

---

# GENO-DIVER (V2)

## Genetic Simulation Toolkit

---

Last Updated on November 18, 2017

Jeremy T. Howard<sup>1</sup>  
Francesco Tiezzi<sup>2</sup>  
Jennie E. Pryce<sup>3,4</sup>  
Christian Maltecca<sup>2</sup>

<sup>2</sup>University of Nebraska, Lincoln, NE, USA

<sup>2</sup>North Carolina State University, Raleigh, NC, USA

<sup>3</sup>Department of Economic Development, Jobs, Transport and  
Resources, Victoria, Australia

<sup>4</sup>La Trobe University, Victoria, Australia.



## Contents

<b>Introduction</b>	<b>1</b>
<b>Disclaimer</b>	<b>3</b>
<b>Download</b>	<b>4</b>
<b>Computing Environment</b>	<b>5</b>
<b>Overview of Program</b>	<b>6</b>
<b>Running the Program</b>	<b>7</b>
<b>Program Parameters</b>	<b>8</b>
General Parameters . . . . .	8
Genome and Marker Information . . . . .	10
QTL Distributions . . . . .	14
Population Parameters . . . . .	17
Selection and Culling Parameters . . . . .	18
Mating Parameters . . . . .	25
Output Options . . . . .	26
<b>Appendix I - Output Information</b>	<b>29</b>
<b>Appendix II - Generation of Effects</b>	<b>35</b>
<b>Appendix III - MaCS Sequence Simulation</b>	<b>38</b>
<b>Literature Cited</b>	<b>41</b>
<b>Example 1: Running the Simulation Software</b>	<b>44</b>
<b>Example 2: Quantitative Trait - Multiple Progeny</b>	<b>47</b>
<b>Example 3: Differential Sire Contribution by Age</b>	<b>49</b>
<b>Example 4 Fitness Trait</b>	<b>52</b>
<b>Example 5 (Quantitative Trait + Fitness Correlated)</b>	<b>54</b>

## Introduction

A variety of disciplines including conservation (McMahon et al. 2014), animal and plant breeding (De los Campos et al. 2013) and human genetics (Yang et al. 2010) are currently making use of large volumes of genetic marker information. In particular, within animal and plant breeding programs the use of genomic information to predict the genetic merit of individuals has become a routine practice (Jonas & Koning 2015; Morrell et al. 2012). This has resulted in a significant increase in the number of individuals within a herd/population having genomic information. Nonetheless, the ability to use this information to efficiently manage agricultural populations at the genomic level, both for preserving genetic diversity and lessening inbreeding depression, remains a challenge for the near future. Although the identification of lethal mutations of large effect segregating within livestock populations with genomic data is possible (VanRaden et al. 2011), optimal mating procedures that minimize the frequency of a large number of lethal and more importantly sublethal mutations across generations have not been fully implemented. As the popularity of genotyping breeding individuals increases across species, the possibility to utilize genomic information from multiple sources to manage genomic information will also increase. Methods that make effective use of information from multiple sources including performance, genome diversity and inbreeding load at the selection and/or mating step are in increasing need. The routine genotyping of individuals is a significant investment from several sectors of agriculture and the need to spread the costs across multiple avenues are of considerable practical interest.

The use of simulation is a low cost alternative to assess and validate proposed methods to predict genetic values or to compare alternative se-

lection or mating strategies across time. A number of simulation programs have been developed that simulate breeding livestock/crops populations, but they primarily focus on testing strategies where the genetic architecture is based solely on a quantitative trait governed by additive effects (Sargolzaei & Schenkel 2009; Hickey et al. 2012; Cheng et al. 2015; Prez-Enciso & Legarra 2016). Currently there is no single self contained software that can simulate complex traits involving both quantitative and fitness components along with the ability to generate complex pedigrees and genomic information ranging from sparse marker to sequence information. Consequently, determining how various selection and management practices impact the fitness and the overall genomic variability of a population undergoing selection for a quantitative trait remains challenging. Furthermore, the precise genetic architecture of complex traits is largely unknown although in general it is assumed that complex traits are affected by variation in a large number of genes, most of which have individually minor effects (Weiss 2008). Knowledge on the demographic history of a population and the extent of linkage disequilibrium is also being generated from dense marker panels or sequence information (McKay et al. 2007; Ai et al. 2013; Porto-Neto et al. 2014). As more information is generated on the genetic architecture of complex traits and variation across the genome, simulation can be employed to more effectively predict how current selection and management practices will impact future generations.

## Disclaimer

This document outlines how to run the Geno-Diver simulation program and describes the parameters utilized within the program. The software is free for any person to use. The authors accept no responsibility for the accuracy of results obtained by using Geno-Diver software.

Geno-Diver is being updated with new tools and functions routinely. If you would like the program to do something that it currently does not do, don't hesitate to contact me at [jthoward@ncsu.edu](mailto:jthoward@ncsu.edu).

Please notify [jthoward@ncsu.edu](mailto:jthoward@ncsu.edu) or [cmaltec@ncsu.edu](mailto:cmaltec@ncsu.edu) if you think the results are not correct or you have encountered a bug. We have written the program in a way for us to reconstruct the simulation scenario you ran to more effectively solve the problem.

The software was profiled using the Valgrind software ([valgrind](#)). In a few instances, calling functions within the Intel MKL library did result in "possibly lost" errors. Although, as outlined in the Intel MKL [user manual](#), new buffers that the library allocates when an Intel MKL function gets called are not deallocated until the program ends. The use of `mkl_free_buffers()` and `mkl_thread_free_buffers()` were called at the end of the program and this fixed the majority of the issues.

The reference for the Geno-Diver software:

- Howard, J.T., F. Tiezzi, J.E. Pryce, C. Maltecca. A combined coalescence forward in time simulator software for pedigreed populations undergoing selection for complex traits. *J. Anim. Breed. Genet.*

The reference for the MaCS software:

- Chen, G. K., P. Marjoram & J. D. Wall. 2009. Fast and flexible simulation of DNA sequence data. *Genome Res.* 19:136-142.

## Download

Source code and executables files are available at the [GIT-HUB](#) account.

### Linux Executable Files:

- macs
- msformatter
- GenoDiver

### Source Code Files:

- ParameterClass.cpp
- Animal.h
- AnimalFun.cpp
- EBV\_Functions.cpp
- Genome\_ROH.cpp
- Genome\_ROH.h
- HaplofinderClasses.cpp
- HaplofinderClasses.h
- MatingDesignClasses.cpp
- MatingDesignClasses.h
- ParameterClass.cpp
- ParameterClass.h
- SelectionCullingFunctions.cpp
- Simulation\_Functions.cpp
- makefile

## Computing Environment

The code is written in the C++11 language using object-oriented techniques. The application has been tested to run on Linux and MAC platforms. The software makes use of two external libraries, Intel MKL and Eigen.

### EIGEN Library:

EIGEN is freely available at: [Eigen Site](#)

Once at the site, download the latest stable release and uncompress it. For example, the current downloaded package is called “eigen-eigen-07105f7124f9.tar”. To use it you just have to place it in the file where all of the other Geno-Diver files are located and uncompress the file. Once you uncompressed the file it will be a folder (e.g. “eigen-eigen-07105f7124f9”). The uncompressed file serves as your path in the makefile outlined below.

### Intel MKL Library:

Intel MKL is a commercial library and is available for purchase. However, there is an opportunity to obtain the Intel MKL library (for Linux) free of charge for non-commercial use at the following website: [Intel MKL Site](#)

The Intel MKL can sometimes be tricky to download and link, but there is a step-by-step protocol within the folders that are downloaded or instructions are at [Intel MKL Guide](#). Depending on the computing system you are running, a guide to linking the Intel MKL libraries are at [Intel MKL Linking Guide](#).

### C++11 Version:

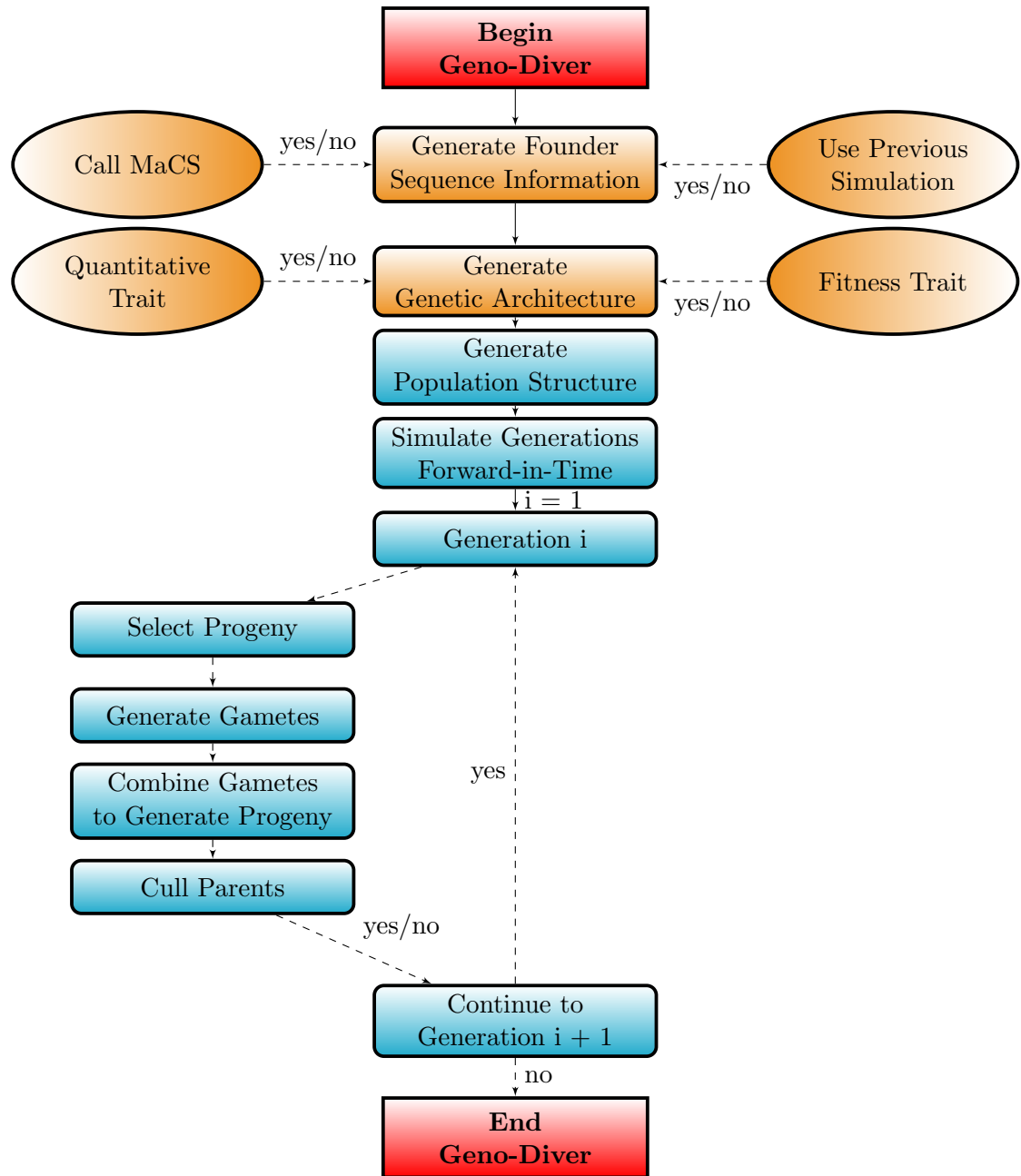
Versions of gcc 4.7 or newer support the C++11 standards. You need to install or update to the correct version of gcc using the standard package manager or installer, depending on what type of OS you are using. Some helpful websites include:

[gcc helper 1](#) [gcc helper 2](#)

### Compiling:

Once both EIGEN, Intel MKL libraries and gcc version 4.7 or newer have been correctly installed and the folders placed in the directory where all of the Geno-Diver source code files are located the last thing you have to do is change the path for EIGEN and Intel MKL libraries. A makefile has been provided for both a Linux and a MAC operating system. In the makefile change lines 13 and 14 to the path which aligns to the updated EIGEN and MKL path. After changing the path, rename the appropriate file to “makefile” and then type “make” on the command line. An executable file called “GenoDiver” is now in your working directory.

