

1 **AlphaSimR: An R-package for Breeding Program**
2 **Simulations**

3

4 R. Chris Gaynor*, Gregor Gorjanc, John M. Hickey

5

6 The Roslin Institute and Royal (Dick) School of Veterinary Studies, University of

7 Edinburgh, Easter Bush Research Centre, Midlothian EH25 9RG, UK

8 Received _____*Corresponding author (chris.gaynor@roslin.ed.ac.uk)

9

10 **Abbreviations:**

11 **Abstract**

12 This paper introduces AlphaSimR, an R package for stochastic simulations of
13 plant and animal breeding programs. AlphaSimR is a highly flexible software package
14 able to simulate a wide range of plant and animal breeding programs for diploid and
15 autopolyploid species. AlphaSimR is ideal for testing the overall strategy and detailed
16 design of breeding programs. AlphaSimR utilizes a scripting approach to building
17 simulations that is particularly well suited for modeling highly complex breeding
18 programs, such as commercial breeding programs. The primary benefit of this scripting
19 approach is that it frees users from preset breeding program designs and allows them to
20 model nearly any breeding program design. This paper lists the main features of
21 AlphaSimR and provides a brief example simulation to show how to use the software.

22

23

24 Introduction

25 This paper introduces AlphaSimR, an R package for stochastic simulations of
26 plant and animal breeding programs. Stochastic simulation is a powerful tool for design
27 and optimization of breeding programs, because it provides a fast, inexpensive method
28 for testing alternative breeding program designs. Simulations have been used to
29 improve both plant breeding programs (e.g.; Lin et al. 2016; Gaynor et al. 2017;
30 Gorjanc et al. 2018) and animal breeding programs (e.g.; Hayes and Goddard 2003;
31 Jenko et al. 2015; Johnsson et al. 2019) as well as to address theoretical concepts in
32 quantitative genetics and breeding (e.g., Gorjanc *et al.* 2015). AlphaSimR has been
33 specifically designed to make simulations more common by providing an easy-to-use
34 and highly flexible software package able to simulate a wide range of plant and animal
35 breeding programs.

36 Stochastic simulations have rarely, if ever, been used to improve breeding
37 programs for many agriculturally important species. This is likely due to the difficulty
38 in setting up and running these simulations. This difficulty is in no small part due to the
39 need for a person with thorough knowledge of breeding programs and computer
40 programming. This person must possess a thorough understanding of the breeding
41 programs they wish to simulate so that they can construct an informative simulation.
42 They must also possess the programming skills needed to develop, run, and evaluate
43 the simulations. The amount of programming skills this person needs to possess is
44 considerable when there are not existing software programs for modeling the specific
45 breeding program of interest. To address this issue, new software is needed that can
46 lower the programming burden and thereby increase the ease of running simulations.

47 AlphaSimR has been specifically designed to make running stochastic
48 simulations of whole breeding programs easier. To accomplish this goal, AlphaSimR
49 provides the ability to run simulations both interactively or via scripts within the R
50 software environment (R Core Team 2019). More specifically, AlphaSimR provides
51 users with a range of R functions that correspond to common operations in a breeding
52 program, such as crossing and selection. This allows users to apply functions
53 representing breeding operations directly to objects that represent populations of
54 animals or plants. The benefit of this approach is that it makes writing simulation code
55 more intuitive, by allowing users to directly translate a description of a breeding
56 program into an R script. It also provides a natural path for learning how to use the
57 software by allowing users to start with simulations of simple breeding programs and
58 gradually progress to more complicated breeding programs. With the respect to
59 learning, simulations are also an invaluable tool to teach students and new professionals
60 about theoretical and practical breeding concepts.

61 AlphaSimR is suitable for simulating a wide range of breeding programs and
62 species. The software models the genomes of both diploid and autopolyploid species.
63 The scripting approach employed by AlphaSimR allows for modeling nearly any
64 breeding program structure, without limiting users to preset designs. AlphaSimR has
65 been heavily optimized for large scale simulations (>1,000,000 individuals), because it
66 is specifically designed for whole breeding program simulations.

