

# MoBPS - Modular Breeding Program Simulator

T. Pook<sup>\*†,1</sup>, M. Schlather<sup>†‡</sup> and H. Simianer<sup>\*†</sup>

<sup>\*</sup>Department of Animal Sciences, Animal Breeding and Genetics Group, University of Goettingen, 37075 Goettingen, Germany, <sup>†</sup>Center for Integrated Breeding Research, University of Goettingen, 37075 Goettingen, Germany, <sup>‡</sup>Stochastics and Its Applications Group, University of Mannheim, 68131 Mannheim, Germany

## ABSTRACT

The R-package MoBPS provides a computationally efficient and flexible framework to simulate complex breeding programs and compare their economic and genetic impact. Simulations are performed on the base of individuals. MoBPS utilizes a highly efficient implementation with bit-wise data storage and matrix multiplications from the associated R-package *miraculix* allowing to handle large scale populations. Individual haplotypes are not stored but instead automatically derived based on points of recombination and mutations. The modular structure of MoBPS allows to combine rather coarse simulations, as needed to generate founder populations, with a very detailed modeling of today's complex breeding programs, making use of all available biotechnologies. MoBPS provides pre-implemented functions for common breeding practices such as optimum genetic contributions and single-step GBLUP but also allows the user to replace certain steps with personalized and/or self-written solutions.

## KEYWORDS

R-package  
simulation  
breeding  
breeding program  
resource management  
population genetics

## INTRODUCTION

Breeding programs aim at improving the genetic properties of livestock and crop populations with respect to productivity, fitness and adaptation. Progress towards the target is limited by the available resources, but also negative effects, such as inbreeding depression or health issues, have to be avoided or at least controlled. Hence, the allocation of resources in a breeding program is a complex optimization problem. Additionally, population history, such as fluctuating population sizes and selection pressures, has an impact on the current genomic architecture and thus the potential for future improvement.

Over the years a variety of simulation tools have been developed to assist breeders to evaluate and optimize their breeding programs. A general problem of simulation studies is that the underlying genomic processes are highly complex and have to be simplified for modeling. In addition, users often have rather different objectives in mind when setting up their simulation studies. Since tools often do not provide the necessary flexibility to execute the specific breeding actions and/or it is not possible to export all nec-

essary outputs, this commonly leads to the use of self-developed solutions that tend to be more error-prone, less sophisticated and computationally inefficient. The functionality of existing software for the simulation of breeding programs ranges from cohort based deterministic simulation that relies on expected gains like ZPLAN+ (Täubert *et al.* 2010) to applications on the base of the stochastic simulation of single individuals such as QMSim (Sargolzaei and Schenkel 2009) and AlphaSim (Faux *et al.* 2016). The functionality of each of these tools highly depends on the intended use. ZPLAN+ (Täubert *et al.* 2010) focuses on the economic impact from a macro-perspective. Since analytic formulas for cohorts are required, it has limitations when simulating complex mating schemes or when focusing on other quantities than genetic or economic gain. QMSim (Sargolzaei and Schenkel 2009) is able to simulate each individual meiosis but is lacking the flexibility in terms of design options in the breeding program, as it is mainly intended for use in population genetics. On the contrary, AlphaSim (Faux *et al.* 2016) provides a lot of flexibility in term of the design of the breeding program, especially for plant breeding and when the number of cohorts in the breeding program is small. However, AlphaSim lacks the efficiency to simulate complex and large scale populations (e.g. as genotypes of all individuals are stored) and does not allow for the export of all potentially relevant results. The interested reader is referred to Sun *et al.* (2011) for an extended review on different simulators used in plant breeding.

Manuscript compiled: Monday 30<sup>th</sup> March, 2020

<sup>1</sup>T. Pook; University of Goettingen, Department of Animal Sciences, Center for Integrated Breeding Research, Animal Breeding and Genetics Group, Albrecht-Thaer-Weg 3, 37075 Goettingen, Germany

Email: torsten.pook@uni-goettingen.de

46 Our goal was to develop a tool that combines the simulation of a  
47 historical population and the evaluation of a subsequent complex  
48 breeding program in a computationally efficient way. The Modular  
49 Breeding Program Simulator (MoBPS) is not only flexible in terms  
50 of parameters and design of breeding programs, but also allows  
51 the user to replace standard procedures of the package with own  
52 ones.