## Overview of Program





## Running the Program

At the current time executable files are available only for Linux operating environments and are located on the GitHub page ([GENO-DIVER](#)). The [website](#) provides multiple examples on how to simulate different scenario's along with snippets of R code to illustrate the results. To run the program place the following executable files in a folder:

- GenoDiver.
- macs.
- msformatter.

Before running the program, the file permissions need to be checked. After verifying the permissions, a parameter file is generated and placed in the same folder as the previous three files. If you are new to the simulation program, multiple examples are on the GitHub page. The simulation program reads the parameter file by searching for keywords that are capitalized and then followed by a colon. Therefore any phrase that does not meet the search criteria is ignored when initializing parameters within the program. Also, if you want to comment out a parameter just add "!!" within the key word and the program will skip over it. For example to skip over the "SEED" parameter just replace it with "SE!!ED" and it won't recognize the parameter any more.

To run the program type `./GenoDiver` and then the name of your parameter file. For example, if the parameter file is named "parameterfile", the program can be ran by typing `./GenoDiver parameterfile`. During the simulation, the program outputs minimal comments on the progress. A more thorough description of the status of the program is printed to the log file (i.e. "log\_file.txt").

**After running the program it is a strongly recommended to check the parameters initialized at the top of the log file within the output folder. The log file contains a large amount of information and is a great tool to ensure that the parameters and outcome of the simulation is what is intended.** If the program is not running correctly the log file should provide knowledge on why and where the simulation crashed or exited.

If you are unsure of the problem an e-mail can be sent to [jthoward@ncsu.edu](mailto:jthoward@ncsu.edu) with a copy of the log file for me to replicate the results.

## Program Parameters

The complete list and description of the parameter files are outlined below and the appendix also contains further discussion and suggestions. Not all of them are required for the program to run. All keywords should be in capital letters and the parameter(s) specified are separated by spaces. Only parameters after the keywords impact the simulation. Lastly, the order in which parameters appear does not matter.

### General Parameters

---

#### **START**

Description: - Determines where to start the simulation program.

Options:

- [sequence](#): Starts at the sequence generation step.
- [founder](#): Skips sequence generation and begins generating the founder population utilizing sequence information from a previous run.

Usage: - "START: sequence".

Type: - Mandatory.

Note:

The use of this option saves time and space due to the large size of sequence information from MaCS. If you are using a different effective population size or MaCS diversity metric you always have to start at the sequence step, while any replicate that has the same MaCS parameters can start at the founder step and use the previously generated haplotypes saved within the output folder. When calling multiple replicates, the simulation for any replicate after the first starts at the founder step by default.

---

#### **SEED**

Description: - Declares the seed number.

Option: - Can be any integer value.

Usage: - "SEED: 1501".

Type: - Optional. Default is the system time.

Note:

- The seed appears at the beginning of the log file if the user does not declare it.
-

## **OUTPUTFOLDER**

Description: - Directory name of files generated from simulation.

Option: - Any valid folder name, but can't have spaces.

Usage: - "OUTPUTFOLDER: GenoDiverFiles".

Type: - Optional. Default is "GenoDiverFiles".

Note:

If the directory contains information from a previous simulation, it may be deleted depending on the value of the START parameter. If the START parameter is "sequence" everything is deleted. If the START parameter is "founder" everything but sequence information is deleted. The sequence files could potentially take up a large amount of memory. Only one sequence file can be generated at a time unless the mac, msformatter and GenoDriver executables are in separate folders.

---

## **THREAD**

Description: - Declares the number of threads used for parallel processing.

Option: - Integer value based on number of cores available.

Usage: - "NTHREAD: 4".

Type: - Optional. Default is 1.

---

## **REPLICATES**

Description: - Declares the number of replicates to produce and the seed is increased by 1 from the previous value for each one.

Option: - Integer value.

Usage: - "NREP: 1".

Type: - Optional. Default is 1.

Note:

- All important replicate files are in a folder within the output folder path called "replicates". Each file has the seed number at the end to distinguish between replicates.

## Genome and Marker Information

---

### **CHR**

Description: - Sets the number of chromosomes to generate.

Option: - Integer from 1 to 31.

Usage: - "CHR: 3".

Type: - Mandatory.

---

### **CHR\_LENGTH**

Description: - Sets the length of each chromosome, in Megabases.

Option: - Integer Value.

Usage: - "CHR\_LENGTH: 150 150 150"

Type: - Mandatory.

Note:

- The number of Chromosome length variables specified has to correspond to chromosome number. If not will exit the program.

---

### **NUM\_MARK**

Description:

- Sets the number of evenly spaced markers for each chromosome.

Option: - Integer Value.

Usage: - "NUM\_MARK: 4000 4000 4000"

Type: - Mandatory.

Note:

- Quantitative and fitness QTL cannot be markers. The numbers of markers per chromosome specified has to correspond to chromosome number. The number of markers has to be less than the available number of mutations within a chromosome generated from MaCS.

---

### **MARKER\_MAF**

Description: - Minimum allele frequency allowed for markers.

Option: - Range from 0.0 to 0.50.

Usage: - "MARKER\_MAF: 0.10"

Type: - Optional. Default is 0.10.

---

## **QTL**

### Description:

- The number of quantitative trait loci (QTL) for each chromosome. The QTL location is generated based on a uniform distribution from 0 to the length of the chromosome.

### Options:

- Integer. Number of QTL and fitness trait loci (FTL) can't exceed 5000.

### Usage:

- "QTL: 50 50 50"

### Type:

- Mandatory.

### Note:

- The number of QTL variables specified has to correspond to chromosome number. If not will exit the program.
- 

## **QUANTITATIVE\_MAF**

Description: - Minimum allele frequency allowed for QTL.

### Value:

- 0.0 to less than 0.5.

### Usage:

- "QUANTITATIVE\_MAF: 0.05"

### Type:

- Optional. Default is 0.05.
- 

## **FIT\_LETHAL**

### Description:

- Lethal FTL can not have a covariance with the quantitative trait. The number of lethal FTL variables specified has to correspond to chromosome number.

### Option:

- Integer. Number of QTL & FTL can't exceed 5000.

### Usage:

- "FIT\_LETHAL: 25 25 25".

### Type:

- Optional. Default is 0.

### Note:

- Lethal FTL can't have covariance with the quantitative trait.
-

## **FIT\_SUBLETHAL**

### Description:

- Sub-lethal FTL can have a covariance with the quantitative trait. The number of sub-lethal FTL variables specified has to correspond to chromosome number.

### Option:

- Integer. Number of QTL & FTL can't exceed 5000.

### Usage:

- "FIT\_SUBLETHAL: 25 25 25".

### Type:

- Optional. Default is 0.

### Note:

- Sub-lethal FTL can have covariance with the quantitative trait.
- 

## **FITNESS\_MAF**

### Description:

- Parameters that describe the maximum FTL frequency allowed and the range in lethal FTL allele frequencies. If the parameter is given, the first two are required, while the last one is optional.

### Options:

- **double (0.0-0.50)**: Maximum allele frequency for the unfavorable allele allowed for lethal FTL.
- **double (0.0-0.50)**: Maximum allele frequency for the unfavorable allele allowed for sub-lethal FTL.
- **double (0.0-0.50)**: Maximum allele frequency range for the lethal FTL.

### Usage:

- "FITNESS\_MAF: 0.02 0.08 0.01".

### Type:

- Optional. Default is 0.02, 0.08 and 0.01.

### Note:

- Using the default values, the lethal FTL frequencies range from 0.01 to 0.02 and the sub-lethal frequencies range from 0.0 to 0.08. If the frequency of either lethal or sub-lethal FTL is too high, a large number of founders and progeny may not make it to breeding age. If the number of individuals is too small to generate the founder generation the simulation exits.
-

## FOUNDER\_HAPLOTYPES

Description:

- The number of haplotypes generated by MaCS. Need to ensure that it is greater than two times the total number of female and male animals needed in the founder population.

Option:

- Integer Value.

Usage:

- "FOUNDER\_HAPLOTYPES: 4000".

Type:

- Optional. Default is based on male & female number.
- 

## HAPLOTYPE\_SIZE

Description:

- The number of markers contained within a non-overlapping haplotype window. Haplotype windows are used to generate haplotype-based relationship matrices and summary statistics.

Option:

- Integer Value.

Usage:

- "HAPLOTYPE\_SIZE: 50".

Type:

- Optional. Default is 50.
- 

## RECOMBINATION

Description:

- The type of distribution that generates the location of recombination events along the genome. The number of recombination events is generated from a Poisson distribution. The rate parameter fixed at 1.0 across all chromosomes.

Options:

- **Uniform**: Recombination sampled from a Uniform distribution from 0 to 1.0.
- **Beta**: Recombination sampled from a Beta distribution (0.5, 0.5). Recombinations occurs more often at the end of the chromosome than towards the middle.

Usage:

- "RECOMBINATION: Uniform".

Type:

- Optional. Default is Uniform.
-

## MUTATION

### Description:

Used in the MaCS software to generate scaled mutation parameter and in the simulation to generate new mutations as generations proceed.

### Options:

- **Mutation Rate (float)**: Probability of a new mutation occurring at a given base pair and follows the infinite-site model. A Poisson distribution with a rate parameter equal to the mutation rate times the length of the chromosome in nucleotides generates the total number of mutations occurring within a new gamete.
- **The proportion of Mutations that can be QTL (double)**: Within the forward-in-time part of the simulation program, represents the number of non-neutral mutations that occurred out of the total number of new mutations. Each type of trait (i.e. quantitative, lethal, sub-lethal) has an equal chance of being chosen.

### Usage:

- "MUTATION: 2.5e-8 0.0".

### Type:

- Optional. Default is 2.5e-8 0.0.

## QTL Distributions

---

## ADD\_QUAN

### Description:

- The parameters for the gamma distribution that generate the un-scaled additive effect for QTL. The effects are scaled such that the sum of the QTL variances in the founder population equals the proportion of the variance that is due to the additive gene action specified. A complete description is in Appendix II.

### Options:

- **Shape**: Shape of gamma distribution.
- **Scale**: Scale of gamma distribution.

### Usage:

- "ADD\_QUAN: 0.4 1.66".

### Type:

- Optional. Default is 0.4 1.66.
-



## **DOM\_QUAN**

### Description:

- The parameters for the normal distribution that generate the degree of dominance (h) for QTL. The effects are scaled such that the sum of the QTL variances in the founder population equals the proportion of the variance that is due to the dominant gene action specified. A complete description is in Appendix II.

### Options:

- [Mean](#): Mean of a normal distribution.
- [Standard Deviation \(SD\)](#): SD of a normal distribution.

### Usage:

- “DOM\_QUAN: 0.1 0.2”.

### Type:

- Optional. Default is 0.1 0.2.
- 

## **LTHA**

### Description:

- The parameters for the gamma distribution that generate the selection coefficients for lethal FTL. A complete description is in Appendix II.

### Options:

- [Shape](#): Shape of gamma distribution.
- [Scale](#): Scale of gamma distribution.

### Usage:

- “LTHA: 70.0 0.01”.

### Type:

- Optional. Default is 70.0 0.01.
- 

## **LTHD**

### Description:

- The parameters for the normal distribution that generate the degree of dominance for the lethal FTL. A complete description is in Appendix II.