67

68 **Methods**

69 AlphaSimR is a large package with an extensive list of features, so we will only
70 describe its main features here. For the sake of brevity, these descriptions are designed
71 to provide an overview of AlphaSimR's functionality and not a detailed accounting of
72 its implementation. First, AlphaSimR's approach to stochastic simulations will be
73 given to provide a high-level overview of how the software works. Then, we will
74 describe a few key elements of this approach before concluding with an overview of
75 AlphaSimR's implementation.

76 *Simulation approach*

77 AlphaSimR uses a simulation approach that combines the coalescent and gene
78 drop methods (Hickey and Gorjanc 2012). The coalescent method is used for
79 backwards-in-time simulations. It is used in AlphaSimR to generate whole-
80 chromosome founder haplotypes. The gene drop method is used for forwards-in-time
81 simulations. It is used in AlphaSimR to create new haplotypes from the original founder
82 haplotypes.

83 *Founder haplotypes*

84 The preferred method for creating founder haplotypes in AlphaSimR is to use
85 the Markovian Coalescent Simulator (MaCS; Chen *et al.* 2009). MaCS is included in
86 AlphaSimR and used to generate founder haplotypes according to either a predefined
87 parameter set for some species, or user supplied parameters. Alternatively, users can
88 create founder haplotypes by importing external data into AlphaSimR or using built-in
89 functions for random sampling. The option to import external data allows users to use
90 other coalescent simulators or real genotypic data, provided the linkage phase and
91 genetic map are known.

92 *Genetic recombination*

93 AlphaSimR creates new haplotypes by modeling genetic recombination during
94 meiosis. A genetic map is used to model the distribution of genetic recombination.
95 AlphaSimR allows for sex-specific genetic maps to represent different recombination
96 rates between sexes. The specifics for modeling meiosis in AlphaSimR depend on
97 whether the species is a diploid or an autopolyploid.

98 For diploid species, AlphaSimR simulates meiosis and genetic recombination
99 according to the gamma model (McPeck and Speed 1995). The gamma model
100 accommodates crossover interference and has been shown to fit real data (e.g. Broman
101 and Weber 2000). The magnitude of crossover interference is controlled by a single
102 parameter that can be adjusted by the user.

103 For autopolyploid species, AlphaSimR simulates meiosis using a combination
104 of bivalent and quadrivalent chromosome pairing. Bivalent or quadrivalent
105 homologous pairs are chosen at random according to a parameter for the probability of
106 quadrivalent pairing, which can be tuned by the user. Bivalent pairs are resolved using
107 the gamma model for diploids. Quadrivalent pairs are resolved according to the model
108 for “cross-type” configurations used in the PedigreeSim software (Voorrips and
109 Maliepaard 2012). This model involves sampling chiasmata positions from a gamma
110 distribution and resolving crossovers by sampling centromeres and working outwards
111 towards the telomeres. This technique models unique features of meiosis in
112 autopolyploids, such as recombinant chromosomes composed of three parental
113 chromosomes and double reductions (Bourke *et al.* 2015).

114 **Traits**

115 Traits in AlphaSimR are classified according to the biological effects they
116 model. The biological effects modeled in AlphaSimR are: **Additive**, **Dominance**,
117 **Epistatic**, and **Genotype-by-environment**. The first letter of each effect is used to derive
118 a name for each trait type under the **ADEG** framework. For example, a trait with only
119 additive effects is called an **A** trait. A trait with additive and dominance effects is called
120 an **AD** trait. AlphaSimR currently supports the following trait types: **A**, **AD**, **AE**, **AG**,
121 **ADE**, **ADG**, **AEG**, and **ADEG**.

122 The modeling of biological effects is based on classic quantitative genetics
123 models. For example, the additive effects are equivalent to additive effects described in
124 a quantitative genetics textbook (e.g. Falconer and Mackay 1996). The modeling of the
125 dominance effects allows for both directional dominance and a variable degree of
126 dominance, ranging from partial dominance to overdominance (Gaynor *et al.* 2018).
127 For autopolyploid species, the modeling of dominance represents digenic dominance.
128 Epistatic effects are modeled as additive-by-additive epistatic effects between discrete
129 pairs of loci. Genotype-by-environment effects are modeled as additive effects whose
130 value is a function of a single environmental covariate.

