 1.

For multi-trait parental selection based on ECV, we employed the lexicographic multi-objective optimization approach described earlier. The tolerances were chosen based on preliminary experiments as follows: let $\tau_{i,c}$ denote the degradation tolerance for optimization objective i in generation c , then we used $\tau_{1,0} = 0.17$, $\tau_{1,1} = 0.05$, $\tau_{1,2} = 0.05$, $\tau_{1,3} = 0.05$ and $\tau_{2,0} = 0.00$, $\tau_{2,1} = 0.00$, $\tau_{2,2} = 0.00$, $\tau_{2,3} = 0.05$. In general, the tolerance parameters can be calibrated to have the desired impact on the model. The results in Fig. 4 show the advantage of using ECV. Despite the negative genetic correlation, ECV was able to increase the desirable allele frequency to $0.70 (\pm 0.02)$, $0.65 (\pm 0.08)$, and $0.72 (\pm 0.01)$, for Trait 1, Trait 2, and Trait 3, respectively. In contrast, the impact of negative correlation between Trait 3 and Trait 1 was most obvious when the phenotypic selection was used, leading to a significant loss of desirable allele for Trait 1 (see Fig. 4a). Similarly, ECV improved phenotypic values of the progeny for all traits simultaneously, whereas no improvement for Trait 1 and Trait 2 was found using phenotypic selection in our simulations when the tolerances were set slightly favoring Trait 3. It is noteworthy that a genomics-informed selection method, GEBV, returned comparable results to ECV for Trait 1. This benefit of GEBV, however, is at the expense of genetic diversity, as shown in Fig. 5. Genomic relatedness (VanRaden 2008) has increased

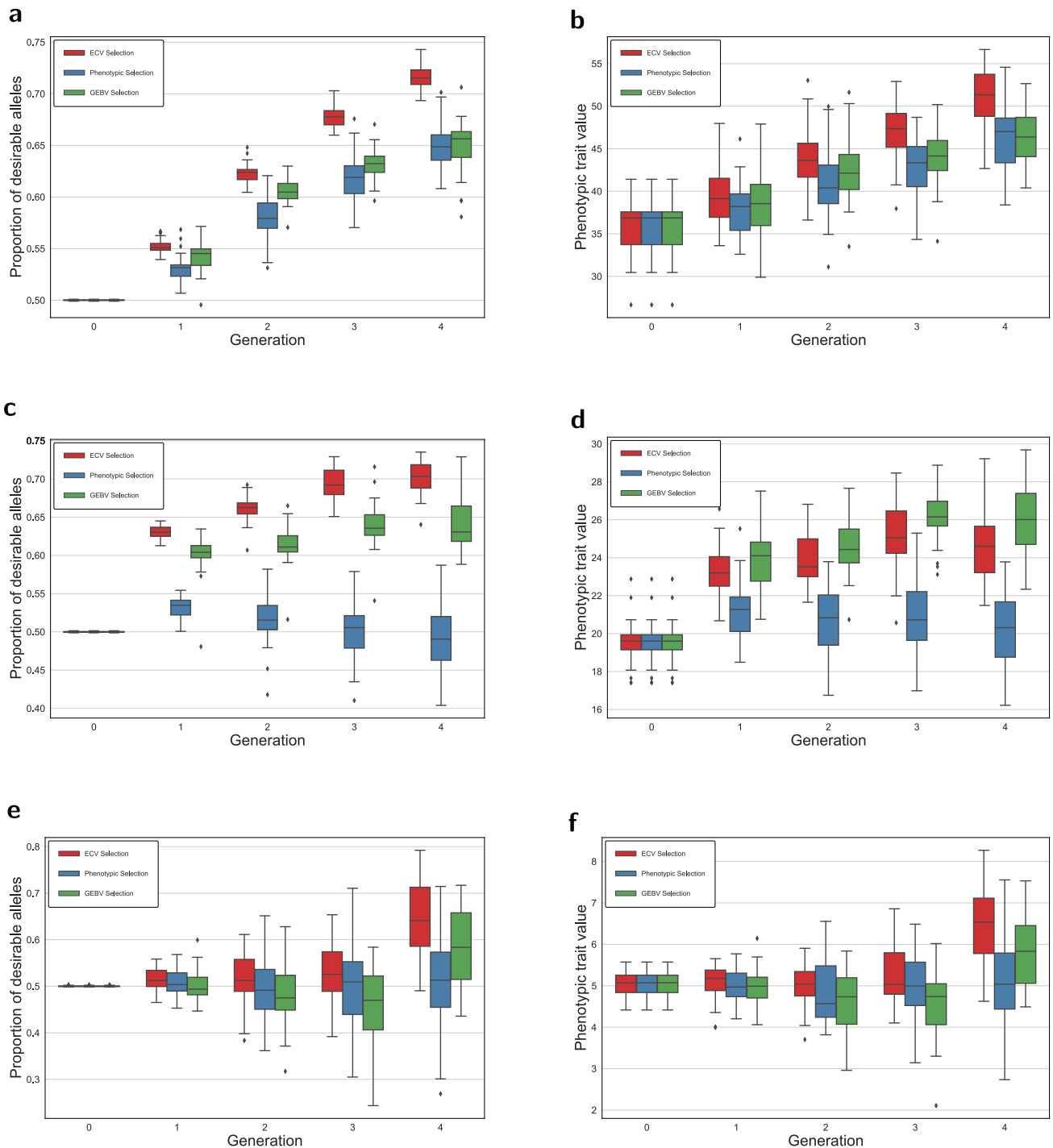


Fig. 4 Performance of ECV, phenotypic, and GEBV selection methods in 30 simulation runs for multiple trait improvement. **a** Proportion of desirable alleles for Trait 3. **b** Phenotypic values for Trait 3. **c** Proportion of desirable alleles for Trait 1. **d** Phenotypic values for Trait 1. **e** Proportion of desirable alleles for Trait 2. **f** Phenotypic values for Trait 2.

noticeably over the four generations using the GEBV selection method.