### Options:

- [Mean](#): Mean of a normal distribution.
- [Standard Deviation \(SD\)](#): SD of a normal distribution

### Usage:

- “LTHD: 0.001 0.0001”.

### Type:

- Optional. Default is 0.001 0.0001.
-

## SUBA

### Description:

- The parameters for the gamma distribution that generate the selection coefficients for sublethal FTL. A complete description is in Appendix II.

### Options:

- [Shape](#): Shape of gamma distribution.
- [Scale](#): Scale of gamma distribution.

### Usage:

- "SUBA: 0.1 0.2".

### Type:

- Optional. Default is 0.1 0.2.
- 

## SUBD

### Description:

- The parameters for the normal distribution that generate the degree of dominance for sublethal FTL. A complete description is in Appendix II.

### Options:

- [Mean](#): Mean of a normal distribution.
- [Standard Deviation \(SD\)](#): SD of a normal distribution

### Usage:

- "SUBD: 0.3 0.1".

### Type:

- Optional. Default is 0.3 0.1.
- 

## COVAR

### Description:

- Determines the relationship between the additive effect of the quantitative trait and selection coefficient for the sub-lethal trait. The relationship is a function of the number of QTL that impact the quantitative and sub-lethal trait and the rank correlation between the effects. A complete description is in Appendix II.

### Options:

- [Proportion of Pleiotropic QTL](#): Proportion of QTL that have both a quantitative and sub-lethal effect.
- [Genetic correlation](#): The rank correlation between QTL effects.

### Usage:

- "COVAR: 0.5 0.2".

### Type:

- Optional. Default is 0.0 0.0.

### Note:

- Care needs to be taken to generate the desired relationship between the quantitative and fitness trait. Fitness values range from 0 to 1 and higher values leading to a lower fitness value.

If the two traits are antagonistic under the scenario of high values being favorable, the correlation should be positive. One just needs to change the favorable direction of the quantitative trait to alter the interpretation.

## Population Parameters

---

### **FOUNDER\_Effective\_Size**

Description:

- Used to generate the population history of the haplotypes generated from MaCS. There are multiple default scenarios and represent a wide range of LD patterns, or one can specify a custom effective population size and population history. Each scenario has a slightly different population history parameter (i.e. “eN”). Lastly, the population history can be a single value with no historical population. A description is in Appendix III.

Options:

- [Ne70](#): A scenario that generates a large amount of short LD.
- [Ne100\\_Scen1](#): A scenario that generates moderate short LD.
- [Ne100\\_Scen2](#): A scenario that generates moderate short LD.
- [Ne250](#): A scenario that generates minimal LD.
- [Ne1000](#): A scenario that generates minimal LD.
- [CustomNe](#): Read in own population history parameters.
- [Any value greater than 1](#): Utilizes the value as the effective population size and no population history assumed.

Usage:

- “FOUNDER\_Effective\_Size: 50”.

Type:

- Mandatory.

---

### **MALE\_FEMALE\_FOUNDER**

Description:

- Sets the number of founder male and female individuals along with the method founders were selected and how many generations it occurred.

Options:

- [Integer value](#): Number of male founders.
- [Integer value](#): Number of female founders.
- [random](#): At the current only random selection can be utilized.
- [Integer value](#): Number of generations conducted.

Usage:

- “MALE\_FEMALE\_FOUNDER: 50 600 random 5”.

Type:

- Mandatory.

---

## **VARIANCE\_A**

### Description:

- Proportion of variance due to additive gene action. Phenotypic variance set at 1.0.

### Option:

- Double value ranging from 0.0 to 1.0.

### Usage:

- "VARIANCE\_A: 0.25".

### Type:

- Mandatory.

---

## **VARIANCE\_D**

### Description:

- Proportion of variance due to dominant gene action. Phenotypic variance set at 1.0.

### Option:

- Double value ranging from 0.0 to 1.0.

### Usage:

- "VARIANCE\_D: 0.05".

### Type:

- Mandatory.

### Note:

- Care needs to be taken in determining the additive and dominance variance and its implications on the number of QTL that display over-dominance or partial-dominance. Parameters to change that impact the dominance variance include the QTL MAF frequency, mean and standard deviation of the normal distribution that generates the degree of dominance parameter and the ratio of additive to dominance variance.

## **Selection and Culling Parameters**

---

## **GENERATIONS**

Description: - Determines the number of generations to simulate.

### Option:

- Any integer value.

### Usage:

- "GENERATIONS: 10".

### Type:

- Mandatory.

### Note:

- The number of generations allowed vary depending on the size of your computer memory. It is advisable to start at a small number of generations and then build up to determine how much memory a scenario will use.

---

## INDIVIDUALS

### Description:

- Determines the number of males and females in each generation and replacement rate for parents each generation. Care should be taken on picking the number of offspring to have enough for the next generation. If the number of animals falls below the input male or female value, the program exits.

### Options:

- **Integer value:** Number of male parents within a generation.
- **Double value (0.0 to 1.0):** Proportion of male parents that are culled and replaced by progeny each generation.
- **Integer value:** Number of female parents within a generation.
- **Double value (0.0 to 1.0):** Proportion of female parents that are culled and replaced by progeny each generation.

### Usage:

- "INDIVIDUALS: 50 0.2 600 0.2".

### Type:

- Mandatory.
- 

## PARITY\_MATES\_DIST

### Description:

- Determines the distribution of the number of mating pairs a sire has for each age group. A Beta distribution, which is parameterized by two parameters, is used to generate the distribution of mating pairs for a given age group. To generate the number of mating pairs by age class, the cumulative distribution function (CDF) is split into quadrants based on the number of age classes that occur within a generation. The total number of mating pairs within an age class is the proportion that falls within the CDF quadrant for a given age class.

### Options:

- Double value ( $> 0.0$ ): Alpha parameter in Beta distribution.
- Double value ( $> 0.0$ ): Beta parameter in Beta distribution.

### Usage:

- "PARITY\_MATES\_DIST: 1 1".

### Type:

- Optional. The default is both parameters being 1, which is very similar to a uniform distribution, such that all age classes have the same proportion of mating pairs.
-

## **PROGENY**

### Description:

- Determines the number of progeny produced by each mating pair. The number of progeny per mating pair may not be the actual number of progeny produced if FTL segregate in the population. The impact of a progeny fitness value on its ability to make it to breeding age is described in detail in Appendix II.

### Option:

- Integer value.

### Usage:

- "PROGENY: 4".

### Type:

- Mandatory.
- 

## **MAXFULLSIB**

### Description:

- Determines the maximum number of full-sib progeny selected within a family. Once the maximum number is reached the full-sib with the lowest selection criteria is no longer selected and the next best animal is selected. This process is repeated across all families until all families are below the value.

### Option:

- Ranges from 1 to number of progeny.

### Usage:

- "MAXFULLSIB: 2".

### Type:

- Optional. Default set at number of progeny.
- 

## **SELECTION**

### Description:

- The metric used to select individuals and favorable direction. If selection is based on optimal contribution selection also need the relationship constrained.

### Options:

- Selection Criteria:
  - **random**: parent selected randomly.
  - **phenotype**: parent selected based on phenotypic value.
  - **true\_bv**: parent selected based on true breeding value.
  - **ebv**: parents selected based on estimated breeding value.
- Direction:
  - **high**: High values are favorable.
  - **low**: Low values are favorable.

### Usage:

- "SELECTION: ebv high".

### Type:

- Mandatory.

### Note:

- For random selection, the direction does not impact the results and therefore the direction value does not matter.
- 

## EBV\_METHOD:

### Description:

- Parameter that specifies how estimated breeding values are generated.

### Options:

- Method:
  - **pblup**: The best linear unbiased prediction (BLUP) of estimated breeding values are obtained by Henderson's (Henderson 1975) mixed model equations. The mixed model equations are outlined below:

$$\begin{bmatrix} X'X & X'Z \\ Z'X & Z'Z + A^{-1}(\frac{\sigma_e^2}{\sigma_a^2}) \end{bmatrix} \begin{bmatrix} \hat{b} \\ \hat{a} \end{bmatrix} = \begin{bmatrix} X'y \\ Z'y \end{bmatrix},$$

where X and Z refers to fixed and random design matrices, which relate records to fixed and random effects,  $A^{-1}$  refers to the inverse of the pedigree based relationship matrix. Fixed effects include the intercept and random effects include the effect of the animal. The mixed model equations are solved by the preconditioned conjugate gradient method or cholesky decomposition of the LHS of the matrix.

- **gblup**: The same as the pblup option, breeding values are obtained by solving the mixed model equations, except the inverse of the pedigree based relationship matrix is replaced by the inverse of the genomic based relationship matrix.

- **rohblup**: The same as the pblup option, breeding values are obtained by solving the mixed model equations, except the inverse of the pedigree based relationship matrix is replaced by the inverse of the run-of-homozygosity based relationship matrix.

- **bayes**: Marker effects are estimated utilizing bayesian whole genome regression models based on the following model:

$$y_i = \mu + \sum_{j=1}^{j=m} SNP_j + \epsilon_i,$$

where y is the phenotype for individual<sub>i</sub>,  $\mu$  is the intercept, SNP is the additive genetic effects that correspond

to allele substitution effects for each marker and  $\epsilon$  is the residual for individual  $i$ . SNP covariates had values of 0 for the homozygote, 1 for the heterozygote and 2 for the alternative homozygote. The intercept is assigned an uninformative prior and the SNP covariates can have 4 possible prior densities: 1.) Gaussian (BayesRidgeRegression); 2.) Scaled-t (BayesA); 3.) Two finite mixture priors: a mixture of a point of mass at zero and a Gaussian slab (BayesC); 4.) Two finite mixture priors: a mixture of a point of mass at zero and a Scaled-t slab (BayesB). Samples were drawn from the posterior density using a Gibbs sampler with scalar updating.

- Usage: - "EBV\_METHOD: pblup".
- Type: - Only mandatory if selection is based on ebv.
- Note:
- If left out and selection is not based on ebv the program will not calculate breeding values. If parameter included and selection is not based on ebv, the program will calculate breeding values.

## BLUP\_OPTIONS:

### Description:

- When ebv are estimated based on 'pblup' or gblup' this parameter can be included to remove older individuals from the analysis, specify solving method and how the inverse of the relationship matrix is calculated.

### Options:

- **Generations (integer)**: Number of generations to trace back from the current group of selection candidates. All animals 'n' generations back and their associated progeny will be included in the analysis. For example, if a value of 1 is specified the selection candidate's parents and the associated parents progeny will only be included in the analysis. The default value is the number of generations simulated (i.e. all animals are used each generation).
- Solver:
  - **direct**: Uses Cholesky decomposition for matrix inversion. If EBV accuracy is needed or only simulating a small number of generations this is advised.
  - **pcg**: Uses the iterative preconditioned conjugate gradient (PCG) method and is faster than direct when the number of animals is large. The EBV accuracy is not calculated (**Default**).



- Genomic Inverse (only applicable if 'gblup' is specified):
  - [cholesky](#): update previous inverse based on the algorithm presented by Meyer et al. (2012) (**Default**).
  - [recursion](#): utilizing the sequential update algorithm presented by Misztal et al. (2014).

Usage: - "BLUP\_OPTIONS: 4 pcg cholesky".

Type: - Optional.

Note:

Parameters that specify inverse calculations are for genomic-based relationships. Calculation of pedigree-based relationships is generated based on Meuwissen & Luo (1992).

---

## G\_OPTIONS:

Description:

- When ebv are estimated based on 'gblup' this parameter can be included to determine how the genomic relationship matrix is calculated.