131 AlphaSimR can simulate multiple traits using any combination of trait types.
132 Each trait is simulated according to a user-defined number of QTL, which can differ
133 between traits. Correlated traits can be simulated, provided they are pleiotropic and
134 belong to the same trait type.

135 AlphaSimR uses a method for sampling QTL effects that is, to the authors'
136 knowledge, unique among stochastic simulation software. Users of AlphaSimR are
137 asked to specify a desired mean and variance, either total or additive, for each trait. The

138 software then samples QTL effects and scales the values for those effects to achieve
139 precisely this mean and variance in a founder population. The benefit of AlphaSimR's
140 approach is that it allows users to set variables relating to the relative levels of
141 dominance or epistasis independently of the founder population's genetic variance. For
142 example, a user can specify the average degree of dominance for QTL controlling a
143 trait independently of the additive genetic variance for this trait.

144 *Variance components*

145 AlphaSimR reports additive, dominance and additive-by-additive epistatic
146 variances for any population. This is done without assuming random mating or linkage
147 equilibrium, so that the values are correct regardless of the population's genetic
148 structure. This allows users to compare simulated populations to real-world data for the
149 sake of benchmarking simulations. AlphaSimR also offers further partitioning of
150 genetic variance into genic variance, covariance due to departures from Hardy-
151 Weinberg equilibrium and covariance due to linkage disequilibrium, as described by
152 Bulmer (1976).

153 *Selection*

154 A wide range of functions are available for modeling selections. These functions
155 allow for selection on multitude of criteria, such as: phenotypes, genetic values,
156 breeding values, or estimated breeding values. Selection can be on one trait or an index
157 of multiple traits. Selections can also be modeled as selection between or within
158 families or over an entire population. AlphaSimR also supports user supplied
159 selections, allowing users to implement their own selection methods, for example
160 optimum contribution selection as in Gorjanc et al. (2018)

161 *Mating and propagation schemes*

162 A wide range of functions are available in AlphaSimR for modeling common
163 mating and propagation schemes. These schemes include: biparental crossing, selfing,
164 clonal propagation, generation of doubled haploid lines, and propagation in open
165 pollinating populations with variable degrees of selfing. AlphaSimR also supports user
166 supplied mating plans.

167 *Genomic prediction*

168 Modeling genomic prediction in breeding programs is one of the main use cases
169 for AlphaSimR. AlphaSimR offers several built-in functions for fitting common
170 genomic prediction models. The built-in functions use mixed model solvers based on
171 the following R packages: rrBLUP (Endelman 2011), EMMREML (Akdemir and
172 Godfrey 2015) and Sommer (Covarrubias-Pazaran 2016). Each solver has been
173 optimized for performance within AlphaSimR and written in C++ using the R packages
174 Rcpp (Eddelbuettel and Francois 2011) and RcppArmadillo (Eddelbuettel and
175 Sanderson 2014). Users can also make use of other R packages or external applications
176 for modeling genomic prediction. This is done by exporting data from an AlphaSimR
177 simulation into another R function or external program for genomic predictions,
178 generating predictions, and importing the predictions back into AlphaSimR objects.

179 *Implementation*

180 Much of AlphaSimR's code has been written in C++ to improve performance.
181 For example, this has been used to implement bitwise storage of genotype data to
182 reduce memory usage and enable multithreading for increased speed. AlphaSimR also
183 improves performance by limiting data storage and calculations, such as variance
184 component calculations, to only those expressly requested by the user. This approach

185 differs from other stochastic simulation programs, including the original AlphaDrop
186 (Hickey and Gorjanc 2012) and AlphaSim (Faux *et al.* 2016), which typically perform
187 all calculations and save all data.

188

189 **Results and Discussion**

190 *Example Simulation*

191 This section will demonstrate AlphaSimR using a simulation of a single
192 breeding cycle for a generic wheat breeding program. The code needed to run this
193 simulation is presented below after a brief description of the breeding program.