Figure 6 displays the boxplots for the three selection intensity scenarios A, B, and C introduced in the Simulation study, focusing on the proportion of desirable alleles as the performance metric. In the early generations, particularly generation 2, our proposed ECV approach outperformed other selection strategies, most evidently for Traits 1 and 3 in the implementation of multiple trait selection. As the generations advanced, the ECV approach continued to excel in Scenarios B and C, resulting in higher

proportion of desirable alleles for Traits 1 and 3. In Scenario A, where selection intensity was higher and ECV selection method was not dominant, the method still yielded replications with a greater proportion of desirable alleles compared to other strategies, despite the genomic relationship constraints inherent in the ECV method. Furthermore, in the last generation under scenario A, the mean (\pm standard deviation) of genomic relatedness over all replications for the ECV, GEBV, and Phenotypic selection approaches are $0.15(\pm 0.04)$, $0.42(\pm 0.10)$, and $0.25(\pm 0.10)$, respectively. These results illustrate the effectiveness

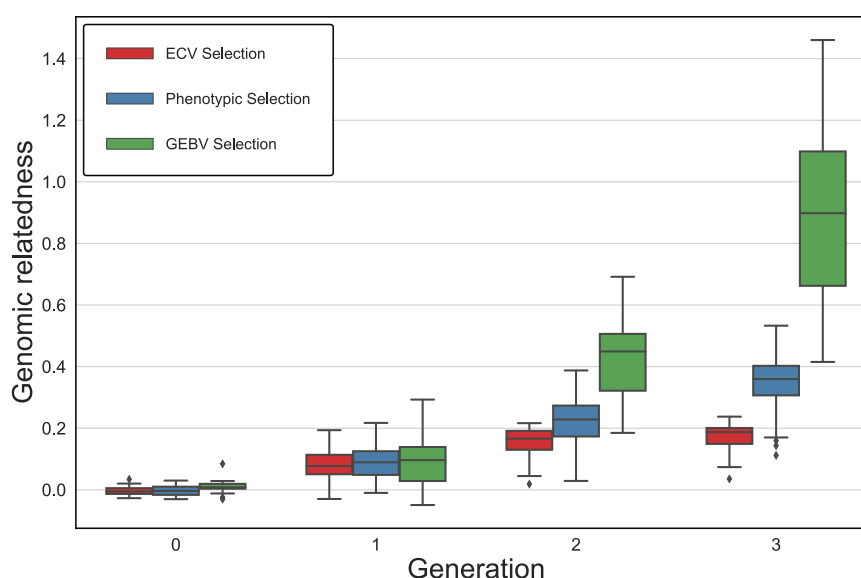


Fig. 5 Genomic relatedness for multiple trait improvement based on ECV, phenotypic and GEBV selection methods.

of our ECV optimization framework, with its explicit constraints limiting genomic relatedness, in managing genomic relatedness over generations when selection intensity is higher, while improving desirable breeding traits. The impact of higher selection intensity, however, led to greater variability in the proportion of desirable alleles in Trait 3 of Scenario A, which could be due to the drift effect of the smaller breeding population size (Turner-Hissong et al. 2020).

DISCUSSION

The principal objective of breeding is to combine as many desirable traits as possible into a genotype that can be distributed to farmers, producers or breeders. For example, in plant breeding the breeders are tasked with developing elite genotypes that display desirable use characteristics including high yields, disease resistance, and are also well-adapted to a range of environmental conditions (Breseghello and Coelho 2013). These desirable characteristics are typically possessed by multiple founders. By mixing and recombining founder genomes, the distribution of these desirable phenotypes observed in the offspring, owing to the segregation of alleles often distributed throughout the genome, allow breeders to identify superior individuals for subsequent breeding, widespread evaluations, and sales.

However, when these desired characteristics differ in variability, heritability, economic importance, and are correlated with other phenotypes and genotypes, effective mating designs capable of improving multiple traits simultaneously can be challenging to identify (Johnson et al. 1988). This breeding process is also ineffective as breeders tend to make hundreds or thousands of crosses, of which only a few are advanced in the subsequent years (Witcombe et al. 2013). Traditionally, these objectives are achieved by breeding from the “best”—the best being determined by their own phenotypic values (Akdemir et al. 2019; Allard 1999). More advanced techniques, such as pedigree-based (Gianola and Fernando 1986; Henderson 1984), marker-based genetic value predictions (Bernardo and Charcosset 2006; Hospital and Charcosset 1997; Lande and Thompson 1990), and mating designs by genomic information (Akdemir and Sánchez, 2016) are also available.