Options:

- How genomic relationship is constructed:
  - [VanRaden](#): Constructed based on  $G=ZZ' / 2\Sigma(pq)$ .
- Allele frequencies calculated:
  - [founder](#): Frequencies calculated based on founder genome.
  - [current population](#): Frequencies calculated based on animals who are selection candidates..

Usage: - "G\_OPTIONS: VanRaden observed".

Type: - Optional.

---

## BAYESOPTIONS:

Description:

- When ebv are estimated based on the "bayes" option, this parameter has to be included to determine which method to utilize and MCMC options.

Options:

- How genomic relationship is constructed:
  - [BayesRidgeRegression](#): Gaussian prior.
  - [BayesA](#): Scaled-t prior.
  - [BayesB](#): Spike with a point of mass at 0 and slab with the slab having a scaled-t prior.
  - [BayesC](#): Spike with a point of mass at 0 and slab with the slab having a gaussian prior.

- **Integer value**: Number of MCMC iterations the first time a method is implemented (Required).
- **Integer value**: Number of MCMC iterations to burn-in the first time a method is implemented (Required).
- **Integer value**: Number of MCMC iterations for any generation after the first time a method is implemented (Required).
- **Integer value**: Number of MCMC iterations to burn-in for any generation after the first time a method is implemented (Required).
- **Integer value**: Thinning rate for MCMC iterations (Required).
- How pie is determined: (Only needed for BayesB & BayesC)
  - **fix**: Fix proportion of SNP with non-null effect.
  - **estimate**: Estimate proportion of SNP with non-null effect.
- **Double value**: Proportion of markers with non-null effect. (Only needed for BayesB & BayesC)
- **Integer value**: Number of generations to trace back from the current group of selection candidates. All animals 'n' generations back and their associated progeny will be included in the analysis. For example, if a value of 1 is specified the selection candidate's parents and the associated parents progeny will only be included in the analysis. The default value is the number of generations simulated (i.e. all animals are used each generation). (Optional).

Usage:

Type:

- "BAYESOPTIONS: BayesC 10000 2000 6000 1000 5 fix 0.10".
- If 'bayes' EBV option is utilized has to be included otherwise can't be in parameter file.

Note:

The number of iterations to run can be investigated by looking at the 'bayes\_mcmc\_samples' files within the output directory. For each generation the files gets replaced by the most recent results.

## CULLING

### Description:

- The metric used to cull individuals and maximum age an animal can remain in the population before being removed due to old age. The direction is the same as the selection direction.

### Options:

- Culling Criteria:
  - [random](#): parent selected randomly.
  - [phenotype](#): parent selected based on phenotypic value.
  - [true.bv](#): parent selected based on true breeding value.
  - [ebv](#): parents selected based on estimated breeding value.
- Age Removed:
  - [Integer](#): Age at which an animal has to be removed from the population.

### Usage:

- "CULLING: ebv 5".

### Type:

- Mandatory.

## Mating Parameters

---

## MATING

Description: - Parameters that decide how animals are mated and the algorithm utilized to optimize matings.

### Options:

- Mating Option:
  - [random](#): Males and female parents are randomly mated.
  - [random5](#): Relationships  $\geq 0.5$  are not allowed to mate otherwise same as random mating option.
  - [random25](#): Relationships  $\geq 0.25$  are not allowed to mate otherwise same as random5 mating option.
  - [random125](#): Relationships  $\geq 0.125$  are not allowed to mate otherwise same as random25 mating option.
  - [minPedigree](#): Minimize parental co-ancestries based on a pedigree relationship matrix.
  - [minGenomic](#): Minimize parental co-ancestries based on a genomic relationship matrix.
  - [minROH](#): Minimize parental co-ancestries based on a ROH relationship matrix.
  - [pos assort](#): Positive assortative mating based on the value animals are being selected on.
  - [neg assort](#): Negative assortative mating based on the value

animals are being selected on.

**-index**: Mating is done based on an index of multiple values and is parameterized by 'MATING\_INDEX' option.

- Optimization Method:

**-simu\_anneal**: Adaptive simulated annealing algorithm to assign mates (Ingber 1993).

**-linear\_prog**: Uses the Hungarian algorithm to find an exact solution to the best possible mating design.

**-sslr**: Sequential selection of least related mates algorithm to assign mates (Pryce et al. 2011).

**-gp**: Genetic programming to assign mating pairs (Chu & Beasley, 1997).

Usage: - "MATING: minPedigree sslr".

Type: - Mandatory. Only the first option is required for random, pos\_assort, and neg\_assort.

Note: For avoidance mating's (i.e. random5, random25, random 125) any coancestry below the threshold is zeroed out and mating are optimized. If utilizing the index mating parameter, the only options are can be utilized are 'gp' or 'sslr'. The 'gp' option can only be used for index mating.

## Output Options

---

### OUTPUT\_LD

Description: - Used to determine if you want to calculate the linkage disequilibrium decay based on the r2 metric for each generation across the genome and surrounding QTL.

Option: - yes or no.

Usage: -"OUTPUT\_LD: no".

Type: - Optional. Default is no.

---

## GENOTYPES

### Description:

- Used to determine if saving genotypes and at what generation to begin saving them. If you do not need all the genotypes, this saves space and reduces running time.

### Options:

- No or yes and if yes provide the generation at which to start exporting genotypes. Generation 0 is founder population.

### Usage:

- "GENOTYPES: yes 0".

### Type:

- Optional. Default is yes and starting at generation 0.
- 

## GENOME\_ROH

### Description:

- Calculates the proportion of the genome in a ROH for each individual. Also the frequency of a SNP being in a ROH and the length of ROH the SNP is contained in is calculated across the genome for a specified generation number.

### Options:

- ROH cutoff
  - Mb (integer): ROH cutoff in Megabases.
- Generation number(s) ROH frequency and length are calculated:
  - Generation (integer): The generations that the ROH statistics across the genome are calculated. Each generation is separated by a space.

### Usage:

- "GENOME\_ROH: 5 5 10 15".

### Type:

- Optional.
- 

## WINDOWQTLVAR

### Description:

- Calculates the additive ( $\sum 2pq[a + d(q - p)]^2$ ) and dominance ( $\sum (2pqd)^2$ ) variance within 1Mb non-overlapping windows for each generation.

### Option:

- yes or no.

### Usage:

- "WINDOWQTLVAR: no".

### Type:

- Optional. Default is no.
-

## **TRAINREFER\_STATS**

### Description:

- Calculates the correlation between estimated breeding value (EBV) and TGV, TBV, and TDD for the training (i.e. anyone other than selection candidates) and validation population (i.e. selection candidates). Also, the bias between EBV and TBV for the training and validation population is calculated. The bias is the regression coefficient of TBV on EBV. Lastly, the maximum pedigree relationship between an individual in the training population and the selection candidates is calculated across generations.

### Option:

- yes or no.

### Usage:

- "TRAINREFER\_STATS: no".

### Type:

- Optional. Default is no.

## Appendix I - Output Information

The folder that contains the information generated by the simulation program contains multiple files, and a description of each one is below.

### Data Summary Files

#### Summary\_Statistics\_DataFrame\_Performance:

Generation: Generation number.

phen: Mean (variance) phenotypic value.

ebv: Mean (variance) estimated breeding value.

gv: Mean (variance) true genotypic breeding value ( $\Sigma (a + d)$ ).

bv: Mean (variance) true genotypic breeding value ( $\Sigma a$ ).

dd: Mean (variance) true dominance deviation ( $\Sigma d$ ).

res: Mean (variance) residual value.

#### Summary\_Statistics\_DataFrame\_Inbreeding:

Generation: Generation number.

ped\_f: Mean pedigree based inbreeding parameter.

gen\_f: Mean genomic relationship diagonal constructed based on Van Raden (2008).

h1\_f: Mean diagonal of haplotype based relationship matrix (Hickey et al. 2012; H1).

h2\_f: Mean diagonal of haplotype based relationship matrix (Hickey et al. 2012; H2).

h3\_f: Mean diagonal of ROH based relationship matrix (Howard et al. Submitted).

homozy: Mean proportion homozygous (i.e.  $1 - \text{homozy} = \text{Observed Heterozygosity}$ ).

PropROH: Mean proportion of the genome in ROH of a given length.

ExpHet: Expected Heterozygosity (i.e.  $\Sigma (1 - p^2 - q^2)$ ).

fitness: Mean multiplicative fitness value of an individual.

homozlethal: Mean number of homozygous FTL classified as lethal.

hetezlethal: Mean number of heterozygous FTL classified as lethal.

homozysublethal: Mean number of homozygous FTL classified as sub-lethal.

hetezsublethal: Mean number of heterozygous FTL classified as sub-lethal.

lethalequiv: Mean lethal equivalents (Lethal equivalents =  $\Sigma s$  for an animal).

#### Summary\_Statistics\_QTL

Generation: Generation number.

Quant\_Founder\_Start: Number of QTL from founder generation segregating.

Quant\_Founder\_Lost: Number of QTL from founder generation fixed.

Mutation\_Quan\_Total: Number of QTL from new mutations segregating.

Mutation\_Quan\_Lost: Number of QTL from new mutations fixed.

Additive\_Var: True additive genetic variance based on  $\Sigma 2pq[a+d(q-p)]^2$ .

Dominance\_Var: True dominance genetic variance based on  $\Sigma (2pqd)^2$ .

Fit\_Founder\_Start: Number of FTL from founder generation segregating.

Fit\_Founder\_Lost: Number of FTL derived from founder generation fixed.

Mutation\_Fit\_Total: Number of FTL derived from new mutations segregating.

Mutation\_Fit\_Lost: Number of FTL derived from new mutations fixed.

Avg\_Haplotypes\_Window: Mean haplotypes contained within a haplotype window.

ProgenyDiedFitness: Number of progeny that died due to fitness.

**LD\_Decay:**

A file that has the average correlation ( $r^2$ ) between two SNP across a range of distances. The average was generated by moving across the genome in 10 Mb blocks and randomly grabbing two SNP and calculating their respective ( $r^2$ ) and placing them in the correct bins based how far they were apart. Within a block 500 pairs of SNP are randomly sampled and once finished the window is shifted by 5 Mb and is repeated until the end of the chromosome. This is conducted within each chromosome. The distances are in the first row and are in Kilobases. Each row after the first row corresponds to the generation, such that line 2 is generation 0, line 3 is generation 1, etc. The formula to calculate the D and ( $r^2$ ) values is below and the subscript refers to either SNP marker 1 or 2.

$$D = (A_1 B_1 \times A_2 B_2) - (A_1 B_2 \times A_2 B_1)$$

$$r^2 = \frac{D^2}{p_1(1 - p_1)p_2(1 - p_2)}$$

**QTL\_LD\_Decay:**

Similar to the "LD\_Decay" file this file estimates the average correlation between a SNP and a QTL within window sizes of 0 - 0.5 Mb, 0.5 - 1.0 Mb, 1.0 - 1.5 Mb, 1.5 - 2.0 Mb and 2.0 - 2.5 Mb. For example if a SNP and QTL were 0.75 Mb apart it would get placed in the 0.5 - 1.0 Mb bin. Each line in the file represents a QTL. If the correlation could not be estimated due to lack of SNP in the window or if SNP are near fixation a '-5' will be produced. The first column represents the chromosome and the second column represents the position in Mb. The last column is the average correlation within a generation across the different window sizes. Correlations within a generation across windows are separated by a ":" and sets of correlations across generations are separated by a "\_".

**Phase\_Persistence:**

The phase was calculated as the square root of ( $r^2$ ) between a SNP and a QTL and was the same sign as D (de Roos et al. 2008). Each row represents a QTL with a particular SNP. The first and second column is the QTL and marker location. The third column determines whether the value is used in the correlation in phase across generation. Once a phase can't be estimated it can no longer be in the correlation calculation. The remaining columns is the phase estimate for each generation.