194 Figure 1 shows a schematic representing the stages of the generic wheat
195 breeding program with a seven-year breeding cycle. In the first year, 200 bi-parental
196 populations are produced by crossing and production of doubled haploid (DH) lines
197 from those bi-parental populations begins. In the second year, the production of DH
198 lines is completed in. In the third year, the DH lines are visually evaluated in a head-
199 row (HDRW) nursery. In the remaining years, lines are selected based on performance
200 in the previous year and evaluated in a yield trial. The yield trials are conducted over
201 the course of three years before selecting a variety to release.

202 The first step is to generate founder haplotypes using MaCS. The founder
203 haplotypes will be used to form the initial parents in the breeding program. Code for
204 simulating the founder haplotypes for 50 inbred individuals is shown below. Each
205 individual will have 21 chromosomes, each with 1000 segregating sites.

```
206 founderPop = runMacS(nInd=50, nChr=21, segSites=1000, inbred=TRUE)
```

207 The second step is to set global parameters. Below is code for setting simulation
208 parameters to model a single trait. The trait models additive genetic effects on 1000 loci
209 per chromosome. The trait is also modeled as having a broad-sense heritability of 0.4
210 for evaluation in a single location.

```
211     SP = SimParam$  
212         new(founderPop)$  
213         addTraitA(1000)$  
214         setVarE(H2=0.4)  
215
```

216 The next step is to simulate each year of the breeding program. In the first year,
217 200 bi-parental populations are produced by crossing the parents formed from the
218 founder haplotypes. This code is presented below. The first line uses the founder
219 haplotypes to form the parents and the second line makes 200 randomly chosen crosses
220 between those parents.

```
221     Parents = newPop(founderPop)  
222     F1 = randCross(Parents, 200)  
223
```

224 In the second and third years, the DH lines are produced and then they are
225 evaluated in the HDRW nursery. The code for both these years is presented below. The
226 first line forms 100 DH lines per F₁ plant. The second line models evaluation in the
227 HDRW nursery for the previously defined additive trait. The broad-sense heritability
228 of this trait is reduced to 0.1 to represent visual selection.

```
229     HDRW = makeDH(F1, 100)  
230     HDRW = setPheno(HDRW, varE=9) #H2=0.1  
231
```

232 In the fourth year, the best entries in the HDRW nursery are selected and
233 evaluated in a preliminary yield trial (PYT). This is modeled with the code below. The
234 first line models selection in the HDRW by selecting the best lines within families. The
235 second line models evaluation of the PYT at one location. The accuracy of this
236 evaluation is based on the broad-sense heritability defined in the simulation parameters.

```
237     PYT = selectWithinFam(HDRW, 5)  
238     PYT = setPheno(PYT)  
239
```

240 In the fifth year, the best PYT entries are selected and evaluated in an advanced
241 yield trial (AYT). This is modeled with the code below. The first line models selection
242 of the best PYT lines. The second line models evaluation of the AYT at four locations,
243 which are represented as reps in the code.

```
244     AYT = selectInd(PYT, 100)  
245     AYT = setPheno(AYT, reps=4)  
246
```

247 In the sixth year, the best AYT entries are selected and evaluated in an elite
248 yield trial (EYT). This is modeled with the code below. The first line models selection
249 of the best AYT lines. The second line models evaluation of the EYT at sixteen
250 locations.

```
251     EYT = selectInd(AYT, 10)  
252     EYT = setPheno(EYT, reps=16)  
253
```

254 In the seventh year, the best performing EYT entry is chosen for release as a
255 variety. This is modeled with the code below.

```
256     Variety = selectInd(EYT, 1)  
257
```

258 The final step is to evaluate the simulation results. This is done by producing a
259 boxplot for the genetic values of entries in stage of the breeding program. The boxplot
260 is shown in Figure 2. The code for generating the boxplot is given below. The first line
261 of code extracts the genetic values for each entry and saves it in a list. The second line