## 53 METHODS

54 Simulations in MoBPS are ultimately based on the simulation of  
55 single individuals. In principle, this allows the user to control each  
56 singular mating and modify recombination or mutation rates for  
57 the respective meiosis. However, breeding programs in MoBPS  
58 can still be constructed in a modular form as a combination of co-  
59 horts and breeding actions. As breeding actions like phenotyping,  
60 selection, aging, or reproduction are typically applied on groups  
61 of individuals, the relevant individuals for each breeding action  
62 can be selected via three different keywords:

- 63 1. gen: all individuals of a certain generation
- 64 2. cohort: group of individuals generated via the same breeding  
65 action
- 66 3. database: all (or in principle specific) individuals of a certain  
67 generation and sex

68 Similar to the gene-flow concept (Hill 1974), a cohort describes  
69 a group of individuals with usually identical characteristics like  
70 age, sex and genetic origin. As the three types of groups can also  
71 be combined, handling of overlapping generations in a breeding  
72 program is possible. Cohorts and breeding actions are defined in a  
73 generic way and are parametrized, so that any breeding program  
74 of arbitrary complexity can be modeled as a suitable sequence of  
75 cohorts and breeding actions.

76 All data for a population is stored in a list that contains general  
77 and individual information. The general part provides information  
78 on the underlying genetics like the physical position of each  
79 marker, allelic variants or structure of the underlying genetic traits.  
80 The individual part contains information that is specific to the  
81 individual. Similarly to simulators like SBVB (Pérez-Enciso et al.  
82 2017), haplotypes are stored for founder individuals only. For all  
83 other individuals only points of recombination and mutation, and  
84 their genetic origins are stored. In particular, haplotypes are not  
85 permanently stored but only derived when needed. Therefore, the  
86 required memory is minimized and only increases slightly with  
87 increasing marker density. When thousands of generations are  
88 simulated it is advisable to classify additional generations as new  
89 founders to reduce the number of recombinations and mutations  
90 to be stored in subsequent generations. The usefulness of this is  
91 highly dependent on the ratio between genome length and number  
92 of markers and on how many new founder genotypes have to be  
93 stored. As the population itself in a breeding program usually  
94 occupies only a small share of the required memory and less than  
95 one hundred generations are considered, benefits here are usually  
96 small, making this only relevant/required in large scale population  
97 genetic studies.

98 Simulation of multiple correlated traits with and without underlying  
99 QTL is supported. Classical additive, dominant and epistatic  
100 or pleiotropic QTL can be defined and any effect structure of mul-  
101 tiple interacting loci is supported. Each locus has to be assigned  
102 with a position in Morgan and different recombination rates for  
103 subgroups (e.g. males/females) are supported. Information on  
104 the number of markers can be manually entered or imported via a

database (Ensembl, (Zerbino et al. 2017)), a map-file (Purcell et al.  
2007) or a vcf-file (Danecek et al. 2011). For common species, exem-  
plary map files are provided in the associated package MoBPSmaps  
(Pook 2019). Genotype data for a base population can be imported  
via PLINK (Purcell et al. 2007) and/or vcf-format (Danecek et al.  
2011), sampled internally or generated by executing prior simu-  
lation in MoBPS and/or other tools (Chen et al. 2009; Sargolzaei  
and Schenkel 2009) to generate the required population structure.  
All breeding actions performed in the simulation can be tracked  
and assigned with costs to derive the expenses of the program.  
Different breeding programs can be compared in terms of their  
economic revenue or other target functions (e.g. development of  
the inbreeding rate) one is interested in.

Common methods for selection such as optimal genetic contri-  
butions (Meuwissen 1997) are implemented and a variety of dif-  
ferent packages for breeding value estimation can be switched  
on. This includes BGLR (Pérez and de los Campos 2014), sommer  
(Covarrubias-Pazaran 2016) and rrBLUP (Endelman 2011), as well  
as an efficient implementation for solving the mixed model (Hen-  
derson 1975) in the traditional GBLUP model (Meuwissen et al.  
2001; VanRaden 2008) that is assuming known heritability and  
is using the R-package RandomFieldsUtils (Schlather et al. 2019)  
for the matrix inversion. Inputs for these packages such as the  
different pedigree and genomic relationship matrices (VanRaden  
2008; Legarra et al. 2014; Martini et al. 2017) can be derived via  
highly efficient and fully-parallelized bit-wise matrix multiplica-  
tions (R-package miraculix (Schlather 2020)). Non of the men-  
tioned packages, however, is required to execute simulations in  
MoBPS. In particular, all functionality of the MoBPS R-package is  
still available when miraculix is not installed, with the downside  
of higher computing times and memory demands.