Beginning with the work of Johnson et al. (1988), mathematical programming approaches have facilitated the improvement of genetic traits through the use of mathematical optimization models that aid breeders in making better decisions in selecting

mating parents. Toro et al. (1991) solved mate selection problems using linear programming techniques and demonstrated the effectiveness of their approach within multiple ovulation and embryo transfer (MOET) breeding schemes for dairy cattle with the help of simulation studies. Jansen and Wilton (1985) addressed the issue of factorial growth in the number of combinations to cross by formulating and solving an integer programming model to improve the overall progeny merit.

Moeiniazade et al. (2019) recently proposed a single-trait optimization of a “look ahead” metric that focuses on a predetermined terminal generation to optimize mating decisions for maximizing expected GEBVs in the terminal generation without explicitly considering the impact of genetic erosion. Amini et al. (2021) further improved this look-ahead framework by prioritizing best individuals for crossing and using multiple prediction algorithms to improve prediction accuracy. These approaches are also complemented by Zhang and Wang (2022) who proposed a “net present value” inspired mechanism for discounting future gains, which values early-term genetic gains more than those anticipated in the future. This was done to overcome a drawback of the original look-ahead scheme by Moeiniazade et al. (2019), which can produce slow genetic gains in the early generations and accelerating more rapidly as we approach the terminal generation.

Byrum et al. (2016) report on their long-term development and quantification of an unbiased genetic gain performance metric, and pioneered its use in evaluating breeding projects as varieties were developed. Byrum et al. (2016) and Byrum et al. (2017) demonstrate the successful commercial use of advanced analytics and operations research tools such as integer linear programming, Monte Carlo simulation, and stochastic optimization by the agriculture industry, which has served to further motivate its broader use in many areas of crop and animal sciences; see also (Byrum 2015, 2016).

Furthermore, when the breeding objective involves more than one trait, a selection index of progeny merit was considered as a linear function of estimated breeding values for each trait by Allaire (1980). In animal breeding, for instance, the genetic merit of calves is estimated as half of the sire’s and half of the dam’s breeding value. An optimization-based procedure for mate selection in animal breeding is introduced by Kinghorn (1998, 2011) based on a mate selection criterion proposed by Kinghorn and Shepherd (1999).