262 creates the boxplot showing the distribution of genetic values for entries in each stage
263 of the breeding program.

```
264     yield = list(Parents=gv(Parents), F1=gv(F1),  
265                 HDRW=gv(HDRW), PYT=gv(PYT),  
266                 AYT=gv(AYT), EYT=gv(EYT),  
267                 Variety=gv(Variety))  
268     boxplot(yield, ylab="Genetic Value")  
269
```

270 **Concluding remarks**

271 AlphaSimR represents a considerable improvement over its predecessor in
272 terms of ease-of-use, flexibility, and computational efficiency (AlphaSim; Faux *et al.*
273 2016). It has been used in a handful of published simulations (Gorjanc *et al.* 2018;
274 Muleta *et al.* 2019; Johnsson *et al.* 2019) as well as numerous unpublished simulations.
275 The largest simulation undertaken in AlphaSimR to date involved over a hundred
276 million individuals (unpublished), a feat that would not be feasible with original
277 AlphaSim.

278 The improvements made to AlphaSimR make it uniquely well suited for
279 simulating whole breeding programs. These types of simulations serve as a valuable
280 tool for aiding strategic decision making within breeding programs. For example,
281 AlphaSimR can be used test the economic value of modifying an existing breeding
282 program. This will be of particular interest to breeding programs considering
283 implementing genomic selection or changing their current implementation. These types
284 of simulations can also be used to optimize selection stages or compare the efficiency
285 of mating strategies.

286 AlphaSimR can be used for a wide range of simulations outside of whole
287 breeding program simulations. For example, AlphaSimR can be used to test QTL

288 mapping strategies or marker imputation strategies. AlphaSimR is also well suited for
289 running simulations that help with teaching quantitative genetics and breeding. This is
290 because students can be quickly taught how to use AlphaSimR for simple simulations,
291 and the software's ability to report variance components, perform genomic evaluations
292 and evaluate accuracy of evaluations against the simulated true values is highly
293 instructive.

294 AlphaSimR is under continuous development with new features being added on
295 a semi-regular basis. Additional planned features include developing standard breeding
296 program blueprints for major species and developing easy-to-use graphical user
297 interfaces for these blueprints. These planned additions should make AlphaSimR even
298 more user-friendly than it currently is.

299

300 **Web resources**

301 AlphaSimR is publicly available on CRAN ([https://CRAN.R-](https://CRAN.R-project.org/package=AlphaSimR)
302 [project.org/package=AlphaSimR](https://CRAN.R-project.org/package=AlphaSimR)). Additional documentation as well as links to
303 graphical user interfaces for specialized applications are available on the AlphaGenes
304 website (<https://alphagenes.roslin.ed.ac.uk/wp/software-2/alphasimr/>). A repository of
305 example simulation scripts for learning to use the software, modeling specific breeding
306 programs, and learning quantitative genetics principles are available on Bitbucket
307 (https://bitbucket.org/hickeyjohnteam/alphasimr_examples).