The simulations in MoBPS are based on two main functions: *creat-  
ing.diploid()* and *breeding.diploid()*. Here, *creating.diploid()* initializes  
the base-line population and *breeding.diploid()* performs breeding  
actions on an existing population list. As a simple example con-  
sider the following script:

```
141 library(MoBPS)
142 pop <- creating.diploid(nsnp=10000, nindi=100,
143   chr.nr=5, chromosome.length=2,
144   n.additive=50, n.dominant=10,
145   var.target=1, name.cohort="Founder")
146 pop <- breeding.diploid(pop, heritability=0.5,
147   new.bv.observation="all")
148 pop <- breeding.diploid(pop, bve=TRUE)
149 pop1 <- breeding.diploid(pop,
150   breeding.size=100,
151   selection.size=c(20,20),
152   selection.criteria="bve",
153   selection.m.cohorts="Founder_M",
154   selection.f.cohorts="Founder_F",
155   name.cohort="Offspring")
156 pop2 <- breeding.diploid(pop,
157   breeding.size=100,
158   selection.size=c(5,20),
159   selection.criteria="bve",
160   selection.m.cohorts="Founder_M",
161   selection.f.cohorts="Founder_F",
162   name.cohort="Offspring")
```

163 Via this code, we first generate a base population containing 100  
164 individuals with 10,000 markers. The underlying genome consists  
165 of 5 chromosomes with a length of 2 Morgan each and equidistant  
166 markers. Furthermore, we generated a single trait that is impacted

167 by 50 purely additive QTLs and 10 dominant QTLs and scale QTL 227  
168 effect to result in a trait with genomic variance of 1. 228

169 In the next step, we initialize a breeding action to generate phe- 229  
170 notypes for all individuals in the population with an assumed 230  
171 heritability of 0.5. Next, a breeding value estimation is performed. 231  
172 Since no cohorts are selected, the last (and only) generation of the 232  
173 population list will be considered for the breeding value estimation. 233  
174 Lastly, we generate 100 offspring by random mating. Here, 234  
175 two scenarios are considered with different selection intensity with 235  
176 the top 5/20 male and top 20 female individuals being used for 236  
177 reproduction, leading to gains of 0.663/0.922 genomic standard 237  
178 deviations and an increase in inbreeding in terms of average kin- 238  
179 ship of 0.0062/0.0156 (Figure 1). In principle, all three breeding 239  
180 actions performed via *breeding.diploid()* could have also been exe- 240  
181 cuted in a joint step. For a full list of all possible breeding actions 241  
182 and available parameters we refer to our user manual (available at 242  
183 <https://github.com/tpook92/MoBPS>). 243

184 For a quick overview of the simulated population, the function 244  
185 *summary()* can be used: 245

```
186 > summary(pop1) 246  
187 Population size: 247  
188 Total: 200 Individuals 248  
189 Of which 100 are male and 100 are female. 249  
190 There are 2 generations 250  
191 and 4 unique cohorts. 251  
192  
193 Genome Info: 252  
194 There are 5 unique chromosomes. 253  
195 In total there are 10000 SNPs. 254  
196 The genome has a total length of 10 Morgan. 255  
197 The genome has a physical size of about: 1 GB 256  
198  
199 Trait Info: 257  
200 There is 1 modelled trait. 258  
201 The trait has underlying QTL 259  
202 The trait is named: Trait 1
```

203 A variety of functions is provided to export required informa-  
204 tion such as the phenotypes (*get.pheno()*), the genotypes (*get.geno()*)  
205 and the pedigree (*get.pedigree()*) for selected individuals from the  
206 population list. These functions are thoroughly described in chap-  
207 ter 7 of the user manual (available at [https://github.com/tpook92/](https://github.com/tpook92/MoBPS)  
208 [MoBPS](https://github.com/tpook92/MoBPS)). In addition, functions to derive rates of inbreeding (*kin-*  
209 *ship.emp()*), development of breeding values (*bv.development()*) or  
210 changes in allele frequency over time (*analyze.population()*) are pro-  
211 vided to further analyze the resulting population list.

