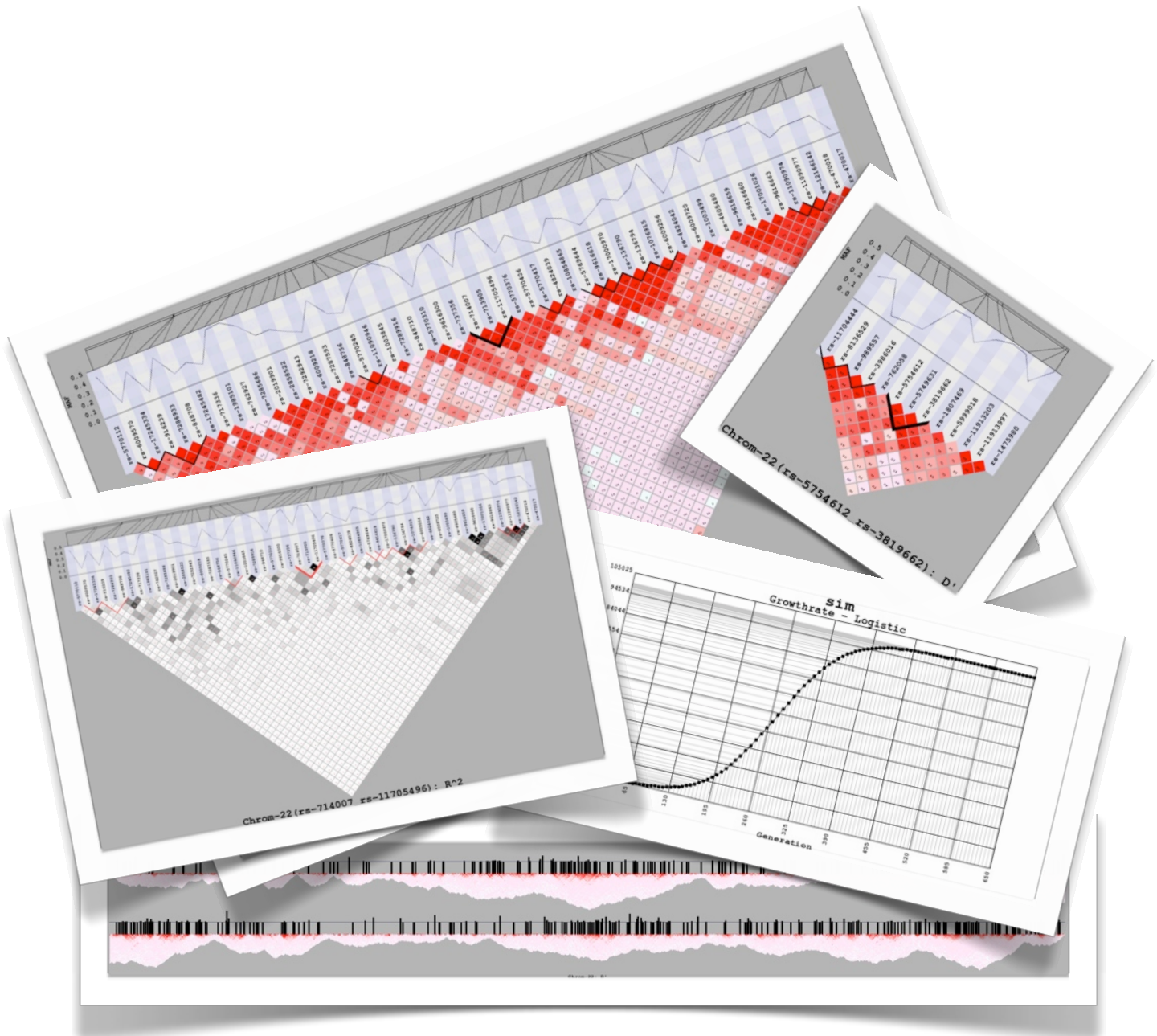


g e n o m e S I M L A  
Reference Manual rev. 1.0.1



genomeSIMLA - A forward time simulation for genetic data

# Table of Contents

<b>Introduction</b>	1
Purpose of this manual	1
Conventions Used	1
Random Numbers	1
Common Parameters	1
Integer	2
Float	2
Index	2
max	2
min	2
On/Off	2
filename	2
label	2
description	2
<b>Using genomeSIMLA</b>	3
Command-Line Arguments	3
genomeSIMLAs config-file [ld   datasets] [-p project-name] [-l Integer] [-d Integer Integer Integer] [-s Integer]	3

Config-file	3
ld (optional)	3
datasets (optional)	3
-p project-name (optional)	3
-l generation-to-load (Integer) (optional)	3
-d first-generation-to-drop generations-between-drops drop-count (optional)	3
-s seed (integer)	3
<b>General Parameters</b>	<b>4</b>
The following parameters control the basic behavior of the application.	4
SEED integer	4
SEED 23125	4
<b>Drop Points</b>	<b>4</b>
FIRST_DROP_POINT integer	4
<code>FIRST_DROP_POINT 500</code>	4
DROP_FREQUENCY integer	4
<code>DROP_FREQUENCY 100</code>	4
DROP_COUNT	4
<code>DROP_COUNT 5</code>	4
<b>Graphical Plot Settings</b>	<b>4</b>
MAX_SNPS_PER_ROW integer	4
<code>MAX_SNPS_PER_ROW 3000</code>	4
BLOCK_REPORT_SIZE integer	4
<code>BLOCK_REPORT_SIZE 30</code>	4
FONT filename	5