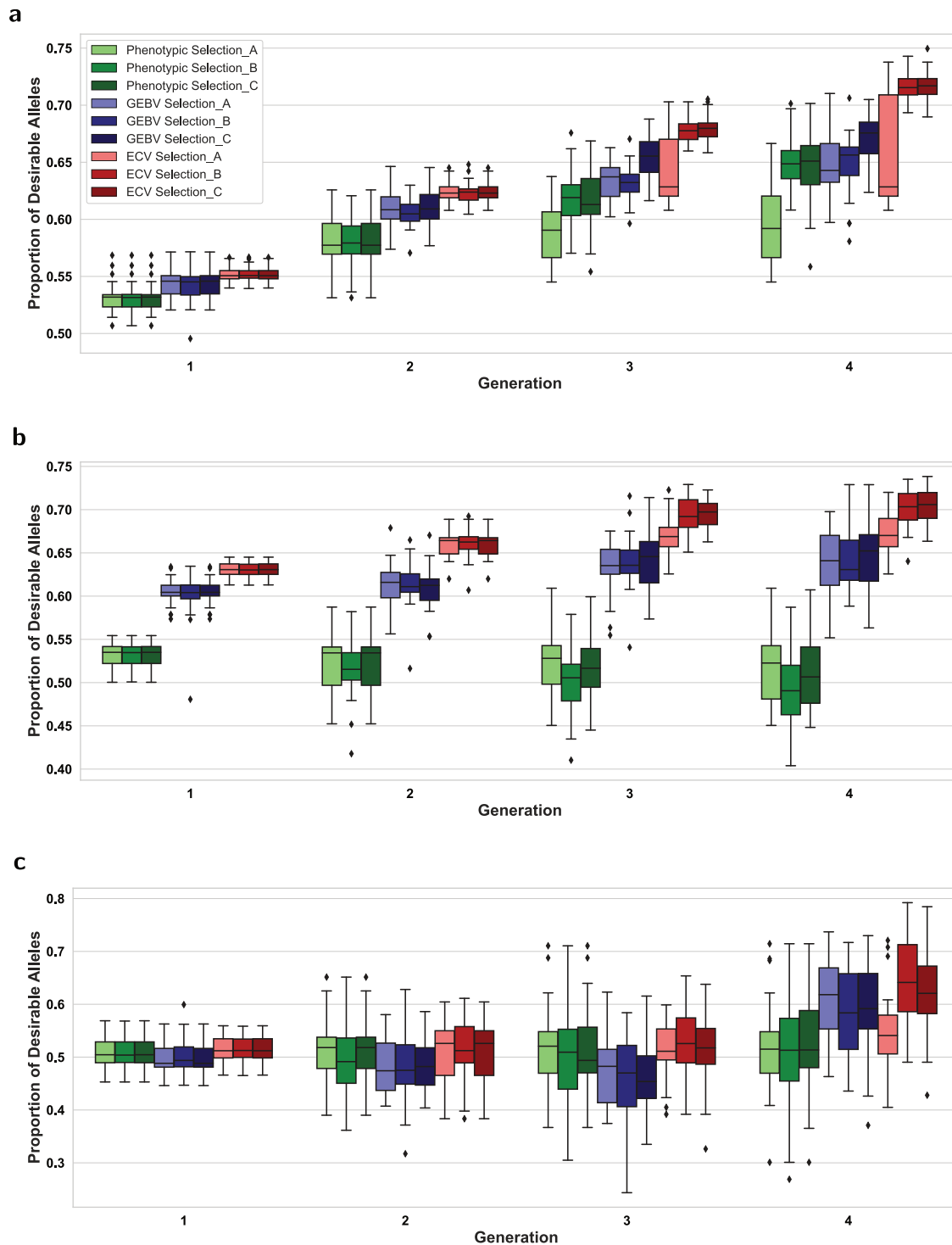


Fig. 6 The effect of selection intensity represented by scenario A (higher intensity), scenario B (intermediate intensity), and scenario C (lower intensity) on the proportion of desirable alleles. **a** Proportion of desirable alleles for Trait 3. **b** Proportion of desirable alleles for Trait 1. **c** Proportion of desirable alleles for Trait 2.

In the genomics era, the parental selection problem has been increasingly addressed with the use of genomic relationships (Sun et al. 2013), heuristic searches for gene pyramiding (De Beukelaer et al. 2015), and by modeling the recombination of desirable alleles as a result of crossing (Han et al. 2017; Moeinizada et al. 2019). For the purpose of introgressing a small number of desirable alleles from a donor to a recipient, Han et al. (2017) proposed an efficient algorithm for calculating the PCV defined as

the probability that a gamete of a random progeny from crossing two genetic individuals would consist only of desirable alleles. In a specific case where the desirable allele for the i -th locus is not present in both parents (denoted by k and k'), such that $L_{i,1}^k = L_{i,2}^k = L_{i,1}^{k'} = L_{i,2}^{k'} = 0$, PCV will conclude that the i -th component of gamete g^3 is zero with probability one and hence $PCV(L^k, L^{k'}, r, a_0) = 0$. In this case, the individuals k and k' will not be selected, regardless whether or not there may be desirable

alleles present in the rest of the genome. While such a result is desirable for the goals of introgressing a small number of alleles for traits like herbicide, disease, or insect resistance, it would be inappropriate for identifying crosses that will have the best opportunity to combine a large number of genetic alleles. Considering the polygenic inheritance of agronomical performance traits (Lynch and Walsh 1998; Scott et al. 2021), the PCV approaches zero for all breeding parents as the number of loci with desirable alleles increases. Consider the following probabilistic inequality (Fréchet inequality):

$$\text{PCV}(L^k, L^{k'}, r, a_0 = 0.5) = \Pr(g_i^3 = 1 \forall i \in [N]) \leq \min_{i \in [N]} \Pr(g_i^3 = 1). \quad (20)$$

Hence, the larger the value of N , the greater the chance $L_{i,1}^k = L_{i,2}^k = L_{i,1}^{k'} = L_{i,2}^{k'} = 0$ for some QTL i . The PCV method could therefore lead to indiscriminate mate selection for traits that have hundreds or thousands of loci with desirable alleles because the PCV value is (nearly) zero for essentially any choice of mates. This observation motivated us to introduce our ECV criterion, especially for breeding targets governed by a large number of genetic loci and for non-introgression projects.

As Fig. 2a shows, our results demonstrated a significantly greater capacity to increase desirable allele frequency compared to the conventional phenotypic selection and the selection done by the genomics-derived GEBV; and, the benefit of using ECV can be realized in as short as two generations. Moreover, the greater range of trait value distribution presents additional opportunities for breeders to identify the superiors for population advancement (see Fig. 2b).

Based on our simulations, we observe that the breeding population has gone from unrelated to essentially full-sibs in three generations of selecting breeding parents based on GEBVs (see Fig. 2c). Compared to the phenotypic selection, GEBV selection might have manifested a rapid increase of relatedness by crossing individuals closely related to the training population (Bassi et al. 2016; Forutan et al. 2018). Though GEBV selection might show a capacity to provide short-term genetic gain, selecting breeding parents solely by GEBVs would lead to undesirable consequences such as loss of genetic diversity, further diminishing long-term genetic gain (Doekes et al. 2018; Jannink 2010).

To ensure the capacity to preserve multiple genetic lineages, ECV allows for the selection of more than one pair of individuals, and while self-crossing was not allowed in this study, our method permitted the same individual to be crossed with multiple breeding parents as long as the genomic relationship of the parents was not greater than ϵ , a parameter that breeders can use to control how much inbreeding is acceptable.

Fundamental to all variety improvement programs is the identification of the most efficient path to reach breeding objectives (Akdemir et al. 2019; Bernardo 2002). However, breeders are usually tasked with combining a suite of traits in addition to yield and growth components. The negative genetic correlations caused by the non-random association of alleles underlying these breeding objectives impose additional challenges, as selecting based on one trait may adversely impact another (Lynch and Walsh 1998). To simultaneously improve multiple traits, phenotype-based selection indices have been widely considered (Hazel et al. 1994; Hazel and Lush 1942; Jannink et al. 2000; Moenizade et al. 2020; Villanueva and Woolliams, 1997). Selecting breeding parents based on a selection index does not necessarily choose the best genetics to recombine; further, since the selection index applied is a weight assignment of target phenotypes, such decisions could result in the loss of beneficial alleles.