**Phase\_Persistence\_Generation:**

Using the output from the Phase\_Persistence file, the correlation between phases across generations was estimated within window sizes of 0 - 0.5 Mb, 0.5 - 1.0 Mb, 1.0 - 1.5 Mb, 1.5 - 2.0 Mb and 2.0 - 2.5 Mb. The first column is the generation and the columns following would be the correlation between the current generation and preceding generations.



**ProgenyParentCorrelationGeneration:**

Generation: Generation number.

Cor\_Parent\_TGV: Correlation between parents ebv and TGV ( $\Sigma (a + d)$ ).

Cor\_Parent\_TBV: Correlation between parents ebv and TBV ( $\Sigma (a)$ ).

Cor\_Parent\_TDD: Correlation between parents ebv and TDD ( $\Sigma (d)$ ).

Bias\_Parent\_TBV: Regression of TBV on parents ebv and the value is the regression coefficient.

Cor\_Progeny\_TGV: Correlation between selection candidates ebv and TGV.

Cor\_Progeny\_TBV: Correlation between selection candidates ebv and TBV.

Cor\_Progeny\_TDD: Correlation between selection candidates ebv and TDD.

Bias\_Progeny\_TBV: Regression of TBV on selection candidates ebv and the value is the regression coefficient.

**AmaxGeneration:**

This file estimates the mean maximum relationship for individuals that were born in a previous generation with selection candidates.

**Summary\_Statistics\_ROH\_Freq:**

The first two columns are the chromosomal and nucleotide position of the SNP and the remaining columns are the frequency of that SNP being in an ROH of the length that was specified for a given generation. A SNP may not be in a window of a given length and therefore is set to -5.

**Summary\_Statistics\_ROH\_Length:**

The first two columns are the chromosomal and nucleotide position of the SNP and the remaining columns are the mean and median length (i.e. "Mean\_Median") of ROH for that given SNP when that SNP is in a ROH for a given generation. A SNP may not be in a window of a given length and therefore is set to -5.

**WindowAdditiveVariance:**

Provides the true additive genetic variance for a given 1-Mb window across generations. The first row is the chromosome and mid-point in the window and the remaining lines are the associated additive genetic variance estimates for a given window.

**WindowDominanceVariance:**

Provides the true dominance genetic variance for a given 1-Mb window across generations. The first row is the chromosome and mid-point in the window and the remaining lines are the associated dominance genetic variance estimates for a given window.

## Data Files

---

### **log\_file.txt:**

This file displays a great deal of information on specifics within each generation and it is advisable that one should look over it after you try a new simulation protocol.

### **Master\_DataFrame:**

A file that contains multiple statistics for individuals that survived.

ID: Identification of individual.

Sire: Sire Identification of individual.

Dam: Dam Identification of individual.

Sex: Sex of individual (0 = male and 1 = female).

Gen: Generation the animal was born.

Age: Age the animal was removed from the population either at the culling or selection stage.

Progeny: Number of progeny.

Dead: Number of dead progeny.

Ped\_F: Pedigree based inbreeding metric.

Gen\_F: Diagonal of genomic based relationship constructed based on Van Raden (2008).

Hap1\_F: Diagonal of haplotype 1 based relationship matrix.

Hap2\_F: Diagonal of haplotype 2 based relationship matrix.

Hap3\_F: Diagonal of ROH based relationship matrix.

Homolethal: Number of homozygous lethal genotypes.

Heterlethal: Number of heterozygous lethal genotypes.

Homosublethal: Number of homozygous sub-lethal genotypes.

Hetersublethal: Number of heterozygous sub-lethal genotypes.

Letequiv: Lethal equivalent value.

Homozy: Proportion of the genome homozygous.

Fitness: Multiplicative Fitness value of the individual.

Phen: Phenotype.

EBV: Estimated Breeding Value.

Acc: Accuracy of EBV (not included yet).

GV: True genotypic value of individual ( $\Sigma (a + d)$ ).

BV: True breeding value of individual ( $\Sigma a$ ).

DD: True dominance deviation of individual ( $\Sigma d$ ).

R: Residual value of individual.

### **Master\_Genotypes:**

A file that contains genotypic information for individuals that survived.

ID: Identification of individual.

Marker: Marker genotypes of individual (0-11; 2-22; 3-12; 4-21).

QTL: QTL genotypes of individual (0-11; 2-22; 3-12; 4-21).

### **Marker\_Map:**

chr: Chromosome location.

pos: Nucleotide position of marker.

### **QTL\_new\_old\_Class:**

A file that contains information on QTL effects and frequency across generations.

Chr: Chromosome location.

Pos: Nucleotide position of QTL.

Additive\_Selective: If it is a QTL it refers to the additive effect and if it is a FTL it refers to the selection coefficient.

Dominance: The dominance effect for the QTL or degree of dominance for FTL.  
Type: Refers to the type of loci (2 = quantitative trait; 4 = fitness lethal; 5 = fitness sub-lethal).  
Gen: Generation at which the mutation occurred.  
Freq: Gene frequency across generations with a “\_” as the delimiter.

### **Low\_Fitness:**

This file has a number of metrics for each individual that did not survive to breeding age.  
Sire: Sire ID of individual.  
Dam: Dam ID of individual.  
Gen: Generation the animal was born.  
Ped\_F: Pedigree based inbreeding metric.  
Gen\_F: Diagonal of genomic based relationship constructed based on Van Raden (2008).  
Hap3\_F: Diagonal of ROH based relationship matrix.  
Homozy: Proportion of the genome homozygous.  
Homo lethal: Number of homozygous lethal genotypes.  
Heter lethal: Number of heterozygous lethal genotypes.  
Homo sublethal: Number of homozygous sub-lethal genotypes.  
Heter sublethal: Number of heterozygous sub-lethal genotypes.  
Letequiv: Lethal equivalent value.  
Fitness: Multiplicative Fitness value of the individual.  
GV: True genotypic value of individual ( $\Sigma (a + d)$ ).  
BV: True breeding value of individual ( $\Sigma a$ ).  
DD: True dominance deviation of individual ( $\Sigma d$ ).  
QTL\_Fitness: QTL genotypes of individual (0-11; 2-22; 3-12; 4-21).

## Supplementary Files

---

### **File: CH\*SNP.txt:**

- Haplotype sequence for each chromosome simulated from MaCS.

### **File: MAP\*.txt:**

- map file corresponding to haplotypes sequence in CH\*SNP.txt.

### **File: FounderGenotypes:**

- Genotypes across chromosomes for each founder. The line number corresponds to the founder ID and the first column represents the row number of the two haplotypes that created the genotype, followed by the genotype string.

### **File: G\_Matrix:**

- The Genomic relationship matrix in binary format.

### **File: Ginv\_Matrix:**

- The inverse of the Genomic relationship matrix in binary format.

### **File: Linv\_Matrix:**

- The inverse of the cholesky matrix from the previous generation that is utilized to construct inverse relationship matrix using the method outlined by Meyer et al. (2013) to

obtain the inverse. This file is in binary format.

**File: Pheno\_GMatrix:**

- Dataframe utilized in constructed genomic relationship matrix.

**File: Pheno\_Pedigree:**

- Used in constructing pedigree relationship matrix.

**File: Previous\_Beta\_PCG:**

- Estimates of solutions for previous generation.

**File: SNPFreq:**

- Frequency of SNP across all chromosomes derived from MaCS.

## Appendix II - Generation of Effects

The generation of effects for the quantitative and fitness traits are important parameters that can have a large impact on the simulation results. The methods to generate effects for both types of traits is similar to previous articles and methods other simulation programs have used. At the current time, the sampling of additive effects is from a gamma distribution and dominance effects are simulated from a normal distribution to generate the covariance between quantitative and fitness traits. The use of alternative distributions to generate effects and allow for covariance to occur between the two traits will occur in the future.

### Quantitative Trait:

The additive effect ( $a$ ), defined as half the difference in genotypic value between alternative homozygotes, is generated from a gamma distribution. The default parameters for the gamma distribution (0.4, 1.66) result in an L-shaped distribution of QTL effects and implies that the majority of effects are small and a few have large effects. The gamma distribution only generates positive values, therefore, with equal probability, one of the two alleles was chosen to be positive or negative based on a binomial distribution ( $p = 0.5$ ).

The dominance effect, defined as the deviation of the value of the heterozygote from the mean of the two homozygotes, was generated using a multistep procedure. Independence between additive and dominance effects is the classical treatment (Falconer & Mackay, 1996) and it is convenient because it allows orthogonality of the additive and dominance estimates. However, this independence is contradictory with the phenomena of inbreeding depression and hybrid vigor that indicates dominance is directional (Lynch & Walsh, 1998) and results from real data (Wellmann & Bennewitz 2011; Wellmann & Bennewitz 2012), which suggest an a priori dependency between additive and dominance effects. Therefore, the degree of dominance ( $h$ ) is sampled from a normal distribution, which allows for the user to vary the proportion of positive or negative dominance effect by altering the mean. Next, dominance effects ( $d$ ) were generated by multiplying the degree of dominance by the absolute value of the additive effect ( $d = h|a|$ ). The use of this simulation method results in the additive and dominance effects to now be dependent on each other. Lastly, the choice of parameters specifying the normal distribution and the minor allele frequency for the quantitative QTL has an impact on the proportion of dominance effects that display partial or over-dominance. The proportion that displays partial or over-dominance is outlined in the log file near the beginning of the log file.

**Fitness Trait:**

The generation of fitness effects was divided into lethal and sub-lethal genetic architectures to allow for full flexibility. The distribution of fitness effects and their associated frequency in the genome have been hypothesized to come from two competing results from the literature. The first one is based on the results obtained by (Mukai et al., 1972) and is what we called the “Mukai scenario”, where mutations are assumed to be numerous and of small effect. The second hypothesis is based on more recent results from mutation-accumulation studies and assume that mutations are considerable less frequent but of larger effect (Caballero & Keightley, 1994; Garcia-Dorado & Caballero, 2000). For both lethal and sub-lethal FTL the fitness was defined as relative fitness and is parameterized by two coefficients and they include the selection coefficient ( $s$ ) and the dominance coefficient ( $h$ ). The  $s$  value measures how much worse the unfit allele is, compared to the fittest allele. The  $h$  value measures the degree of dominance that the heterozygote shows regarding the reduced fitness compared to the unfit homozygote (Wright 1931). The normalization procedure forces the fittest homozygote genotype to have a value of 1, and the other homozygote genotype has a value of  $1 - s$ . Lastly, heterozygote genotypes have a fitness value of  $1 - hs$ .

The selection coefficient was generated from a gamma distribution with different parameters for lethal and sublethal. The logfile outlines the mean selection coefficient for the lethal and sub-lethal FTL. As a reference when altering the shape and scale parameter, the mean of a gamma distribution is the shape  $X$  scale.

The dominance coefficient was generated from a normal distribution with different parameters for the lethal and sublethal. The absolute value of the sample is taken as the dominance coefficient. The logfile outlines the mean dominance coefficient for the lethal and sub-lethal FTL. As a reference when altering the shape and scale parameter, the mean of a gamma distribution is the shape  $X$  scale.

The fitness of an individual was then calculated as the multiplicative effect of each fitness genotype across both lethal and sub-lethal FTL with a maximum value of 1 and minimum of 0. A value closer to 1 has a higher fitness and is more likely to survive. In order to simulate environmental stochasticity, a random number was generated from a uniform distribution between 0 and 1 and compared with the fitness value for an individual. If the fitness value was less than the random value from the uniform distribution, the individual did not survive to breeding age and if it was greater than or equal to the animal survived to breeding age.