FONT ./FreeMonoBold.ttf	5
<b>CSS_FILENAME filename</b>	5
CSS_FILENAME ../genomesimla.css	5
<b>LD_REPORT_BUFFER_SIZE On/Off</b>	5
LD_BUFFER_SIZE 50	5
<b>WRITE_LD_REPORT On/Off</b>	5
WRITE_LD_REPORT On	5
<b>DRAW_RSQUARED_PLOTS On/Off</b>	6
<b>DRAW_DPRIME_PLOTS On/Off</b>	6
DRAW_RSQUARED_PLOTS off	6
<b>MAX_SNP_DISTANCE integer</b>	6
MAX_SNP_DISTANCE 500000	6
<b>CLOSE_POOLS_BETWEEN_DROPS On/Off</b>	6
CLOSE_POOLS_BETWEEN_DROPS On	6
<b>FAST_LD_POOL_SIZE integer</b>	6
FAST_LD_POOL_SIZE 3000	6
<b>FAST_LD_PLOT_SIZE</b>	6
FAST_LD_PLOT_SIZE 10000	6
NO_FASTLD	6
<b>Locus Generation</b>	7
<b>Block Based Locus Generation</b>	7
<b>DEFAULT_BLOCK min max float float float float</b>	7
DEFAULT_BLOCK 5 10 0.01 0.015 0.00001 0.000025	7
<b>ADD_CHROMOSOME integer label</b>	8
ADD_CHROMOSOME 5 chromosome_1	8

ADD_BLOCK chr_idx snp_idx float float float float float	8
ADD_BLOCK 5 10 0.000001 0.00001 0.5	8
Other Block Related Settings	8
DEFAULT_ALLELE_FREQ float float	8
DEFAULT_ALLELE_FREQ 0.1 0.5	8
File Based Chromosome Configuration	8
Locus File Format	8
Locus Miscellany	9
ALLELE_FREQUENCY chr_idx snp_idx float float	9
ALLELE_FREQUENCY 1 5 0.25 0.75	9
Population Control	10
GROWTH_RATE LINEAR initial_population variation growth_rate	10
GROWTH_RATE LINEAR 30000 0.05 10.0	10
GROWTH_RATE EXPONENTIAL initial_population variation growth_rate	10
GROWTH_RATE EXPONENTIAL 700 0.05 0.3	10
GROWTH_RATE LOGISTIC initial_population variation growth_rate carry- ing_capacity	11
GROWTH_RATE LOGISTIC initial_population variation growth_rate carrying_ca- capacity	11
GROWTH_RATE RICHARDS initial_population variation growth_rate carry- ing_capacity time_of_max_growth polarity	11
GROWTH_RATE LOGISTIC initial_population variation growth_rate carrying_ca- capacity time_of_max_growth polarity	11
General Growth rate parameters	11
MAX_POOL_SIZE integer	11
MAX_POOL_SIZE 90000	11

MIN_POOL_SIZE integer	11
MIN_POOL_SIZE 1500	11
TARGET_POP_SIZE integer	11
TARGET_POP_SIZE 100000	11
<b>Dataset Generation</b>	<b>12</b>
Case/Control	12
DATASET CC label affected unaffected genotype_error phenocopy missing	12
DATASET CC sample-01 500 500 0.05 0.1 0.15	12
Pedigree Data	12
DATASET PED label genotype_error phenocopy missing	12
DATASET PED family-01 0.05 0.1 0.15	12
DATASET FAMTYPE affected unaffected extra_sibs number_of_families	12
DATASET FAMTYPE 1 1 1 250	12
DATASET FAMTYPE 1 0 0 150	13
DATASET FAMTYPE 2 1 0 75	13
DATASET FAMTYPE 1 0 3 50	13
General Data-set Configuration Parameters	13
DATASET_COUNT integer	13
DATASET_COUNT 500	13
BINARY_DATASETS Yes/No	13
BINARY_DATASETS Yes	13
USE_STD_PEDIGREE_HEADER on/off	13
USE_STD_PEDIGREE_HEADER On	13
<b>Locus Searching</b>	<b>14</b>

LOCUS_SELECTOR label float float float integer integer integer de- scription	14
LOCUS_SELECTOR rare_loci 0.2 0.15 0.23 4 2 10 The following loci are moderately rare and appear in a block	14
ADD_REGION label snp_start snp_stop	14
ADD_REGION rare_loci rs-321412 rs-543231	14
MAX_LOCI_PER_CHROM_REPORTED Integer	14
MAX_LOCI_PER_CHROM_REPORTED 50	14
<b>Disease Modeling</b>	<b>15</b>
<b>Penetrance Table Disease Models</b>	<b>15</b>
DEFINE_MODEL PENTABLE INDEX pen-file chrom-id snp-id [chrom-id snp- id, ...]	15
DEFINE_MODEL PENTABLE INDEX disease.pen 1 5	15
DEFINE_MODEL PENTABLE LABEL pen-file snp-label [snp-label, ...]	15
DEFINE_MODEL PENTABLE LABEL disease.pen RL5	15
<b>Penetrance File Configuration</b>	<b>15</b>
FREQ_THRESHOLD Float	15
FREQ [AaBbCcDd....etc] float	15
FREQ A 0.2	15
FREQ a 0.8	15
<b>PENTABLE</b>	<b>16</b>
model-identification penetrance	16
AABB 0.171	16
AABb 0.155	16
<b>Purely Epistatic Models with simPEN</b>	<b>16</b>

DEFINE_MODEL SIMPEN INDEX simpen-cfg chrom-id snp-id [chrom-id snp-id, ...]	16
DEFINE_MODEL SIMPEN INDEX disease.simpenn 1 5	16
DEFINE_MODEL PENTABLE LABEL simpen-cfg snp-label [snp-label, ...]	16
DEFINE_MODEL PENTABLE simpen-cfg disease.penn RL5	16
simPEN File Configuration	16
HERIT float	16
HERIT 0.01	16
HERITWEIGHT float	16
HERITWEIGHT 10	16
ODDSRATIO float	16
ODDSRATIO 1.25	16
ODDSWEIGHT float	17
ODDSWEIGHT 1	17
MARGVAR float	17
MARGVAR 0.0000001	17
MARGWEIGHT float	17
MARGWEIGHT 100	17
PENTARGET float	17
PENTARGET 0.15	17
GEN 15000	17
POPSIZE 1000	17
DEMES 100	17
MUTATE 0.01	17
CROSS 0.6	17