In this study, the proposed ECV framework is based on an allele transmission process. Rather than relying on the phenotypes of breeding parents, ECV identifies the crosses with the highest likelihood of transmitting desirable alleles from pairs of parents to

the progeny. In the case where multiple traits need to be considered simultaneously, ECV seeks the optimal combination of alleles for all target phenotypes ordered by their importance, while maintaining a customizable tolerance such that QTLs with antagonistic pleiotropic gene action could remain in consideration before the final breeding recommendation is made. Figures 4 and 6 showed that despite the negative correlation between Traits 1 and 3, ECV was able to increase desirable allele frequency for all traits in our simulation studies. In addition, as seen from Fig. 5, the inbreeding coefficient in the progeny was regulated as ECV was optimized with the tolerance constraint on the genomic relatedness between breeding parents. As genotyping has become routine in breeding programs (Bentley et al. 2022; Hayes and Goddard, 2010), the application of this constraint ought to be considered to mitigate the multiple trait scenario in Fig. 4, where the gain might be built at the expense of genetic diversity (Fig. 5), a phenomenon also found in index selection methods (Akdemir et al. 2019). If practical considerations favor breeding parents to be selected from a narrow genetic pool, the constraint could be moved to the objective as a penalty term.

Breeding programs develop elite genotypes that often demonstrate similar essential genomic profiles of desirable end-use characteristics, agronomical attributes, disease resistance packages, as well as adaptation to the target environment. Breeding among the elites can produce new variability as the source of new cultivars with minimal risk of introducing undesirable features. This variation may eventually be exhausted, and new genes and alleles must be introduced. Identifying beneficial alleles from un-adapted material itself has been described as searching for a needle in a haystack (Pixley et al. 2014). Introgressing these novel alleles can also be risky because the unwanted alleles in exotic germplasm may disrupt essential allele combinations (Willcox et al. 2022); and, it requires a higher institutional cost due to a greater number of crosses and longer breeding cycle needed to achieve the breeding objectives (Neyhart et al. 2019; Snelling et al. 2019). Based on our simulations, we reckon that ECV can be an option.

Beyond animal and cereal crop breeding, we suspect that implementing optimization-based methods like ECV could be advantageous to breeding of genetically diverse, long-generation, and slow reaction, cross-pollinated species, such as conifers. Tree breeders generally establish open-pollinated seed orchards for selection (White et al. 2007), and several mating designs have been proposed (Namkoong, 1976; Zobel and Talbert 1984), among which the polycross is considered as one of the most cost-effective (Kumar et al. 2007; Lenz et al. 2020). The ability to design the pollen pool while managing inbreeding with ECV will provide the capacity to rapidly increase desirable allele frequencies and, at the same time, avoid severe inbreeding depression for conifer species (Berry and Evans 2014; Mwansa et al. 2002; Snelling et al. 2019).

The conceptual framework we have introduced and our results show that adopting multi-objective optimization tools from operations research to solve breeding problems is highly advantageous (Beans 2020; Cameron et al. 2017; Kusmec et al. 2021). Several improvements or extensions suggested next should also be considered. When the pool of genetic diversity increases, solving integer programming problems for ECV will require further development to account for different distributions of crossover events in different crosses (Dreissig et al. 2019; Jabbari et al. 2019; Nachman 2002; Stapley et al. 2017). Furthermore, as multi-parent populations like MAGIC (multi-parent advanced generation inter-cross) have become a means to provide germplasm for breeding programs (Scott et al. 2020), there is a need to expand optimization frameworks such as ECV to consider multiple parental lineages, which might also help guide the polycross mating design in forestry (Frandsen 1940; Lambeth et al. 2001).

Our proposed methodology relies on the the underlying genetic information of the breeding population, such as QTLs and genetic association of desirable traits. While affordable large-scale genotyping and phenotyping technologies are becoming accessible to breeding programs (Bassi et al. 2024; Reynolds et al. 2020), large breeding populations necessitate extensive genomic information, which can be computationally demanding. Moreover, integer linear programming is NP-hard in general, making it challenging to solve very large-scale problems to optimality. In the case of mate selection, the size of the population and the number of genes directly influence the computational time, which implies that massive datasets could make obtaining optimal solutions unrealistic for practical applications. In such circumstances, we may consider modifying our approach to solving the integer linear program by employing decomposition techniques to address the large-scale instances and likely settle for sub-optimal (but good quality) feasible solutions.

Our simulations also indicate that the ECV selection framework may result in higher performance variability when the selection intensifies in earlier generations. To alleviate this issue, we could either relax the genomic relatedness constraint (smaller ϵ) in earlier generations or intensify selection only in advanced generations. Further, care must also be taken in choosing the degradation tolerances τ_ℓ for each trait $\ell \in [M]$ in the lexicographic multi-trait ECV optimization framework, which will entail computational expenditure in terms of preliminary computational experiments, which could become challenging at larger scales.