**Covariance Between Traits:**

The correlation between the quantitative trait and the fitness trait can be due linkage or pleiotropy. Setting the COVAR parameters both to 0 re-

sults in linkage to be the only possible source of correlation between fitness and quantitative traits. Setting the COVAR parameters to a value greater than 0 results in a pleiotropic correlation between the additive effects for the quantitative trait and the selection coefficient for the sub-lethal fitness traits. The scaling of quantitative traits results in the additive effects for the quantitative trait to change and therefore covariance was generated based on Trivariate Reduction algorithm. The Trivariate Reduction algorithm only allows the correlation to be positive. For example, high values for the quantitative trait would result in the two traits being antagonistic based on a positive correlation. One just needs to change the favorable direction of the quantitative trait to alter the interpretation.

Trivariate Reduction for Gamma1 ( $a_1, b_1$ ) and Gamma2 ( $a_2, b_2$ )

- Correlation ( $\rho$ ) Bounded between:  $0 \leq \rho \leq \min(a_1, a_2) / \sqrt{a_1 a_2}$ .

- Steps:

- 1.) Generate  $Y_1 \sim \text{gamma}(a_1 - \rho\sqrt{a_1 a_2}, 1)$
- 2.) Generate  $Y_2 \sim \text{gamma}(a_2 - \rho\sqrt{a_1 a_2}, 1)$
- 3.) Generate  $Y_3 \sim \text{gamma}(\rho\sqrt{a_1 a_2}, 1)$
- 4a.) Generate Value for Gamma1:  $b_1(Y_1 + Y_3)$
- 4b.) Generate Value for Gamma2:  $b_2(Y_2 + Y_3)$

The  $Y_3$  value generate the covariance between the two traits. For FTL that have a covariance with the quantitative trait the  $Y_2$  value is sampled for each FTL within an iteration and the rank correlation is calculated. Once the rank correlation gets within a 1.5 percent of the value specified it then generates the selection coefficient and dominance values using the current iterations  $Y_2$  values.

## Appendix III - MaCS Sequence Simulation

The MaCS program (Markovian Coalescence Simulator; Chen et al. 2009) generates the founder genome. Before using the program, it is advisable to understand the coalescent process and a good review paper is Hudson (1991) and Chapter 5 of Charlesworth & Charlesworth (2010). We have chosen to employ a coalescent simulator (specifically MaCS) in this step due to the flexibility of the approach in generating haplotype sequences for a wide range of population scenarios in terms of the size and structure of the ancestral population across time and genome scenarios with varying mutation and recombination rates. We have chosen the default scenarios to resemble options specified by AlphaSim (Hickey & Gorjanc, 2012) as they represent typical agricultural species LD patterns.

There are 5 default scenarios that represent a range of linkage disequilibrium (LD) patterns that can be called within the simulation, as outlined in the figure below. The five scenarios are called by specifying either “Ne70”, “Ne100\_Scen1”, “Ne100\_Scen2”, “Ne250” or “Ne1000” after the `FOUNDER_Effective.Size` parameter in the parameter file. The user can input a custom effective population size and historical population parameters by using “CustomNe” as the parameter. An easy way to generate custom parameters for MaCS is to utilize a default scenario that resembles the pattern you are wanting and to change the effect population size of the population and determine how the LD pattern changes. If this is specified the program looks for a file called “CustomNe” within the folder where the execution of the program occurred. The file should contain two rows, with the first one being the effective population size parameter and the last one being the historical population size parameters. Lastly, if the program only reads an integer value, then the value is the effective population size with no population history.

The generation of sequence information may take some time to compute. The files generated from the program may be large. Due to this, it is advisable only to generate sequence data once for a given scenario and then adjust narrow-sense heritability, broad sense heritability, selection or mating parameters and start with generating the founder generation.



The default scenarios are below. An illustration of the LD decay associated with each scenario is on the following page:

Ne70:

- Effective population size = “70”.
- Historical population parameters: “-eN 0.18 0.71 -eN 0.36 1.43 -eN 0.54 2.14 -eN 0.71 2.86 -eN 0.89 3.57 -eN 1.07 4.29 -eN 1.25 5.00 -eN 1.43 5.71”.

Ne100\_Scen1:

- Effective population size = “100”.
- Historical population parameters: “-eN 0.06 2.0 -eN 0.13 3.0 -eN 0.25 5.0 -eN 0.50 7.0 -eN 0.75 9.0 -eN 1.00 11.0 -eN 1.25 12.5 -eN 1.50 13.0 -eN 1.75 13.5 -eN 2.00 14.0 -eN 2.25 14.5 -eN 2.50 15.0 -eN 5.00 20.0 -eN 7.50 25.0 -eN 10.00 30.0 -eN 12.50 35.0 -eN 15.00 40.0 -eN 17.50 45.0 -eN 20.00 50.0 -eN 22.50 55.0 -eN 25.00 60.0 -eN 50.00 70.0 -eN 100.00 80.0 -eN 150.00 90.0 -eN 200.00 100.0 -eN 250.00 120.0 -eN 500.00 200.0 -eN 1000.00 400.0 -eN 1500.00 600.0 -eN 2000.00 800.0 -eN 2500.00 1000.0”.

Ne100\_Scen2:

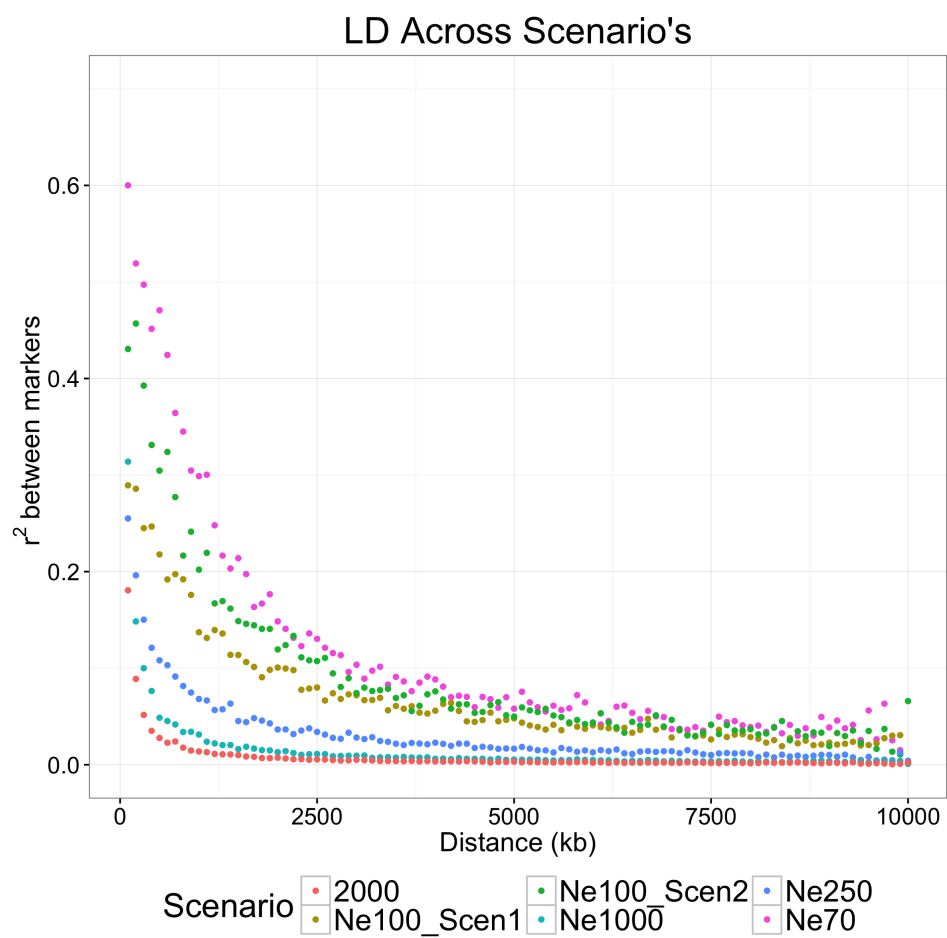
- Effective population size = “100”.
- Historical population parameters: “-eN 50.00 200.0 -eN 75.00 300.0 -eN 100.00 400.0 -eN 125.00 500.0 -eN 150.00 600.0 -eN 175.00 700.0 -eN 200.00 800.0 -eN 225.00 900.0 -eN 250.00 1000.0 -eN 275.00 2000.0 -eN 300.00 3000.0 -eN 325.00 4000.0 -eN 350.00 5000.0 -eN 375.00 6000.0 -eN 400.00 7000.0 -eN 425.00 8000.0 -eN 450.00 9000.0 -eN 475.00 10000.0”.

Ne250:

- Effective population size = “250”.
- Historical population parameters: “-eN 0 1.04 -eN 0 1.08 -eN 0 1.12 -eN 0 1.16 -eN 0.01 1.2 -eN 0.03 1.6 -eN 0.05 2.0 -eN 0.1 2.8 -eN 0.2 4.8 -eN 0.3 5 -eN 0.4 5.2 -eN 0.5 5.4 -eN 0.6 5.6 -eN 0.7 5.7 -eN 0.8 5.8 -eN 0.9 5.9 -eN 1 6 -eN 1 4 -eN 2 8 -eN 3 10 -eN 4 12 -eN 5 14 -eN 6 16 -eN 7 18 -eN 8 20 -eN 9 22 -eN 10 24 -eN 20 28 -eN 40 32 -eN 60 36 -eN 80 40 -eN 100 48 -eN 200 80 -eN 400 160 -eN 600 240 -eN 800 320 -eN 1000 400”.

Ne1000:

- Effective population size = “1000”.
- Historical population parameters: “-eN 0.50 2.00 -eN 0.75 2.50 -eN 1.00 3.00 -eN 1.25 3.20 -eN 1.50 3.50 -eN 1.75 3.80 -eN 2.00 4.00 -eN 2.25 4.20 -eN 2.50 4.50 -eN 5.00 5.46 -eN 10.00 7.37 -eN 15.00 9.28 -eN 20.00 11.19 -eN 25.00 13.10 -eN 50.00 22.66 -eN 100.00 41.77 -eN 150.00 60.89 -eN 200.00 80.00”.



## Literature Cited

- Aguilar, I., I. Misztal, A. Legarra & S. Tsuruta. 2011. Efficient computation of the genomic relationship matrix and other matrices used in single-step evaluation. *J. Anim. Breed. Genet.* 128:422-428.
- Ai H., L. Huang, J. Ren. 2013. Genetic diversity, linkage disequilibrium and selection signatures in chinese and Western pigs revealed by genome-wide SNP markers. *PLoS One* 8(2):e56001.
- Caballero A. & P. Keightley. 1994. A pleiotropic nonadditive model of variation in quantitative traits. *Genetics* 138:883-900.
- de los Campos, G., J. M. Hickey, R. Pong-Wong, H. D. Daetwyler, & M. P. L. Calus. 2013. Whole-Genome Regression and Prediction Methods Applied to Plant and Animal Breeding. *Genetics*, 193, 327-345.
- de Roos, A. P. W., B. J. Hayes, R. J. Spelman, & M. E. Goddard. 2008. Linkage Disequilibrium and Persistence of Phase in Holstein-Friesian, Jersey and Angus Cattle. *Genetics*, 2008, 1503-1512.
- Charlesworth, B., & D. Charlesworth. 2010. *Elements of Evolutionary Genetics*. Greenwood Village, Colorado, USA: Roberts and Company.
- Chen, G. K., P. Marjoram & J. D. Wall. 2009. Fast and flexible simulation of DNA sequence data. *Genome Res.* 19(1):136-42.
- Cheng H., D. J. Garrick, R. L. Fernando. 2015. XSim: simulation of descendants from ancestors with sequence data. *G3* 5:1415-1417.
- Falconer D. S. & T. F. S. Mackay. *Introduction to quantitative genetics*. 4th ed. New York, NY: Longman Scientific and Technical; 1996.
- Garcia-Dorado, A. & A. Caballero. 2000. On the average degree of dominance of deleterious spontaneous mutations. *Genetics*, 155, 1991-2001.
- Hickey, J. M., & G. Gorjanc. 2012. Simulated Data for Genomic Selection and Genome-Wide Association Studies Using a Combination of Coalescent and Gene Drop Methods. *G3* 2:425-427.
- Howard J. T., F. Tiezzi, Y. Huang, K.A. Gray & C. Maltecca. The use of alternative genomic metrics in nucleus herds to manage the diversity of purebred and crossbred animals. *Genet. Sel. Evol.* 48:91.
- Hudson R. R. 1990. Gene genealogies and the coalescent process. *Oxford Surveys in Evolutionary Biology*. 7:1-45.