SUBMODELS ON	17
UPDATE 100	17
LOCI 2	18
FREQ 0.2 0.8	18
<b>Main Effects and Interactions with SIMLA</b>	<b>18</b>
<pre> DEFINE_MODEL SIMLA INDEX simla-cfg float int chrom-id snp-id MIN   MAJ float float [...] </pre>	18
<pre> DEFINE_MODEL SIMLA INDEX interactions.simla 0.05 2 1 5 MIN 0.26 0.0 </pre>	18
<pre> DEFINE_MODEL SIMLA LABEL simla-cfg float int snp-label MIN   MAJ float float [...] </pre>	18
<pre> DEFINE_MODEL SIMLA LABEL interactions.simla 0.05 2 RL5 MIN 0.26 0.0 </pre>	18
<b>SIMLA configuration file</b>	<b>18</b>
<pre> 1x2x3 0.26 </pre>	18

# Introduction

## Purpose of this manual

Contained within this manual are details for configuring and running the application, genomeSIMLA. If this is your first time to use the software, we highly recommend that you take a few minutes to download and work through one or more tutorials. Then, once familiar with the capabilities of the software, users can refer to this guide when making changes to the basic configuration settings.

## Conventions Used

There are two conventions used throughout this document. These text conventions are intended to help distinguish examples from configuration parameters.

## Random Numbers

genomeSIMLA uses an open source implementation of the mersenne-twister pseudo random number generator available at <http://agner.org/random>. When using genomeSIMLA to generate data, the following should be kept in mind in order to ensure that products are as reproducible as possible:

- At the beginning of execution of any kind (population initialization, generational advancement, dataset extract, etc) the random seed will be set.
- 

## Configuration details are listed first in bold- left aligned with the rest of the text.

The first word(s) are the keywords which specify what is being changed. Each keyword (or phrase) has some number of parameters. These are listed in the order they should appear in the configuration line. In some cases, parameters can be repeated or are optional. Those are denoted inside [].

Configuration details are generally followed immediately by an example line:

*This is an example*

Examples show how an actual entry would look and are followed by some descriptive information to help the user understand how the example would affect genomeSIMLA's runtime.

## Common Parameters

There are a number of parameters which are used commonly across multiple configuration settings. In order to simplify the descriptions of the various properties of each command, we'll describe those properties here, and just refer to them as if they were a type.

**Integer**

Parameters specified in this way just simply refer to a whole number. In general, these values should be equal to or greater than 0, except when specified otherwise.

**Float**

Values specified as float are decimal values.

**Index**

If a parameter is listed as an index, it refers to the index, starting at 1 the user wishes to select.

**max**

This is generally an integer value representing the upper bound of some value. In some cases, such as minor allele frequency, it might represent a floating point value.

**min**

This is generally an integer value representing the lower bound of some value. In some cases, such as minor allele frequency, it is possible that it represents a floating point value.

**On/Off**

These parameters accept a boolean, Yes/No type setting. Users can use ON/OFF or YES/NO to set them.

**filename**

When a configuration refers to a file for input or output, the filename is generally used. This can be either a fully qualified path (such as /home/torstees/wga) or it can be specified as a path relative to the directory where the application was run (such as ../data/goodfilename). It can also be just a plain filename as long as the file itself is available from the directory in which the application was run.

**label**

A label refers to a parameter whose value can be any text string without whitespace. These labels are generally used for reporting but in many cases are used to determine filenames. As a result, users should avoid using unusual characters in the string that could possibly cause problems with filenames. Because spaces and tabs are used to separate each parameter on a given line, labels can not contain spaces.

**description**

A description is a chunk of text that can contain spaces. It will always be at the very end of a line and is generally optional.

# Using genomeSIMLA

Except in very specific cases, generating data-sets with genomeSIMLA is a multistep process. At the very least, users must run genomeSIMLA forward through time performing at least 1 drop along the way. It is this drop that the user's data-sets will be drawn from. In addition to generational advancement and data-set production, genomeSIMLA can pick up from a specified generation and advance further through time, or perform "complete" LD analysis.

To control genomeSIMLA in this way, we offer a small number of different parameters to give the user control over genomeSIMLA's behavior. It is important to note that a handful of these parameters must appear in a certain order (those that lack a -T flag, where T is some parameter designator).

## Command-Line Arguments

**genomeSIMLAs config-file [ld | datasets] [-p project-name] [-l Integer] [-d Integer Integer Integer] [-s Integer]**

### Config-file

Specifies the filename to be used to control genomeSIMLAs overall specific behavior. If the configuration is available from within the current working directory, the filename alone is sufficient. If the filename exists in another directory, a fully qualified or relative path should be provided along with the filename itself.

### ld (optional)

When the ld command is present, no generational advancement will be performed, and complete LD analysis will be performed on the specified pool. If no generation is specified via the -l flag, generation 0 is assumed. All other commands are ignored in the presence of this flag.

### datasets (optional)

When the datasets command is present, no generational advancement will be performed, and data-sets will be drawn from the specified pool. If no generation is specified via the -l flag, generation 0 is assumed. All other commands are ignored in the presence of this flag.

### -p project-name (optional)

Specifying a project name allows the user to override the nature behavior of using the name of the configuration file as the base name for all of the products generated by execution. This can include a relative or fully qualified path, as long as a base filename is present (i.e. data/affy or /home/torstees/simulated\_data/affy ) All files generated will start with this string.