DATA AVAILABILITY

The data and codes are available at: <https://github.com/transgenomicsosu/ECV>.

REFERENCES

- Akdemir D, Beavis W, Fritsche-Neto R, Singh AK, Isidro-Sánchez J (2019) Multi-objective optimized genomic breeding strategies for sustainable food improvement. *Heredity* 122(5):672–683
- Akdemir D, Sánchez JI (2016) Efficient breeding by genomic mating. *Front Genet* 7:210
- Ali M, Zhang L, DeLacy I, Arief V, Dieters M, Pfeiffer WH (2020) Modeling and simulation of recurrent phenotypic and genomic selections in plant breeding under the presence of epistasis. *Crop J* 8(5):866–877
- Allaire F (1980) Mate selection by selection index theory. *Theor Appl Genet* 57(6):267–272
- Allard RW (1999) *Principles of plant breeding*. John Wiley & Sons
- Amini F, Franco FR, Hu G, Wang L (2021) The look ahead trace back optimizer for genomic selection under transparent and opaque simulators. *Sci Rep* 11(1):4124
- Bassi FM, Bentley AR, Charmet G, Ortiz R, Crossa J (2016) Breeding schemes for the implementation of genomic selection in wheat (*triticum* spp.). *Plant Sci* 242:23–36
- Bassi FM, Sanchez-Garcia M, Ortiz R (2024) What plant breeding may (and may not) look like in 2050? *Plant Genome* 17(1):e20368
- Beans C (2020) Inner workings: Crop researchers harness artificial intelligence to breed crops for the changing climate. *Proc Natl Acad Sci* 117(44):27066–27069
- Bentley A, Chen C, D'Agostino N (2022) Genome wide association studies and genomic selection for crop improvement in the era of big data. *Front Genet* 13:873060
- Bernardo R (2002) *Breeding for quantitative traits in plants*. Stemma Press, Woodbury, Minnesota, USA
- Bernardo R, Charcosset A (2006) Usefulness of gene information in marker-assisted recurrent selection: A simulation appraisal. *Crop Sci* 46(2):614–621
- Berry DP, Evans R (2014) Genetics of reproductive performance in seasonal calving beef cows and its association with performance traits. *J Anim Sci* 92(4):1412–1422
- Breseghele F, Coelho ASG (2013) Traditional and modern plant breeding methods with examples in rice (*oryza sativa* L.). *J Agric Food Chem* 61(35):8277–8286
- Byrum J (2015) Agriculture: Fertile ground for analytics and innovation. *OR/MS Today* 42(6):28–32
- Byrum J (2016) Optimizing crop management: “Smart” application of fertilizer illustrates payoff in using analytical tools to enhance crop yields and improve the environment. *OR/MS Today* 43(3):26–30
- Byrum J, Beavis B, Davis C, Doonan G, Doubler T, Kaster V (2017) Genetic gain performance metric accelerates agricultural productivity. *Interfaces* 47(5):442–453
- Byrum J, Davis C, Doonan G, Doubler T, Foster D, Luzzi B (2016) Advanced analytics for agricultural product development. *Interfaces* 46(1):5–17
- Cameron JN, Han Y, Wang L, Beavis WD (2017) Systematic design for trait introgression projects. *Theor Appl Genet* 130(10):1993–2004
- Cassandro M (2020) Animal breeding and climate change, mitigation and adaptation. *J Anim Breed Genet* 137(2):121–122
- De Beukelaer H, De Meyer G, Fack V (2015) Heuristic exploitation of genetic structure in marker-assisted gene pyramiding problems. *BMC Genet* 16(1):1–16
- Doekes HP, Veerkamp RF, Bijma P, Hiemstra SJ, Windig JJ (2018) Trends in genome-wide and region-specific genetic diversity in the Dutch-Flemish Holstein-Friesian breeding program from 1986 to 2015. *Genet Select Evolut* 50(1):1–16
- Donald CT (1968) The breeding of crop ideotypes. *Euphytica* 17(3):385–403
- Dreissig S, Mascher M, Heckmann S (2019) Variation in recombination rate is shaped by domestication and environmental conditions in barley. *Mol Biol Evolution* 36(9):2029–2039
- Endelman JB (2011) Ridge regression and other kernels for genomic selection with R package rrBLUP. *Plant Genome* 4(3):250–255
- Forutan M, Ansari Mahyari S, Baes C, Melzer N, Schenkel FS, Sargolzaei M (2018) Inbreeding and runs of homozygosity before and after genomic selection in north american holstein cattle. *BMC Genomics* 19(1):1–12
- Frandsen H (1940) Some breeding experiments with timothy. *Imp Agr Bur Jt Pub* 3:80–92
- Gianola D, Fernando RL (1986) Bayesian methods in animal breeding theory. *J Anim Sci* 63(1):217–244
- Gurobi Optimization, LLC Gurobi optimizer reference manual. <https://www.gurobi.com>. Accessed 26 May 2024 (2024)
- Han Y, Cameron JN, Wang L, Beavis WD (2017) The predicted cross value for genetic introgression of multiple alleles. *Genetics* 205(4):1409–1423
- Hayes B, Goddard M (2010) Genome-wide association and genomic selection in animal breeding. *Genome* 53(11):876–883
- Hayes BJ, Bowman PJ, Chamberlain AJ, Goddard ME (2009) Invited review: genomic selection in dairy cattle: progress and challenges. *J Dairy Sci* 92(2):433–443
- Hazel L, Dickerson G, Freeman A (1994) The selection index—then, now, and for the future. *J Dairy Sci* 77(10):3236–3251
- Hazel L, Lush JL (1942) The efficiency of three methods of selection. *J Heredity* 33(11):393–399
- Heffner EL, Lorenz AJ, Jannink J-L, Sorrells ME (2010) Plant breeding with genomic selection: gain per unit time and cost. *Crop Sci* 50(5):1681–1690
- Henderson CR (1984) *Applications of linear models in animal breeding*. University of Guelph, Guelph, ON, Canada
- Hesser LF (2006) *The man who fed the world: Nobel Peace Prize laureate Norman Borlaug and his battle to end world hunger: An authorized biography*. Leon Hesser
- Hospital F, Charcosset A (1997) Marker-assisted introgression of quantitative trait loci. *Genetics* 147(3):1469–1485
- Jabbari K, Wirtz J, Rauscher M, Wiehe T (2019) A common genomic code for chromatin architecture and recombination landscape. *PLoS One* 14(3):e0213278
- Jannink J-L (2010) Dynamics of long-term genomic selection. *Genet Select Evolut* 42(1):1–11
- Jannink J-L, Orf J, Jordan N, Shaw R (2000) Index selection for weed suppressive ability in soybean. *Crop Sci* 40(4):1087–1094
- Jansen G, Wilton J (1985) Selecting mating pairs with linear programming techniques. *J Dairy Sci* 68(5):1302–1305
- Johnson BE, Dauer JP, Gardner CO (1988) A model for determining weights of traits in simultaneous multitrait selection. *Appl Math Model* 12(6):556–564
- Kinghorn BP (1998) Mate selection by groups. *J Dairy Sci* 81:55–63
- Kinghorn BP (2011) An algorithm for efficient constrained mate selection. *Genet Select Evolut* 43(1):1–9
- Kinghorn BP, Shepherd RK (1999) Mate selection for the tactical implementation of breeding programs. *Proc Assoc Adv Anim Breed Genet* 13:130–133
- Kumar S, Gerber S, Richardson T, Gea L (2007) Testing for unequal paternal contributions using nuclear and chloroplast ssr markers in polycross families of radiata pine. *Tree Genet Genomes* 3(3):207–214
- Kusmec A, Zheng Z, Archontoulis S, Ganapathysubramanian B, Hu G, Wang L (2021) Interdisciplinary strategies to enable data-driven plant breeding in a changing climate. *One Earth* 4(3):372–383