## 212 RESULTS AND DISCUSSION

213 The package MoBPS is completely written in R (R Core Team 2017)  
214 so that all functionalities for genetic applications are platform in-  
215 dependent. The R-packages miraculix (Schlather 2020) can be  
216 activated in MoBPS and leads to more efficient data storage and  
217 shorter simulation times. In particular vector multiplications with  
218 genetics data (0,1,2) are performed via bitwise operations on a  
219 whole register (128/256 bit) using SSE2/AVX2. Computing times  
220 are similar to the ones in PLINK (Purcell et al. 2007) with one fourth  
221 of the memory usage. The interested reader is referred to Schlather  
222 (2020) for extended benchmark-testing on miraculix. The storage  
223 of founder genotypes is 4 times as efficient as in AlphaSim (Faux  
224 et al. 2016) and 32/64 times as efficiency as the use of integer/dou-  
225 ble variables to store haplotypes.  
226 Even though basically all information regarding each individual is

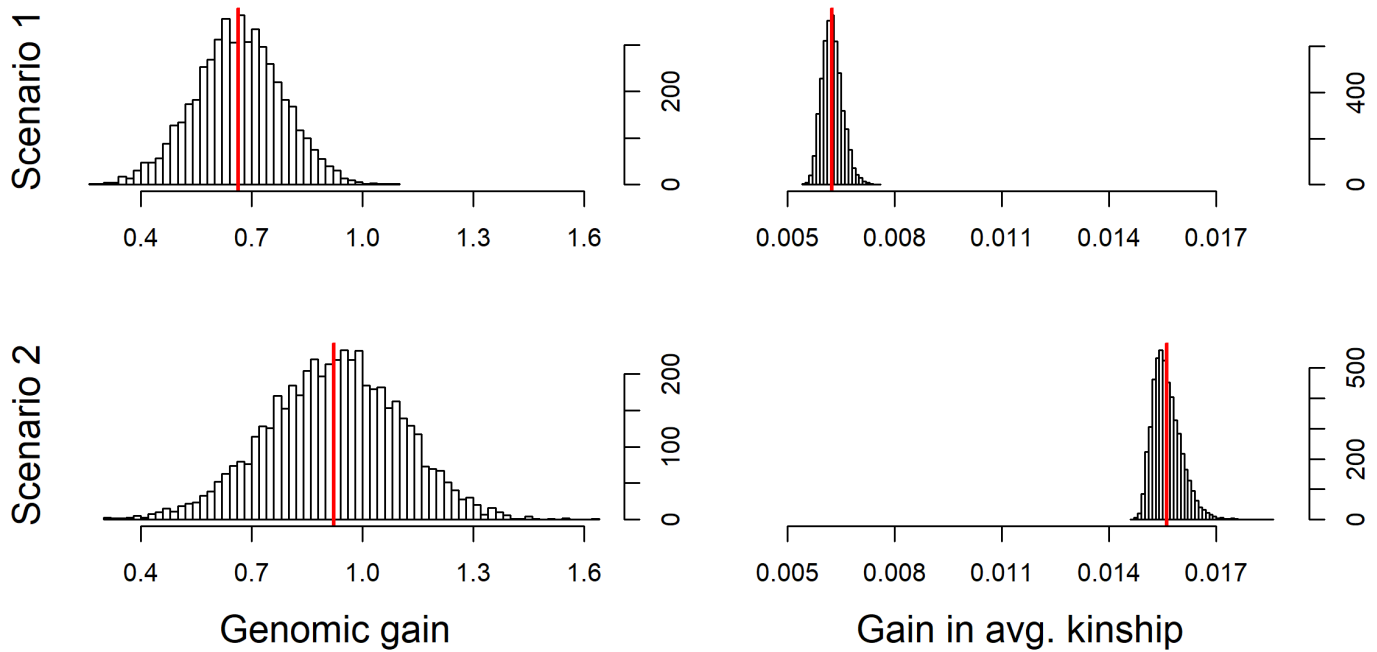
stored, the required memory in MoBPS is still relatively low as a  
highly efficient storage structure is used. Haplotypes of founders  
and details on the origin of the segment between points of recom-  
bination are stored bitwise. E.g. the simulation of 20 generations  
with 50,000 cows with 50,000 markers and breeding value estima-  
tion via GBLUP takes 26.2 hours using 24 cores on a server cluster  
with Intel E5-2650 (2X12 core 2.2GHz) processors. At peak, 65 GB  
of memory was used. The main share of this was required for the  
storage of the genomic relationship matrix whereas the resulting  
population list, containing more than a million individuals, only  
had a size of about 0.44 GB. The biggest proportion of the comput-  
ing time is used for breeding value estimation (25.3 hours, 96.4%).  
The generation of new animals took 55 minutes (3.5%, 304 animals  
per second using a single core). All other parts needed negligible  
computing time (132 seconds, 0.1%). Computing times for most  
parts (except breeding value estimation) increase linearly with the  
number of individuals. This highly efficient storage structure there-  
fore also allows for the simulation of historical populations with  
thousands of generations and undergone population dynamics  
such as genetic bottlenecks, migration or mutational drift.  
The flexible and efficient environment of MoBPS allows for the  
simulation of a variety of different and potential large-scale breed-  
ing programs. For exemplary scripts of more complex breeding  
programs we refer to the user manual. Exemplary simulations are  
given for the effect of gene editing in a cattle breeding program  
(Simianer et al. 2018), the simulation of a multi-parent advanced  
generation intercross in maize (Pook et al. 2019), an introgression  
scheme in chicken (Ha et al. 2017) and the generation of a base  
population with a hard sweep. A further advantage of MoBPS  
compared to other simulation tools is its flexible structure that  
allows the user to substitute single steps of the breeding program  
with a customized approach. For this consider the following exam-  
ple to execute one owns breeding value estimation:

```
260 genos <- get.geno(pop, gen=1)  
261 y <- get.pheno(pop, gen=1)  
262 indi_names <- colnames(genos)  
263 # Execute one owns function to perform  
264 # the breeding value estimation  
265 y_hat <- own.method.for.bve(genos, y)  
266  
267 # Enter BVEs in the population-list  
268 pop <- insert.bve(pop,  
269 bves = cbind(indi_names, y_hat))
```