### -l generation-to-load (Integer) (optional)

Specifies the generation to load. This assumes that a previous run has been completed, and pools at the specified generation were created.

### -d first-generation-to-drop generations-between-drops drop-count (optional)

This allows the user to override the drop configuration found in the configuration file.

### -s seed (integer)

Allows the user to override the seed specified in the configuration.

# General Parameters

The following parameters control the basic behavior of the application.

## SEED integer

```
SEED 23125
```

Sets the seed for all random number calls. Seeds can range from 0 - 4.2 billion.

## Drop Points

Drop points are points in simulated time (generations) where the entire contents of the pool(s) is written to disk and analyzed. Reports are produced in HTML format to help the user to interpret the current state of the pool. Drop points are designed to allow the user to track the state of LD within the population. If the population is large enough, any drop point can be the source for dataset generation.

The reports initially written for a given generation are done using sampling. Prior to selecting loci for modeling diseases, users are expected to extract detailed reports from the generation of interest.

## FIRST\_DROP\_POINT integer

```
FIRST_DROP_POINT 500
```

Sets the first drop point to be performed at generation 500.

## DROP\_FREQUENCY integer

```
DROP_FREQUENCY 100
```

Causes genomeSIMLA to drop every N generations once it has reached the initial drop point. The example says to drop every 100 generations.

## DROP\_COUNT

```
DROP_COUNT 5
```

Indicates the total number of drops to be performed (including the initial drop). If we look at all three of the previous DROP related examples, genomeSIMLA would perform 5 drops at generations: 500, 600, 700, 800 and 900.

It should be noted that the calling parameters can change how drop points are interpreted.

## Graphical Plot Settings

The following general parameters control various aspects of the graphical reporting during a given drop.

## MAX\_SNPS\_PER\_ROW integer

```
MAX_SNPS_PER_ROW 3000
```

In order to render a general overview of an entire chromosome, it is necessary to set a maximum number of SNPS which can be drawn on a single row. This is not a hard setting- more of a suggestion. genomeSIMLA will distribute the SNPS evenly on all rows in order to avoid having a small chunk at the bottom.

## BLOCK\_REPORT\_SIZE integer

```
BLOCK_REPORT_SIZE 30
```

Determines the number of detailed blocks that are reported. Each report takes disk space as well as time to generate. If your configuration is set to use 22 chromosomes with a report size of 30, there will likely be over 1200 charts drawn and a fair amount of information added to the final report. However, if the settings are too small, it might be more difficult to find the preferred SNP.

Each *Block Report* consists of 2 graphs. The first is expected to be smaller (fewer SNPS on either side of the block of interest). The second is generally zoomed out to show a bit more of the surrounding SNPS.

#### FONT filename

```
FONT ./FreeMonoBold.ttf
```

genomeSIMLA requires access to a true type font in order to write labels and details onto the graphical portions of the reports. This font should be available to genomeSIMLA during **execution time**. If the file can't be found, there will be a large amount of warnings rendered to STDOUT and none of the graphs will have any textual information on them- but execution will continue.

#### CSS\_FILENAME filename

```
CSS_FILENAME ../genomesimla.css
```

In order to make the reporting flexible, each report refers to a stylesheet which contains the necessary information about shading spacing and other information. An example stylesheet is provided with the application as well as each of the examples. Users are welcome to change this to suit their needs- and should be aware that editing the stylesheet does not require anything be done with genomeSIMLA. However, it is necessary that the stylesheet be found, as stated in the configuration file when the reports are read.

The example above indicates that there will be a file named genomesimla.css that resides in the directory above the one in which the report is read from. In other words:

If there is a report named: */home/torstees/genomesimla/data/test1.index.50.html*

And I used the setting from the example above, the following file **must** exist:  
*/home/torstees/genomesimla/genomesimla.css*.

If I copy the report (s) to a new filesystem, I should be sure to copy the style-sheet to the parent directory of the new location.

If the style-sheet isn't found, the report will just be harder to read.

#### LD\_REPORT\_BUFFER\_SIZE On/Off

```
LD_BUFFER_SIZE 50
```

This sets the number of SNPS around the block in the *detailed plot*. In the example above, 50 SNPS on either side of a block are drawn.

#### WRITE\_LD\_REPORT On/Off

```
WRITE_LD_REPORT On
```

Causes genomeSIMLA to produce a complete complete report of the pairwise LD values. Users should be aware that this file can be very large.

### **DRAW\_RSQUARED\_PLOTS On/Off**

### **DRAW\_DPRIME\_PLOTS On/Off**

*DRAW\_RSQUARED\_PLOTS Off*

By default, both RSquared and DPrime plots are drawn. If the user wants to save time and disk space, they can opt to one or both charts off.

### **MAX\_SNP\_DISTANCE integer**

*MAX\_SNP\_DISTANCE 500000*

This allows the user to determine far apart SNPs can be before genomeSIMLA decides to calculate LD values. Lowering this value from the default (500K) can speed up LD processing.

### **CLOSE\_POOLS\_BETWEEN\_DROPS On/Off**

*CLOSE\_POOLS\_BETWEEN\_DROPS On*

In general, it is assumed that genomeSIMLA will be used to produce very large populations (~1 million unique chromosomes) with genomes that approach 500K. In order to manage this on a single computer, genomeSIMLA must close pools down when they aren't currently in use. This frees up valuable memory- allowing us to do this without gigabytes of ram. However, it can slow down processing when only one or two small chromosomes in use.

### **FAST\_LD\_POOL\_SIZE integer**

*FAST\_LD\_POOL\_SIZE 3000*

Sets the number of chromosomes that are used in sampled plots. If this is larger than the population, the entire population is used.

### **FAST\_LD\_PLOT\_SIZE**

*FAST\_LD\_PLOT\_SIZE 10000*

Sets the max number of SNPs used in sampled LD production. If the size is greater than a given chromosome, the entire chromosome is used.