- Lambeth C, Lee B-C, O'Malley D, Wheeler N (2001) Polymix breeding with parental analysis of progeny: an alternative to full-sib breeding and testing. *Theor Appl Genet* 103(6):930–943
- Lande R, Thompson R (1990) Efficiency of marker-assisted selection in the improvement of quantitative traits. *Genetics* 124(3):743–756
- Lenz P, Nadeau S, Azaiez A, Gérard S, Deslauriers M, Perron M (2020) Genomic prediction for hastening and improving efficiency of forward selection in conifer polycross mating designs: an example from white spruce. *Heredity* 124(4):562–578
- Lynch M, Walsh B (1998) *Genetics and analysis of quantitative traits*. Sinauer Sunderland, MA
- Meuwissen TH, Hayes BJ, Goddard M (2001) Prediction of total genetic value using genome-wide dense marker maps. *Genetics* 157(4):1819–1829
- Miettinen K (2012) *Nonlinear multiobjective optimization*, volume 12. Springer Science & Business Media
- Miettinen K, Hakanen J, Podkopaev D (2016). Interactive nonlinear multiobjective optimization methods. In Greco, S., Ehrgott, M., and Figueira, J. R., editors, *Multiple criteria decision analysis: State of the art surveys*, pages 927–976. Springer New York, New York, NY (2016)
- Moeiniazade S, Hu G, Wang L, Schnable PS (2019) Optimizing selection and mating in genomic selection with a look-ahead approach: An operations research framework. *G3: Genes, Genomes, Genet* 9(7):2123–2133
- Moeiniazade S, Kusmec A, Hu G, Wang L, Schnable PS (2020) Multi-trait genomic selection methods for crop improvement. *Genetics* 215(4):931–945
- Mwansa P, Crews Jr D, Wilton J, Kemp R (2002) Multiple trait selection for maternal productivity in beef cattle. *J Anim Breed Genet* 119(6):391–399
- Nachman MW (2002) Variation in recombination rate across the genome: evidence and implications. *Curr Opin Genet Dev* 12(6):657–663
- Namkoong G (1976) A multiple-index selection strategy. *Silvae Genet* 25:5–6
- Neyhart JL, Lorenz AJ, Smith KP (2019) Multi-trait improvement by predicting genetic correlations in breeding crosses. *G3: Genes, Genomes, Genet* 9(10):3153–3165
- Oldenbroek K, van der Waaij L (2015) Textbook animal breeding and genetics for bsc students. *Centre for Genetic Resources The Netherlands and Animal Breeding and Genomics Centre*, page 245
- Pixley K, Hearne S, Willcox M et al (2014). Seeds of discovery: characterizing and utilizing maize genetic resources for germplasm diversification. *Maize for Food, Feed, Nutrition and Environmental Security*, page 61
- Reynolds M, Chapman S, Crespo-Herrera L, Molero G, Mondal S, Pequeno DN (2020) Breeder friendly phenotyping. *Plant Sci* 295:110396
- Riley R, Chapman V (1958) Genetic control of the cytologically diploid behaviour of hexaploid wheat. *Nature* 182(4637):713–715
- Rincent R, Charcosset A, Moreau L (2017) Predicting genomic selection efficiency to optimize calibration set and to assess prediction accuracy in highly structured populations. *Theor Appl Genet* 130(11):2231–2247
- Sawaragi Y, Nakayama H, Tanino, T (1985) *Theory of multiobjective optimization*. Elsevier
- Scott MF, Fradgley N, Bentley AR, Brabbs T, Corke F, Gardner KA (2021) Limited haplotype diversity underlies polygenic trait architecture across 70 years of wheat breeding. *Genome Biol* 22(1):1–30
- Scott MF, Ladejobi O, Amer S, Bentley AR, Biernaskie J, Boden SA (2020) Multi-parent populations in crops: a toolbox integrating genomics and genetic mapping with breeding. *Heredity* 125(6):396–416
- Simmonds NW (1995) The relation between yield and protein in cereal grain. *J Sci Food Agriculture* 67(3):309–315
- Snelling WM, Kuehn LA, Thallman RM, Bennett GL, Golden BL (2019) Genetic correlations among weight and cumulative productivity of crossbred beef cows. *J Anim Sci* 97(1):63–77
- Stapley J, Feulner PG, Johnston SE, Santure AW, Smadja CM (2017) Variation in recombination frequency and distribution across eukaryotes: patterns and processes. *Philos Trans R Soc B: Biol Sci* 372(1736):20160455
- Sun C, VanRaden P, O'Connell J, Weigel K, Gianola D (2013) Mating programs including genomic relationships and dominance effects. *J Dairy Sci* 96(12):8014–8023
- Toro M, Silió L, Pérez-Enciso M (1991) A note on the use of mate selection in closed moet breeding schemes. *Anim Sci* 53(3):403–406
- Turner-Hissong SD, Mabry ME, Beissinger TM, Ross-Ibarra J, Pires JC (2020) Evolutionary insights into plant breeding. *Curr Opin Plant Biol* 54:93–100
- VanRaden PM (2008) Efficient methods to compute genomic predictions. *J Dairy Sci* 91(11):4414–4423
- Velu G, Singh RP (2013) Phenotyping in wheat breeding. *Phenotyping for plant breeding: applications of phenotyping methods for crop improvement*, pages 41–71
- Villanueva B, Woolliams J (1997) Optimization of breeding programmes under index selection and constrained inbreeding. *Genet Res* 69(2):145–158