- Jonas, E. & D. -J. de Koning. 2015. Genomic Selection Needs to Be Carefully Assessed to Meet Specific Requirements in Livestock Breeding Programs. *Front. Genet.*, 6, 49.
- Lynch, M. & B. Walsh. 1998. Genetics and analysis of quantitative traits. Vol. 1. Sunderland, MA: Sinauer.
- McKay, S. D., R. D. Schnabel, B. M. Murdoch, L. K. Matukumalli, J. Aerts, W. Coppieters, D. Crews, E. Dias Neto, C. A. Gill, C. Gao, H. Mannen, P. Stothard, Z. Wang, C. P. Van Tassell, J. L. Williams, J. F. Taylor, S. S. Moore. 2007. Whole genome linkage disequilibrium maps in cattle. *BMC Genet.* 8:74.
- McMahon, B. J., E. C. Teeling, J. Hoglund. 2014. How and Why Should We Implement Genomics into Conservation? *Evol. Appl.*, 7, 999-1007.
- Meuwissen, T. H. E., & Z. Luo. 1992. Computing inbreeding coefficients in large populations. *Genet Sel Evol.* 24(4): 305-313.
- Meyer, K., B. Tier & H. U. Graser. 2013. Technical note: Updating the inverse of the genomic relationship matrix. *J. Anim. Sci.* 91:2583-2586.
- Misztal, I., A. Legarra & I. Aguilar. 2014. Using recursion to compute the inverse of the genomic relationship matrix.
- Morrell, P. L., E. S. Buckler, J. Ross-Ibarra. 2012. Crop genomics: advances and applications. *Nat. Rev. Genet.* 13, 85-96.
- Mukai, T., S. I. Chigusa, L. E. Mettler, J. F. Crow. 1972. Mutation rate and dominance of genes affecting viability in *Drosophila melanogaster*. *Genetics* 72:333-355.
- Pérez-Enciso M., & A. Legarra. 2016. A combined coalescence gene-dropping tool for evaluating genomic selection in complex scenarios (ms2gs). *J. Anim. Breed. Genet.* 133(2):85-91.
- Porto-Neto L. R., T. S. Sonstegard, G. E. Liu, D. M. Bickhart, M. V. Da Silva, M. A. Machado, Y. T. Utsunomiya, J. F. Garcia, C. Gondro, C. P. Van Tassell. 2013. Genomic divergence of zebu and taurine cattle identified through high-density SNP genotyping. *BMC Genomics.* 14:876.
- Sargolzaei, M. & F. S. Schenkel. 2009. QMSim: a large-scale genome simulator for livestock. *Bioinformatics*, 25, 680-1.
- VanRaden P. M. 2008. Efficient Methods to Compute Genomic Predictions. *J. Dairy Sci.*, 91, 4414-4423.

- VanRaden P. M., K. M. Olson, D. J. Null, J. L. Hutchison. 2011. Harmful recessive effects on fertility detected by absence of homozygous haplotypes. *J. Dairy Sci.*, 94, 6153-6161.
- Wellmann, R. & J. Bennewitz. 2011. The contribution of dominance to the understanding of quantitative genetic variation. *Genet Res (Camb)* 93:139?154.
- Wellmann, R. & J. Bennewitz. 2012. Bayesian models with dominance effects for genomic evaluation of quantitative traits. *Genet Res (Camb)* 94:21?37.
- Wright, S. 1931. Evolution in Mendelian populations. *Genetics* 16:97-159.
- Yang, J., B. Benyamin, B. P. McEvoy, S. Gordon, A. K. Henders, D. R. Nyholt, P. A. Madden, A. C. Heath, N. G. Martin, G. W. Montgomery, M. E. Goddard, P. M. Visscher. 2010. Common SNPs explain a large proportion of the heritability for human height. *Nat Genet.*, 42, 565-569.

## Example 1: Running the Simulation Software

### Quantitative Trait - Single Progeny

```
-----| Running the Program Example |-----
-| General |-
START: sequence
SEED: 1501
-| Genome & Marker |-
CHR: 3
CHR_LENGTH: 150 150 150
NUM_MARK: 4000 4000 4000
QTL: 50 50 50
-| Population |-
FOUNDER_Effective_Size: Ne70
MALE_FEMALE_FOUNDER: 50 400 random 0
VARIANCE_A: 0.35
VARIANCE_D: 0.05
-| Selection |-
GENERATIONS: 15
INDIVIDUALS: 50 0.2 400 0.2
PROGENY: 1
SELECTION: ebv high
EBV_METHOD: pblup
CULLING: ebv 5
-| Mating |-
MATING: random125 simu.anneal
```

To run the program place, "GenoDiver", "macs" and "msformatter" executable files in the folder where the program will run. Before running the program, the file permissions need to be checked. After verifying the permissions, a parameter file, outlined above, will need to be generated and placed in the same folder as the previous three executable files. A parameter file can be generated using any text editor or downloaded from the Geno-Diver GitHub page.

The simulation program reads the parameter file by searching for keywords that are capitalized and followed by a colon. Therefore any phrase that does not meet the search criteria is ignored when initializing parameters within the program. Also, if you want to comment out a parameter just add "!!" within the key word and the program will skip over it. For example to skip over the "SEED" parameter just replace it with "SE!!ED"

and it won't recognize the parameter any more.

To run the program type `./GenoDiver` and the name of your parameter file. For example, if the parameter file is named `"parameterfile"`, the program is run by typing in `./GenoDiver parameterfile`. During the simulation, the program outputs minimal comments on the progress. A more thorough description of the program's status is printed to the log file (i.e. `"log_file.txt"`).

### Parameter File Summary

Sequence information is generated for three chromosomes with a length of 150 Megabases. The genome simulated has a high degree of short-range LD ( $N_e70$ ). The SNP panel contains 12,000 marker (i.e. 4,000 markers per chromosome). For each chromosome, 50 randomly placed QTL were generated, and no FTL mutations were generated. The narrow and broad sense heritability for the trait is 0.35 and 0.40, respectively. The phenotypic variance is by default set at 1.0, and therefore the residual variance is 0.6. The founder population consisted of 50 males and 400 females. For each generation, a total of 50 males and 400 females are in the population. A total of 10 and 80 (0.2 replacement rate) male and female parents, respectively, are culled and replaced by new progeny each generation. A total of 15 generations were simulated. Animals with a high estimated breeding value (EBV) were selected or culled each generation. The EBV are solved using an animal model based on a pedigree-based relationship matrix. Each mating pair produced one progeny. Parents that had pedigree-based relationships greater than 0.125 were not avoided, and this was optimized based on the simulated annealing method.

### Overview of Results

After the program has finished it is strongly recommended to check the parameters initialized at the top of the log file (i.e. `"log_file.txt"`) within the output folder. The log file contains a large amount of information and is a great tool to ensure that the parameters and outcome of the simulation is what is intended. If the program is not running correctly the log file should provide knowledge on why and where the simulation crashed or exited. The files are by default placed in the `"GenoDiverFiles"` directory. If the `"OUTPUTFOLDER"` option is utilized, the files will be placed in the user-specified directory. Below is a screenshot of the files that are generated from the simulation software based on the parameter file outlined above.

A number of files are generated, but only a few are needed to generate summary statistics on the simulation program and include:

- **Master\_DataFrame**: File with phenotype, inbreeding and pedigree information across all animals.
- **Master\_Genotypes**: File with genotype information for each animal.
- **QTL\_new\_old\_Class**: File with information for each QTL/FTL mutation.
- **Marker\_Map**: Location of markers.
- **Summary\_Statistics\_DataFrame\_Performance**: Summary statistics by generation on performance metrics.
- **Summary\_Statistics\_DataFrame\_Inbreeding**: Summary statistics by generation on inbreeding metrics.
- **Summary\_Statistics\_QTL**: Summary statistics by generation QTL/FTL metrics.

#### Overview of Results:

- When simulating dominance be sure to check to see if the number that display over-dominance or partial dominance is what you were wanting. In this case with the given narrow and broad sense heritability, the proportion of QTL displaying over-dominance is small (around 15 %).
- The table below outlines the change in multiple parameters as the generations proceed

Generation	Phenotype	EBV	Pedigree Inbreeding	Genomic Inbreeding
0	0.000	0.000	0	0.978
1	0.004	0.004	0	0.997
2	0.308	0.254	0	0.999
3	0.459	0.476	0	0.999
4	0.765	0.752	0.0008	1.007
5	1.030	1.014	0.0022	1.018
6	1.207	1.214	0.0078	1.021
7	1.507	1.463	0.0143	1.047
8	1.677	1.670	0.0200	1.058
9	1.887	1.871	0.0251	1.066
10	2.099	2.063	0.0296	1.080
11	2.246	2.256	0.0350	1.099
12	2.519	2.465	0.0389	1.111
13	2.617	2.619	0.0435	1.123
14	2.765	2.791	0.0485	1.148
15	2.965	2.971	0.0510	1.153



## Example 2: Quantitative Trait - Multiple Progeny

```
-----|   Multiple Progeny Quantitative Trait   |-----
-| General |-
START: sequence
SEED: 1501
-| Genome & Marker |-
CHR: 3
CHR_LENGTH: 150 150 150
NUM_MARK: 4000 4000 4000
QTL: 50 50 50
-| Population |-
FOUNDER_Effective_Size: Ne70
MALE_FEMALE_FOUNDER: 25 200 random 0
VARIANCE_A: 0.35
VARIANCE_D: 0.05
-| Selection |-
GENERATIONS: 15
INDIVIDUALS: 25 0.2 200 0.2
PROGENY: 8
MAXFULLSIB: 2
SELECTION: ebv high
EBV_METHOD: pblup
CULLING: ebv 5
-| Mating |-
MATING: random125 simu_anneal
```

The parameter file outlined below illustrates how to simulate a prolific species (i.e. Swine) with a threshold on the minimum number of progeny to select within each family. Once the program has finished, inspection of the log file will provide details on the impact of the “MAXFULLSIB” option. Within the log file, after the “Begin ebv Selection of offspring” section, the number of times “n” number of siblings were selected within a family is outlined for each generation. This provides an overview of the degree of co-selection of siblings that occur based on the simulation design outlined above. If this option would not have been included, the maximum number of siblings that can be selected within a family would be 8. Given the low narrow-sense heritability of the trait, the number of siblings selected within a family would be higher if this parameter wasn’t included.