### **NO\_FASTLD**

*NO\_FASTLD*

Unlike most other options, this takes no parameters. When it is encountered, genomeSIMLA will skip the sampled LD calculations and instead render only complete LD plots. Please be aware that this will increase the length of time during a growth scan where you are unsure where the best generation to draw from lies. It should be used only when you either know ahead of time what your LD will look like, or you want to get statistics on the entire population at every drop.

# Locus Generation

genomeSIMLA allows for the creation of *chromosomes* in two different ways. The first involves the description of one or more block types and populating a chromosome randomly with one or more of these blocks. The other approach uses a locus description file containing positional SNP names and positional information.

## Block Based Locus Generation

Block based generation is the production of chromosomes using completely random draws. The user specifies one or more block configurations which will be applied randomly to create the loci on a given chromosome. There are 3 elements involved in this process:

**Block Definitions:** These describe 4 things:

- Min/Max number of snps that can be associated with the “block”
- Min/Max recombination fraction for the first SNP (how far away is that SNP from the previous SNP on the chromosome)
- Min/Max recombination fraction for each of the containing snps.
- Probability this block will be drawn

When a chromosome draws a block definition to be used to construct a set of loci, it will randomly draw the number of SNPs based on the block Min/Max value. Then, for each SNP, it will determine the distance between each SNP and its predecessor. All but the first SNP use the second set of Min/Max recombination values. The first SNP is drawn from the first set. This allows the user to space the block further out from the SNPs in front of it (or not).

**Default Block:** When a chromosome is deciding which block definition to use next, it uses the probabilities associated with the blocks. It is possible for the sum to be less than 1.0. The difference between the sum and 1.0 is the probability that the default block will be used. The default block is common to ALL chromosomes, and should be defined before any other blocks or chromosomes. With the exception of probability, the default block has the same parameters as regular blocks.

**Chromosome:** To create a block based chromosome, users will use the ADD\_CHROMOSOME command, indicating how many blocks to draw and possibly giving it a label. The user then applies blocks to the chromosome by using the ADD\_BLOCK command.

When the draws are made, the blocks associated with a given chromosome will be drawn based on their probability (including the possibility of using the default block, if necessary).

**DEFAULT\_BLOCK min max float float float float**

```
DEFAULT_BLOCK 5 10 0.01 0.015 0.00001 0.000025
```

The default block is used when the sum of a given chromosome’s blocks probabilities don’t sum up to 1.0. Otherwise, it is the same as a regular block.



The first two parameters specify the minimum and maximum number of SNPs will be created. The next two represent the range of distance this block falls from the previous SNP on the chromosome. The last two represent the range of distances of SNPs within the block itself.

DEFAULT\_BLOCK should be set prior to the definition of any chromosomes (and thus, any other blocks).

#### **ADD\_CHROMOSOME integer label**

```
ADD_CHROMOSOME 5 chromosome_1
```

Adds a new chromosome to the genome. The first parameter represents how many blocks to draw and the last (optional) parameter is the label that will be used in naming files and on the reports. The example above will create a chromosome with 5 blocks named "chromosome\_1".

#### **ADD\_BLOCK chr\_idx snp\_idx float float float float float**

```
ADD_BLOCK 5 10 0.0001 0.0002 0.000001 0.00001 0.5
```

Adds a block definition to the most recently defined chromosome (using ADD\_CHROMOSOME). The example above will create a block that ranges from 5 to 10 SNPs. 0.0001 and 0.0002 represent the chance of a cross-over event occurring between the previous SNP (if one exists) and the first in the block. This effectively describes how far away from that last SNP the block is.

The next two describe the chance of a recombination occurring between any two SNPs found inside the block itself.

The last parameter is the probability this block will be drawn.

#### **Other Block Related Settings**

##### **DEFAULT\_ALLELE\_FREQ float float**

```
DEFAULT_ALLELE_FREQ 0.1 0.5
```

This allows the user to define min/max allele frequencies to be used during the configuration of a new block based chromosome.

#### **File Based Chromosome Configuration**

File based chromosome files have all of the information necessary to simulate a chromosome. There are two reasons one would use such files:

1. To mimic one or more region from a real genome.
2. To precisely control a region's density as part of a research project.

The Ritchie Lab has made a set of these files available which represent a large portion of the Affymetrix 500K coverage. As we produce others, it is expected they will be made available for use as well. These files allow for SNPs to be distributed very similarly to real human assays- though, the actual LD patterns will depend largely on the generations the data-sets were extracted from as well as the random seed used.

#### **Locus File Format**

Line #1 is just a line used to describe which chromosome the file was derived from. This is ignored when the file is read

Line #2: Indicates the number of loci contained within the file. genomeSIMLA doesn't actually parse that number out- so again, this line is not used for reading.

Line #3: Column Headers. This is for the user's benefit...and is not used during reading

Line #4...N+3: Each line describes a single locus. The following represent the 6 columns that should be present (in the order listed). Each column must have a value for each line, and should be separated by whitespace (multiple spaces or tabs is fine).

Col #1: Label This is usually the RS Number. However it can be any label one wants to use. All SNPS must have unique labels.

Col #2: Freq Allele 1: Allele 1's allele frequency

Col #3: Freq Allele 2: Allele 2's allele frequency

Col #4: Recombination Fraction: Chance that an odd number of recombinations took place between this SNP and the previous SNP in the genome.

Col #5: Position: This is the physical position on the chromosome (relative to the beginning of that chromosome...NOT the genome) These values should be in base pairs.

Col #6: Description (optional): This is just a note that can be added. Currently, this isn't used anywhere.

The last line of the file should be an empty line (the last entry should contain a return character).