Even though a simulation study can never fully reflect reality  
and is relying on model assumptions, the use of a simulation study  
comes with major benefits and still allows the user to draw im-  
portant conclusions. In contrast to reality the underlying truth in  
a simulation study is known, and therefore new methods can be  
thoroughly evaluated and compared to existing ones. Furthermore,  
the effects of particular breeding actions on a variety of output  
dimension can be assessed and compared. This in turn can be  
used to derive an ideal resource allocation and optimize poten-  
tially highly complex breeding scenarios in a setting that can be  
evaluated multiple times and without constrains both in terms of  
money and time.

## 282 WEB RESOURCES

283 An executable version of MoBPS and the associated R-packages  
284 miraculix (Schlather 2020), RandomFieldsUtils (Schlather et al.  
285 2019) and MoBPSmaps (Pook 2019) for Windows and Linux are  
286 freely available at <https://github.com/tpook92/MoBPS>. This directory



**Figure 1** Resulting genetic gain and increase in inbreeding with low (Scenario 1) and high (Scenario 2) selection pressure for 5,000 simulation runs each.

287 also contains an comprehensive user manual explaining the func- 318  
 288 tionality of all input parameters and utility functions in MoBPS. A 319  
 289 frozen version of the R-packages MoBPS (v1.4.87), MoBPSmaps 320  
 290 (v0.1.7), miraculix (v0.9.10), RandomFieldsUtils (v0.5.9), and our 321  
 291 user manual at submission are also provided there. The MoBPS 322  
 292 R-package can be directly installed within your R session via fol- 323  
 293 lowing commands:

```
294 install.packages("devtools")
295 devtools::install_github("tpook92/MoBPS",
296   subdir="pkg")
```

## 297 ACKNOWLEDGMENTS

298 This package was developed in the context of the European 329  
 299 Union's Horizon 2020 Research and Innovation Program under 330  
 300 grant agreement n°677353 IMAGE.

## 301 LITERATURE CITED

302 Chen, G. K., P. Marjoram, and J. D. Wall, 2009 Fast and flexible 337  
 303 simulation of dna sequence data. *Genome Research* **19**: 136–142. 338  
 304 Covarrubias-Pazaran, G., 2016 Genome-assisted prediction of 339  
 305 quantitative traits using the r package sommer. *PLOS ONE* **11**: 340  
 306 e0156744. 341  
 307 Danecek, P., A. Auton, G. Abecasis, C. A. Albers, E. Banks, *et al.*, 342  
 308 2011 The variant call format and vcftools. *Bioinformatics* **27**: 343  
 309 2156–2158. 344  
 310 Endelman, J. B., 2011 Ridge regression and other kernels for ge- 345  
 311 nomic selection with r package rrblup. *The Plant Genome* **4**: 346  
 312 250–255. 347  
 313 Faux, A.-M., G. Gorjanc, R. C. Gaynor, M. Battagin, S. M. Edwards, 348  
 314 *et al.*, 2016 Alphasim: Software for breeding program simulation. 349  
 315 *The Plant Genome* **9**: doi:10.3835/plantgenome2016.02.0013. 350  
 316 Ha, N.-T., T. Pook, C. Dierks, S. Weigend, R. Preisinger, *et al.*, 2017 351  
 317 A simulation approach to optimize breeding programs with 352

application to the introgression of the blue egg color into a high 318  
 performing layer line. 10th European Symposium on Poultry 319  
 Genetics **2017**: 107.