### Parameter File Summary

Sequence information is generated for three chromosomes with a length of

150 Megabases. The genome simulated has a high degree of short-range LD ( $N_e70$ ). The SNP panel contains 12,000 marker (i.e. 4,000 markers per chromosome). For each chromosome, 50 randomly placed QTL were generated, and no FTL mutations were generated. The narrow and broad sense heritability for the trait is 0.10 and 0.105, respectively. The phenotypic variance is by default set at 1.0, and therefore the residual variance is 0.895. The founder population consisted of 25 males and 200 females. For each generation, a total of 25 males and 200 females are in the population. A total of 5 and 40 (0.2 replacement rate) male and female parents, respectively, are culled and replaced by new progeny each generation. A total of 10 generations were simulated. Animals with a high estimated breeding value (EBV) were selected or culled each generation. The EBV are solved using an animal model based on a pedigree-based relationship matrix. Each mating pair produced eight progeny. Within each full-sib family a maximum of 2 progeny can be selected. Parents that had pedigree-based relationships greater than 0.125 were not avoided, and this was optimized based on the simulated annealing method.

## Overview of Results

Once the program has finished, inspection of the log file will provide details on the impact of the “MAXFULLSIB” option. Within the log file, after the “Begin ebv Selection of offspring” section, the number of times “n” number of siblings were selected within a family is outlined for each generation. This provides an overview of the degree of co-selection of siblings that occur based on the simulation design outlined above. If this option would not have been included, the maximum number of siblings that can be selected within a family would be 8. Given the low narrow-sense heritability of the trait, the number of siblings selected within a family would be higher if this parameter wasn’t included.

### Example 3: Differential Sire Contribution by Age

```

-----| Differential Sire Contribution by Age |-----
-| General |-
START: sequence
SEED: 1501
-| Genome & Marker |-
CHR: 3
CHR_LENGTH: 150 150 150
NUM_MARK: 4000 4000 4000
QTL: 50 50 50
-| Population |-
FOUNDER_Effective_Size: Ne70
MALE_FEMALE_FOUNDER: 50 400 random 0
VARIANCE_A: 0.35
VARIANCE_D: 0.05
-| Selection |-
GENERATIONS: 15
INDIVIDUALS: 50 0.2 400 0.2
PROGENY: 1
PARITY_MATES_DIST: 2.0 1.0
SELECTION: ebv high
EBV_METHOD: pblup
CULLING: ebv 5
-| Mating |-
MATING: random125 simu_anneal

```

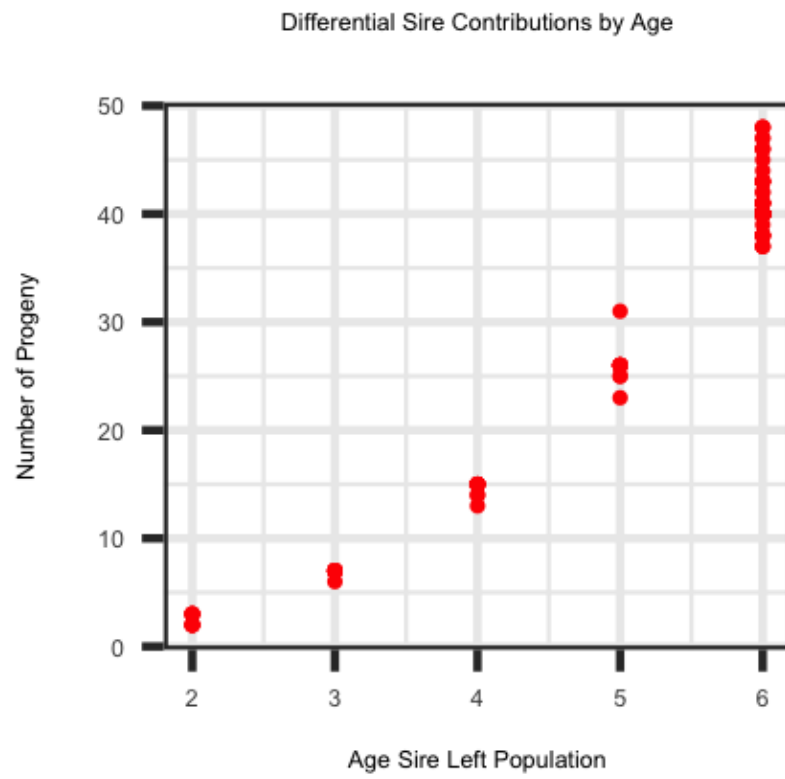
The parameter file outlined above illustrates how to simulate a population where older animals are assigned a larger number of mating pairs compared to younger animals. This type of scenario is generated by utilizing the “PARITY\_MATES\_DIST” parameter with the following values “2.0 1.0”. These two values are utilized to generate the distribution of mating pairs. A Beta distribution, which is parameterized by two parameters, is used to generate the distribution of mating pairs. A beta distribution was utilized in order to allow for a wide range of mating scenarios. For example, to generate a scenario where younger animals are assigned a larger number of mating pairs compared to older animals, the parameters need to be changed to “1.0 2.0”. The number of mating pairs by age class are generated by splitting the cumulative distribution function (CDF) into quadrants based on the number of age classes that occur within a generation. The total number of mating pairs within an age class is the proportion that falls within the CDF quadrant for a given age class.

## Parameter File Summary

Sequence information is generated for three chromosomes with a length of 150 Megabases. The genome simulated has a high degree of short-range LD ( $N_e70$ ). The SNP panel contains 12,000 marker (i.e. 4,000 markers per chromosome). For each chromosome, 50 randomly placed QTL were generated, and no FTL mutations were generated. The narrow and broad sense heritability for the trait is 0.35 and 0.40, respectively. The phenotypic variance is by default set at 1.0, and therefore the residual variance is 0.6. The founder population consisted of 50 males and 400 females. For each generation, a total of 50 males and 400 females are in the population. A total of 10 and 80 (0.2 replacement rate) male and female parents, respectively, are culled and replaced by new progeny each generation. A total of 15 generations were simulated. Animals with a high estimated breeding value (EBV) were selected or culled each generation. The EBV are solved using an animal model based on a pedigree-based relationship matrix. Each mating pair produced one progeny. The mating distribution was skewed so that older animals had more mating pairs than younger animals. Parents that had pedigree-based relationships greater than 0.125 were not avoided, and this was optimized based on the simulated annealing method.

## Overview of Results

Once the program has finished, inspection of the log file will provide details on the impact of the "PARITY\_MATES\_DIST" option. Within the log file (lines 148-169), the mating distribution CDF is illustrated and is outlined below. When generating the number of matings for a given age class for each generation the CDF outlined below is split into quadrants based on the number of age classes that occur within a generation. A potential reasons for putting this into a simulation is to generate some sires with a large number of progeny while other sires have very few progeny. The sires that generate a large number of progeny would then have a large impact on the genome of future generations which may include the spread of a lethal/sublethal mutation that the sire(s) carry. The plot below depicts the non-linear relationship of number of progeny left by a sire across different ages at which a sire left the herd.



## Example 4 Fitness Trait

```
-----|   Simulating a Fitness Trait   |-----
-| General |-
START: sequence
SEED: 1500
-| Genome & Marker |-
CHR: 3
CHR.LENGTH: 150 150 150
NUM.MARK: 4000 4000 4000
QTL: 0 0 0
FIT.LETHAL: 20 20 20
FIT.SUBLETHAL: 200 200 200
-| Population |-
FOUNDER_Effective_Size: 500
MALE_FEMALE_FOUNDER: 150 600 random 0
VARIANCE_A: 0.0
VARIANCE_D: 0.0
COVAR: 0.5 0.2
-| Selection |-
GENERATIONS: 25
INDIVIDUALS: 50 0.2 400 0.2
PROGENY: 1
SELECTION: random high
CULLING: random 10
-| Mating |-
MATING: random
-| Output Options |-
OUTPUT_LD: yes
GENOTYPES: no
```

The parameter file below illustrates how to simulate a fitness trait only. When simulating a fitness trait it is important to ensure that you have enough founder individuals because a portion will not make it to breeding age. Therefore, an extra 100 males and 200 females were added to the founder population to ensure enough individuals are available to generate the breeding population. This will also impact the replacement rate because if enough progeny aren't available to remain at the chosen male and female population size the simulation will exit. A full description of how the fitness value of an individual impacts its ability to make it to breeding age is described in the ["QTL/FTL Distribution Parameters"](#) link.

## Parameter File Summary

Sequence information is generated for three chromosomes with a length of 150 Megabases. The genome simulated has a high degree of short-range LD (250). The SNP panel contains 12,000 marker (i.e. 4,000 markers per chromosome). For each chromosome, 20 lethal and 200 sub-lethal mutations were generated. The quantitative trait has a broad sense heritability of 0.0 and therefore an animals phenotype is only a function of random environmental deviations with a variance of 1.0. The founder population consisted of 150 males and 600 females. For each generation, a total of 50 males and 400 females are in the population. A total of 10 and 80 (0.2 replacement rate) male and female parents, respectively, are culled and replaced by new progeny each generation. Starting at the first generation animals are randomly selected or culled each generation. The maximum number of generations an animal can remain in the breeding population is 10. Each mating pair produced one progeny and parents were mated at random.

### Overview of Results:

- The number of progeny that died due to fitness and the number of the FTL purged from the population is in Summary\_Statistics\_QTL file.
- The log file contains the mean selection coefficient and degree of dominance for the lethal and sublethal FTL.
- When simulating fitness effects, a certain proportion of the progeny will die due to fitness and if the population does not have enough progeny to stay at the value given the program exits. Therefore, careful consideration of the number and magnitude of fitness effects along with the number of progeny produced per mating pair needs to be carefully considered when constructing the parameter file.

## Example 5 (Quantitative Trait + Fitness Correlated)

```

-----| Quantitative and Fitness Trait |-----
-| General |-
START: sequence
SEED: 1501
-| Genome & Marker |-
CHR: 3
CHR.LENGTH: 150 150 150
NUM_MARK: 4000 4000 4000
QTL: 100 100 100
FIT.LETHAL: 15 15 15
FIT.SUBLETHAL: 100 100 100
-| Population |-
FOUNDER_Effective_Size: Ne70
MALE_FEMALE_FOUNDER: 100 350 random 0
VARIANCE_A: 0.20
VARIANCE_D: 0.05
COVAR: 0.5 0.2
-| Selection |-
GENERATIONS: 20
INDIVIDUALS: 50 0.2 250 0.2
PROGENY: 1
SELECTION: ebv high
EBV_METHOD: pblup
CULLING: ebv 5
-| Mating |-
MATING: random

```

The parameter file outlined below illustrates how to simulate a quantitative trait that has a proportion (i.e. 50%) of the quantitative trait loci (QTL) also having a fitness effect. The relationship between the QTL with a fitness effect has a positive correlation of 0.20, but the two traits are antagonistic based on the way the fitness value of a loci is parameterized. High selection coefficients (s) result in the unfavorable homozygote genotype to be less fit and is described in detail in the 'QTL/FTL Distributions' link. As outlined in the previous example, it is important to add a few extra animals in the founder population since a portion will die. If the number of male or female founder animals is smaller than the population size it will exit out of the program.



### Parameter File Summary

Sequence information is generated for three chromosomes with a length of 150 Megabases. The genome simulated has a high degree of short-range LD ( $N_e70$ ). The SNP panel contains 12,000 marker (i.e. 4,000 markers per chromosome). For each chromosome, 150 randomly placed QTL were generated. Also, 15 lethal and 100 sub-lethal mutations were generated. The narrow and broad sense heritability for the quantitative trait is 0.20 and 0.25, respectively. The phenotypic variance is by default set at 1.0, and therefore the residual variance is 0.75. Half of the QTL also have a fitness effect and the correlation between the additive QTL effects and sub-lethal selection coefficients is 0.20. The founder population consisted of 100 males and 350 females. For each generation, a total of 50 males and 250 females are in the population. A total of 10 and 80 (0.2 replacement rate) male and female parents, respectively, are culled and replaced by new progeny each generation. Starting at the first generation animals with a high EBV were selected or culled each generation. The EBV are estimated using a pedigree-based BLUP utilizing all the animals. Each mating pair produced one progeny. Parents were mated at random.