It should be noted that when genomeSIMLA sets up the loci, allele 1 is ALWAYS the minor allele, regardless of the locus' frequency in the file. This is only important if a user were to draw datasets from a pool at generation 0. Their interpretation of 'A' and 'a' could be different from the way genomeSIMLA. When drawing data-sets from generation 0, 'A' is ALWAYS the minor allele.

Also, allele frequencies are not exact, even in large populations. When one is setting up a disease model for generation 0, it is recommended to let genomeSIMLA create the pool, drop generation 0 (it defaults to this) and assign model loci based on allele frequencies found in the locus file generated during the initialization.

### Locus Miscellany

**ALLELE\_FREQUENCY chr\_idx snp\_idx float float**

```
ALLELE_FREQUENCY 1 5 0.25 0.75
```

This sets the frequency of allele 1 of Snp # 5 on chromosome 1 to 25% and the second allele to 75%.

# Population Control

Currently, there is a single population in genomeSIMLA, though each individual could have several different chromosomes. This population is grown using one of several growth rates. During a generational advancement, individuals are drawn (with replacement) from the current population, mated using Hardy-Weinburg mating and added to the new pool until it reaches it's target size.

Growth rates share many parameters. Below is a list of parameters that are used in each of the growth curves.

initial_population	integer	The population that is created at generation 0.
variation	float	This value is used to simulate imperfect growth curves. It represents the percentage of fluctuation around the curve's value at a given generation. The amount of fluctuation is actually +/- 1/2 the variation- so it is possible that the population at generation N+1 be smaller than at N.
growth_rate	float	This is the rate of growth. While it is applied differently for each model, the higher the growth-rate, the faster the growth.
carrying_capacity	integer	This is used in logistic style growths and specifies the ceiling of the growth curve. As the population approaches this value, it becomes less and less exponential in nature until it becomes static.
time_of_max_growth	integer	Used only in Richard's Logistic, this parameter effectively moves the exponential part of an S curve about on the X axis (in the direction of the generation specified).
polarity	float	Used only in Richard's Logistic, this parameter affects the "draw" of the curve toward the carrying capacity.

To set up a growth rate, the user should configure one of the following:

## **GROWTH\_RATE LINEAR initial\_population variation growth\_rate**

```
GROWTH_RATE LINEAR 30000 0.05 10.0
```

This is just a straight line that grows by growth\_rate each generation

## **GROWTH\_RATE EXPONENTIAL initial\_population variation growth\_rate**

```
GROWTH_RATE EXPONENTIAL 700 0.05 0.3
```

This is just a basic exponential growth based on the growth rate specified.

### **GROWTH\_RATE LOGISTIC initial\_population variation growth\_rate carrying\_capacity**

*GROWTH\_RATE LOGISTIC initial\_population variation growth\_rate carrying\_capacity*

This is considered to be one of the preferred models for describing growth rates. The carrying capacity represents the peak potential (which could be caused by various reasons). For our needs, it is the size of pool required for drawing data-sets.

### **GROWTH\_RATE RICHARDS initial\_population variation growth\_rate carrying\_capacity time\_of\_max\_growth polarity**

*GROWTH\_RATE LOGISTIC initial\_population variation growth\_rate carrying\_capacity time\_of\_max\_growth polarity*

Richard's logistic is just an enhanced logistic curve- with two parameters capable of determining when growth starts and just how steep the growth will be.

By pushing the time\_of\_max\_growth forward, the population hovers at initial\_population for some amount of time. This small population will produce rich LD patterns which tend to be carried forward in time once growth begins. However, this small population increases the risk of fixing alleles dramatically.

### **General Growth rate parameters**

#### **MAX\_POOL\_SIZE integer**

*MAX\_POOL\_SIZE 90000*

Sets the hard upper limit for population size. Every generation is compared against this value and can NEVER exceed it, regardless of variation nor growth curve details. It is very important for non-logistic growth rates, especially exponential, where growth could occur very fast and cause memory problems.

#### **MIN\_POOL\_SIZE integer**

*MIN\_POOL\_SIZE 1500*

Sets the semi-hard lower limit for population size. This is evaluated for all generations other than 0. So, it is possible to set the initial population to below MIN\_POOL\_SIZE and cause a hard spike in population at generation 1.

#### **TARGET\_POP\_SIZE integer**

*TARGET\_POP\_SIZE 100000*

If TARGET\_POP\_SIZE is greater than 0, genomeSIMLA will use the value as a hard limit for advancement- once it reaches the specified population size all advancement and population growth will cease.

# Dataset Generation

The entire purpose of genomeSIMLA is data. genomeSIMLA is capable of generating 2 types of data-sets: case/control and basic pedigrees. genomeSIMLA can produce any number of different data-sets and guarantees that, in a truly diverse population, no single individual will be used in any data-sets generated during a single run.

Both pedigree and case/control data-sets allow for the use of a label. This label is used as part of the filename, and allows the user to quickly recognize different data-sets. These labels can have any character except slashes and spaces.

Both types of data-sets can have the following types of error:

genotype_error	float	Exact Portion of SNPs which are not derived via cross-over. This error is applied evenly across SNPs
phenocopy	float	Percentage of the affected individuals in a given data-set whose affected status was determined not by the chosen model.
missing	float	Percentage of SNPs that will be missing.

## Case/Control

### **DATASET CC label affected unaffected genotype\_error phenocopy missing**

```
DATASET CC sample-01 500 500 0.05 0.1 0.15
```

This line creates datasets with 500 affected, 500 unaffected each with 5% genotype error and 15% missing data. Of the 500 affected individuals, 50 of them will not have been evaluated with the model.

## Pedigree Data

Pedigree data is slightly more complicated because you can specify multiple types of family structures to be added to your dataset. The affected/unaffected numbers simply describe the number of children in those categories.

### **DATASET PED label genotype\_error phenocopy missing**

```
DATASET PED family-01 0.05 0.1 0.15
```

This sets up a framework for datasets with 5% genotype error and 15% missing data. 10% of all affected children will not have been evaluated with the disease model.