Henderson, C. R., 1975 Best linear unbiased estimation and predic- 320  
 tion under a selection model. *Biometrics* pp. 423–447. 321

Hill, W. G., 1974 Prediction and evaluation of response to selection 322  
 with overlapping generations. *Animal Science* **18**: 117–139. 323

Legarra, A., O. F. Christensen, I. Aguilar, and I. Misztal, 2014 Single 324  
 step, a general approach for genomic selection. *Livestock Science* 325  
**166**: 54–65. 326

Martini, J. W. R., N. Gao, D. F. Cardoso, V. Wimmer, M. Erbe, 327  
*et al.*, 2017 Genomic prediction with epistasis models: On the 328  
 marker-coding-dependent performance of the extended gblup 329  
 and properties of the categorical epistasis model (ce). *BMC Bioin-* 330  
 formatics **18**: 3. 331

Meuwissen, T. H. E., 1997 Maximizing the response of selection 332  
 with a predefined rate of inbreeding. *Journal of animal science* 333  
**75**: 934–940. 334

Meuwissen, T. H. E., B. J. Hayes, and M. E. Goddard, 2001 Predic- 335  
 tion of total genetic value using genome-wide dense marker 336  
 maps. *Genetics* **157**: 1819–1829. 337

Pérez, P. and G. de los Campos, 2014 Genome-wide regression & 338  
 prediction with the bgrr statistical package. *Genetics* pp. 483– 339  
 495. 340

Pérez-Enciso, M., N. Forneris, G. de Los Campos, and A. Legarra, 341  
 2017 Evaluating sequence-based genomic prediction with an 342  
 efficient new simulator. *Genetics* **205**: 939–953. 343

Pook, T., 2019 Mobps: Modular breeding program simulator: Avail- 344  
 able at <https://github.com/tpook92/mobps>; r-package version 345  
 1.4.14. 346

Pook, T., M. Schlather, G. de los Campos, M. Mayer, C. C. Schoen, 347  
*et al.*, 2019 Haploblocker: Creation of subgroup specific haplo- 348  
 type blocks and libraries. *Genetics* pp. 1045–1061. 349

Purcell, S., B. Neale, K. Todd-Brown, L. Thomas, M. A. R. Ferreira, 350  
*et al.*, 2007 Plink: A tool set for whole-genome association and 351  
 352



353 population-based linkage analyses. *The American Journal of*  
354 *Human Genetics* **81**: 559–575.

355 R Core Team, 2017 R: A language and environment for statistical  
356 computing.

357 Sargolzaei, M. and F. S. Schenkel, 2009 Qmsim: A large-scale  
358 genome simulator for livestock. *Bioinformatics* **25**: 680–681.

359 Schlather, M., 2020 Efficient calculation of the genomic relationship  
360 matrix. *bioRxiv* p. 2020.01.12.903146.

361 Schlather, M., R. Furrer, and M. Kroll, 2019 Randomfieldsutils:  
362 Utilites for the simulation and analysis of random fields: Avail-  
363 able at <https://github.com/tpook92/mobps>; r-package version  
364 0.5.9.

365 Simianer, H., T. Pook, and M. Schlather, 2018 Turning the page -  
366 the potential of genome editing in breeding for complex traits  
367 revisited. *World Congress on Genetics Applied to Livestock* p.  
368 190.

369 Sun, X., T. Peng, and R. H. Mumm, 2011 The role and basics of  
370 computer simulation in support of critical decisions in plant  
371 breeding. *Molecular breeding* **28**: 421–436.

372 Täubert, H., F. Reinhardt, and H. Simianer, 2010 Zplan+, a new  
373 software to evaluate and optimize animal breeding programs.  
374 *World Congress on Genetics Applied to Livestock* p. 950.

375 VanRaden, P. M., 2008 Efficient methods to compute genomic pre-  
376 dictions. *Journal of Dairy Science* **91**: 4414–4423.

377 Zerbino, D. R., P. Achuthan, W. Akanni, M. R. Amode, D. Barrell,  
378 *et al.*, 2017 Ensembl 2018. *Nucleic acids research* **46**: D754–D761.