308

309 **Acknowledgments**

310 The authors acknowledge the financial support from Bayer Crop Science and
311 the Gates Foundation.

312

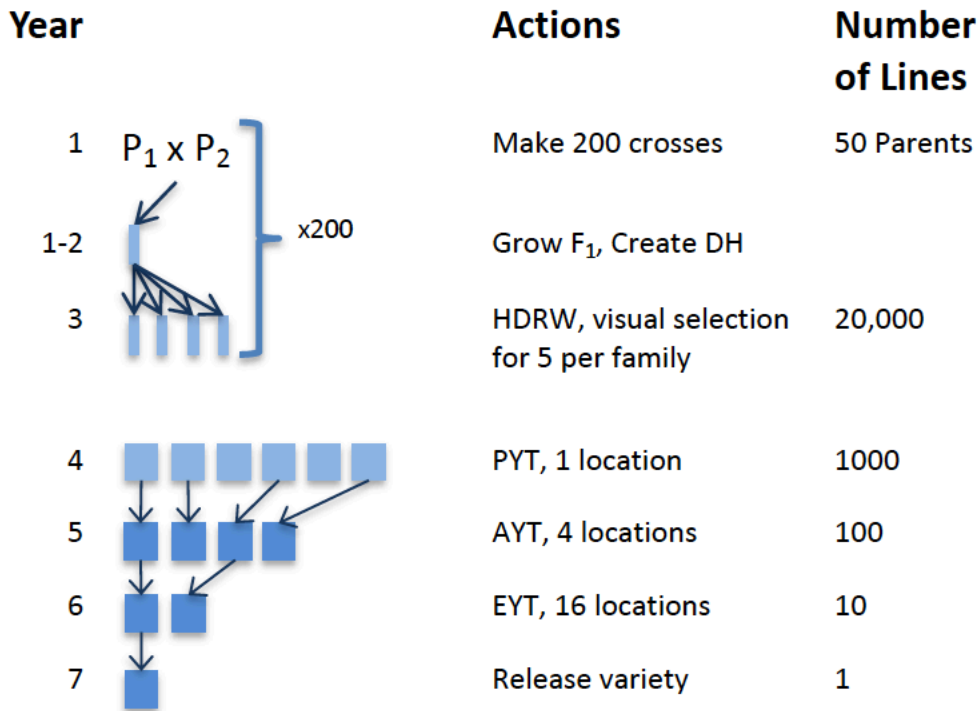
313 **References**

- 314 Akdemir, D., and O. U. Godfrey, 2015 *EMMREML: Fitting Mixed Models with*
315 *Known Covariance Structures*. [https://CRAN.R-](https://CRAN.R-project.org/package=EMMREML)
316 [project.org/package=EMMREML](https://CRAN.R-project.org/package=EMMREML)
- 317 Bourke, P. M., R. E. Voorrips, R. G. F. Visser, and C. Maliepaard, 2015 The Double-
318 Reduction Landscape in Tetraploid Potato as Revealed by a High-Density
319 Linkage Map. *Genetics* 201: 853–863.
- 320 Broman, K. W., and J. L. Weber, 2000 Characterization of Human Crossover
321 Interference. *Am. J. Hum. Genet.* 66: 1911–1926.
- 322 Bulmer, M. G., 1976 The effect of selection on genetic variability: a simulation study.
323 *Genet. Res.* 28: 101–117.
- 324 Chen, G. K., P. Marjoram, and J. D. Wall, 2009 Fast and flexible simulation of DNA
325 sequence data. *Genome Res.* 19: 136–142.
- 326 Covarrubias-Pazarán, G., 2016 Genome-Assisted Prediction of Quantitative Traits
327 Using the R Package sommer. *PLoS ONE* 11:.
- 328 Eddelbuettel, D., and R. Francois, 2011 Rcpp: Seamless R and C++ Integration. *J.*
329 *Stat. Softw. Artic.* 40: 1–18.
- 330 Eddelbuettel, D., and C. Sanderson, 2014 RcppArmadillo: Accelerating R with high-
331 performance C++ linear algebra. *Comput. Stat. Data Anal.* 71: 1054–1063.
- 332 Endelman, J. B., 2011 Ridge Regression and Other Kernels for Genomic Selection
333 with R Package rrBLUP. *Plant Genome J.* 4: 250.
- 334 Falconer, D. S., and T. F. C. Mackay, 1996 *Introduction to Quantitative Genetics*.
335 Longman, Harlow, UK.