This just sets up the data-set framework. Until you add family types to it, the data-sets will be empty.

### **DATASET FAMTYPE affected unaffected extra\_sibs number\_of\_families**

```
DATASET FAMTYPE 1 1 1 250
```

This sets up a type of family which will be added to the data-set. A given data-set can have as many different types of families as the user needs. The affected/unaffected counts represent the number of affected/unaffected children a given family MUST have. The number of extra sibs indicates that a random number from 0 to extra\_sibs will be added to the family. All children will be evaluated for status, however, by adding extra sibs, you can have larger families which vary by the number of affected siblings.

It is perfect acceptable to have 0 extra\_sibs or unaffected sibs. Affected sibs MUST be greater than or equal to 1.

```
DATASET FAMTYPE 1 0 0 150
```

This would add 150 trios to the data-set.

```
DATASET FAMTYPE 2 1 0 75
```

This would add 75 AAU families to the data-set.

```
DATASET FAMTYPE 1 0 3 50
```

This would add 50 families with between 1 and 4 children with at least 1 affected sib in them.

Parent's status is evaluated and written to the dataset, but is not considered for determining whether or not the family will be included into the dataset.

*Performance Note: It is important to note that pedigrees with more than 1 affected individual can be computationally difficult. Children are created by actually crossing over the parents just like is done during generational advancement. If the children don't meet the necessary family shape, all individuals are thrown away. For instance, for a disease with a prevalence of 0.1, it would take the production of almost 1,000,000 families before we found a family with 2 affected sibs and it gets worse as you add in more required affected sibs. Most data-sets can be generated in a few minutes, but be aware of the possibility of long delays for large numbers of affected sibs and rare disease models.*

## General Data-set Configuration Parameters

### DATASET\_COUNT integer

```
DATASET_COUNT 500
```

This indicates the number of files that will be created per data-set. In this example, all data-sets created by this configuration would result in 500 unique files.

### BINARY\_DATASETS Yes/No

```
BINARY_DATASETS Yes
```

This compresses data-sets dramatically, allowing whole genome size data-sets to occupy a minimal amount of disk-space. This format was developed in-house, and won't be supported by any products other than those produced at the lab here (and only now are we beginning to implement it in our own applications).

If you are interested in the format, we will make the format available on the wiki in the near future. In the meantime, feel free to contact us at [genomeSIMLA@chgr.mc.vanderbilt.edu](mailto:genomeSIMLA@chgr.mc.vanderbilt.edu).

This is currently not supported for pedigree datasets.

### USE\_STD\_PEDIGREE\_HEADER on/off

```
USE_STD_PEDIGREE_HEADER On
```

When on, all pedigree data-sets will have 10 column headers. When off, the header count will be 6 columns.

# Locus Searching

The main goal for genomesimla is the production of realistic data-sets. These might be very large, and choosing disease loci can be a daunting task when presented with over 200,000 possible loci. To make the task as easy as possible, genomeSIMLA can limit the loci presented and present them in a sorted fashion where the topmost SNP shown most closely matches the user's specifications.

The following commands are used to set up searches. Users can have as many searches as they like...even if they don't need them all for setting up their models.

A search describes 3 qualities: minor allele frequency ranges, types of blocks the SNP is contained within and location. Each of the ranges contains three pieces. Target Min and Max. Currently, LOCUS\_SELECTOR has 2 ranges: minor\_allele\_frequency and block\_size.

```
LOCUS_SELECTOR label float float float integer integer integer description  
LOCUS_SELECTOR rare_loci 0.2 0.15 0.23 4 2 10 The following loci are  
moderately rare and appear in a block
```

This creates a new search called *rare\_loci* which will only contain SNPs whose minor allele frequency is between 0.15 and 0.23 and are found in blocks of up to 10 SNPs large.

The SNPs will be ranked so that those that are closest to a minor allele frequency of 0.2 and in blocks with 4 SNPs will be ranked first. Notice that the min/max values are not evenly distributed around the target. The score is ranked on the relative distance from the target for that particular arm. So, a SNP with a minor allele frequency with just a bit larger MAF than 0.2 would score very similarly to one that had a MAF of just under 0.23.

The block size and MAF weights are different. MAF is more currently weighted higher than block size. However, each block a SNP is found in will add more to its final score- meaning it will rise higher in the report.

The description of a LOCUS\_SELECTOR gets used in the locus report. Be descriptive as is necessary. However, there can be no newline characters in it. Spaces are allowed, though.

```
ADD_REGION label snp_start snp_stop  
ADD_REGION rare_loci rs-321412 rs-543231
```

This adds a region to the selector, *rare\_loci*. This region is bounded by the two SNPs rs-321412 and rs-543231. Both SNPs must be found and exist on the same chromosome.

By default, all searches are performed over the entire genome. However, if a user wishes to restrict the region to search, they can do so by adding regions (once you add a single region, it will only search the regions that have been added). To add a whole chromosome, simply add the first and last SNP geographically.

```
MAX_LOCI_PER_CHROM_REPORTED Integer  
MAX_LOCI_PER_CHROM_REPORTED 50
```

Instructs genomeSIMLA to report at most N loci per chromosome for each sector described. Loci are ranked according to how well they fit the criterion, which could possibly be rather extensive. By setting this value to a reasonable number, the locus report can be kept at a manageable size. Setting the value to -1 will catch all possible matching loci.

# Disease Modeling

genomeSIMLA comes with 3 options for modeling affection status: User generated Penetrance tables, simPEN (purely epistatic models) and SIMLA (main effect + interactions). Each method requires it's own configuration details.