- Visscher PM, Haley CS, Thompson R (1996) Marker-assisted introgression in backcross breeding programs. *Genetics* 144(4):1923–1932
- White TL, Adams WT, Neale DB (2007) *Forest genetics*. CABI
- Willcox MC, Burgueño JA, Jeffers D et al (2022) Mining alleles for tar spot complex resistance from cimmyt's maize germplasm bank. *Frontiers in Sustainable Food Systems*, page 297 (2022)
- Witcombe JR, Gyawali S, Subedi M, Virk DS, Joshi KD (2013) Plant breeding can be made more efficient by having fewer, better crosses. *BMC Plant Biol* 13(1):1–12
- Zhang Z, Wang L (2022) A look-ahead approach to maximizing present value of genetic gains in genomic selection. *G3: Genes, Genomes, Genet* 12(8):jkac136
- Zobel B, Talbert J (1984) *Applied forest tree improvement*. Wiley New York

ACKNOWLEDGEMENTS

Funding for this work was supported by grants from the Oklahoma Wheat Research Foundation (for CC), Oklahoma Center for the Advancement of Science and Technology (OCAST) award number PS15-011-2 and PS19-004 for CC. This work was completed utilizing the High-Performance Computing Center facilities of Oklahoma State University at Stillwater, and also in part by the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National Science Foundation grant number ACI-1548562. Specifically, it used the Bridges system, which is supported by NSF award number ACI-1445606, at the Pittsburgh Supercomputing Center (PSC) under the resource allocation MCB-180177. The authors are grateful to the anonymous reviewers for their careful reading of our original manuscript, their constructive criticism, and providing detailed and thoughtful comments that helped us improve this manuscript.

AUTHOR CONTRIBUTIONS

BB, JB, and CC were responsible for the conceptualization of the study. PA developed the theoretical results and computer implementations as part of his thesis. PA and CC performed the analysis and wrote the original draft. All authors contributed to interpreting results, providing feedback, and editing and approving the manuscript.

COMPETING INTERESTS

The authors declare no competing interests.

RESEARCH ETHICS STATEMENT

No approval of research ethics committees was required because no experimental work was conducted; only computer simulations were used in this study.

ADDITIONAL INFORMATION

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41437-024-00697-y>.

Correspondence and requests for materials should be addressed to Charles Chen.

Reprints and permission information is available at <http://www.nature.com/reprints>

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



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