- 336 Faux, A.-M., G. Gorjanc, R. C. Gaynor, M. Battagin, S. M. Edwards *et al.*, 2016
337 AlphaSim - software for the simulation of breeding programs. *The Plant*
338 *Genome Journal*.
- 339 Gaynor, R. C., G. Gorjanc, A. R. Bentley, E. S. Ober, P. Howell *et al.*, 2017 A Two-
340 Part Strategy for Using Genomic Selection to Develop Inbred Lines. *Crop Sci.*
341 *57*: 2372–2386.
- 342 Gaynor, R. C., G. Gorjanc, and J. M. Hickey, 2018 Dominance in stochastic
343 simulations of animal breeding programs, in *Proceedings of the 11th World*
344 *COngress on Genetics Applied to Livestock Production*, Auckland.
- 345 Gorjanc, G., P. Bijma, and J. M. Hickey, 2015 Reliability of pedigree-based and
346 genomic evaluations in selected populations. *Genet. Sel. Evol.* 47:.
- 347 Gorjanc, G., R. C. Gaynor, and J. M. Hickey, 2018 Optimal cross selection for long-
348 term genetic gain in two-part programs with rapid recurrent genomic selection.
349 *Theor. Appl. Genet.* 131: 1953–1966.
- 350 Hayes, B., and M. E. Goddard, 2003 Evaluation of marker assisted selection in pig
351 enterprises. *Livest. Prod. Sci.* 81: 197–211.
- 352 Hickey, J. M., and G. Gorjanc, 2012 Simulated Data for Genomic Selection and
353 Genome-Wide Association Studies Using a Combination of Coalescent and
354 Gene Drop Methods. *G3 GenesGenomesGenetics* 2: 425–427.
- 355 Jenko, J., G. Gorjanc, M. A. Cleveland, R. K. Varshney, C. B. A. Whitelaw *et al.*,
356 2015 Potential of promotion of alleles by genome editing to improve
357 quantitative traits in livestock breeding programs. *Genet. Sel. Evol.* 47: 1–14.
- 358 Johnsson, M., R. C. Gaynor, J. Jenko, G. Gorjanc, D.-J. de Koning *et al.*, 2019
359 Removal of alleles by genome editing (RAGE) against deleterious load.
360 *Genet. Sel. Evol.* 51: 14.

- 361 Lin, Z., N. O. I. Cogan, L. W. Pembleton, G. C. Spangenberg, J. W. Forster *et al.*,
362 2016 Genetic Gain and Inbreeding from Genomic Selection in a Simulated
363 Commercial Breeding Program for Perennial Ryegrass. *Plant Genome* 9:.
364 McPeck, M. S., and T. P. Speed, 1995 Modeling interference in genetic
365 recombination. *Genetics* 139: 1031–1044.
366 Muleta, K. T., G. Pressoir, and G. P. Morris, 2019 Optimizing Genomic Selection for
367 a Sorghum Breeding Program in Haiti: A Simulation Study. *G3 Genes*
368 *Genomes Genet.* 9: 391–401.
369 R Core Team, 2019 *R: A Language and Environment for Statistical Computing*. R
370 Foundation for Statistical Computing, Vienna, Austria.
371 Voorrips, R. E., and C. A. Maliepaard, 2012 The simulation of meiosis in diploid and
372 tetraploid organisms using various genetic models. *BMC Bioinformatics* 13:
373 248.
374
375

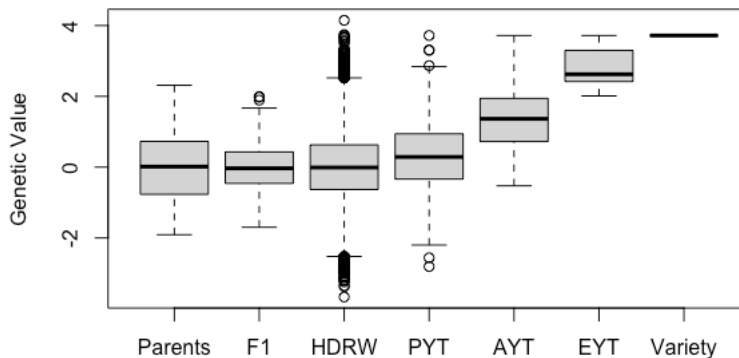
376 Figure 1. An overview of the variety development cycle for the example wheat
 377 breeding program. A variety is developed over the course of seven years. The
 378 steps in the development cycle are: making bi-parent crosses, forming doubled
 379 haploid (DH) lines, visually select lines grown in headrows (HDRW), evaluate
 380 lines in a preliminary yield trial (PYT), evaluate lines in an advanced yield trial
 381 (AYT), evaluate lines in an elite yield trial (EYT), and release a variety.



382

383

384 Figure 2. The distribution of genetic values in one replicate of the example
 385 breeding program. Separate boxplots are given for each stage of the breeding
 386 program.



387

388