## Penetrance Table Disease Models

Users can use predefined penetrance tables to assign status to models. The only requirement is that the user specify the allele frequencies associated with each possible allele associated with each model locus. This is to help ensure that the appropriate meaning of a given cell is being applied. genomeSIMLA will not proceed to use a model if the actual allele frequencies differ too much from those specified in the configuration.

Users indicate to genomeSIMLA that a penetrance based model is to be used using a line similar to one of the two lines below.

**DEFINE\_MODEL PENTABLE INDEX pen-file chrom-id snp-id [chrom-id snp-id, ...]**

```
DEFINE_MODEL PENTABLE INDEX disease.pen 1 5
```

**DEFINE\_MODEL PENTABLE LABEL pen-file snp-label [snp-label, ...]**

```
DEFINE_MODEL PENTABLE LABEL disease.pen RL5
```

Both lines do the same thing. The first tells genomeSIMLA to load the penetrance table in disease.pen and use Locus 5 on chromosome 1 as the single disease locus. The contents of the specified penetrance table must match the number of model loci specified on the configuration line. Otherwise, genomeSIMLA will generate an error (or worse, become confused and generate misleading data-sets.)

The second example simply uses labels to specify which loci are to be associated with the disease model. Block based chromosomes (see BLOCK\_DEFINITION) are labeled RLN where N is a number between 1 and however many loci there are associated with all chromosomes being simulated. Otherwise, the labels are based on information found inside the locus files that were used to populate the simulation (most likely, this will be an RS number). There can be no duplicately labeled SNPs.

### Penetrance File Configuration

A small number of parameters make up the configuration details of a penetrance file. All but the threshold must be fully specified even if the value is 0.0.

**FREQ\_THRESHOLD Float**

Specifies maximum allowed variation from the allele frequencies that will be tolerated before execution is halted.

**FREQ [AaBbCcDd....etc] float**

```
FREQ A 0.2
```

```
FREQ a 0.8
```

Using letter notation for specifying penetrance cells, this command allows the user to tell genomeSIMLA what the intended frequency for a given allele should be. The user **MUST** specify all alleles that are expected to be involved in the given model.



## PENTABLE

This just indicates to genomeSIMLA That the various penetrances are about to follow.

### model-identification penetrance

```
AABB 0.171  
AABb 0.155
```

Each possible combination must be present, regardless if it's value is anything other than 0.0.

Penetrance tables should be written as a separate file from the main configuration.

## Purely Epistatic Models with simPEN

simPEN is a method for using a Genetic Algorithm (GA) to evolve purely epistatic models. With few exceptions, the configuration details are considered to be beyond the scope of this document, however, a few details will be covered, such as those that specify target odds ratios and heritability.

Users indicate to genomeSIMLA that a penetrance based model is to be used using a line similar to one of the two lines below.

### DEFINE\_MODEL SIMPEN INDEX simpen-cfg chrom-id snp-id [chrom-id snp-id, ...]

```
DEFINE_MODEL SIMPEN INDEX disease.simpen 1 5
```

### DEFINE\_MODEL PENTABLE LABEL simpen-cfg snp-label [snp-label, ...]

```
DEFINE_MODEL PENTABLE simpen-cfg disease.pen RL5
```

Both lines do the same thing with the exception of how they specify which loci are to be associated with the model.

The configuration file for simpen must be a separate file. When deciding which weights are most appropriate, users should keep in mind that the values themselves can differ drastically, and a weight of 1 differs in effectiveness for a value whose target is 0.1 than that of a value whose target is 0.000001. The values used in the following examples were determined to be reasonable starting points for obtaining good results from the simPEN module.

### simPEN File Configuration

A small number of parameters make up the configuration details of a simpen configuration file.

#### HERIT float

```
HERIT 0.01
```

Specifies that the target heritability will be 0.01.

#### HERITWEIGHT float

```
HERITWEIGHT 10
```

Adjusting the weight determines how important the attribute is.

#### ODDSRATIO float

```
ODDSRATIO 1.25
```

**ODDSWEIGHT float**

*ODDSWEIGHT 1*

**MARGVAR float**

*MARGVAR 0.000001*

Set the target Marginal variance. This determines how pure of an epistatic model you want. The higher the values, the more likely there will be main effects.

**MARGWEIGHT float**

*MARGWEIGHT 100*

**PENTARGET float**

*PENTARGET 0.15*

Specifies the target prevalence of the *disease*.

The following parameters are associated with the GA portion of simPEN. For more information about how to use these parameters, please see the simPEN user's manual (available at the genomeSIMLA [website](#)). References to pool sizes, generations, populations and mutation below are completely unrelated to the simulation being performed by the forward time simulation.

**GEN 15000**

Number of generations to be tried before ending.

**POPSIZE 1000**

The search population.

**DEMES 100**

Multiple pools of penetrance values

**MUTATE 0.01**

Frequency of mutation

**CROSS 0.6**

Rate of cross over (1 per "genome")

**SUBMODELS ON**

Turning this on will possibly catch a pattern where a smaller model contained within a larger model does exist with enough strength as to represent a potential problem.

**UPDATE 100**

Specifies how many generations between progress is reported

The following two parameters are legacy and have no effect on the production of valid models. However, the error checking currently requires them to be present. Just use these values to satisfy the error checking code for now.

**LOCI 2**

**FREQ 0.2 0.8**

## Main Effects and Interactions with SIMLA

SIMLA is a simulation program that allows the researcher to specify varying levels of both linkage and linkage disequilibrium among and between markers and disease loci. SIMLA was specifically designed for the simultaneous study of linkage and association methods in extended pedigrees, but the penetrance specification algorithm can also be used to simulate samples of unrelated individuals (e.g., cases and controls).

Users indicate to genomeSIMLA that a SIMLA based model is to be used using a line similar to one of the two lines below.

**DEFINE\_MODEL SIMLA INDEX simla-cfg float int chrom-id snp-id MIN|MAJ float float [...]**

```
DEFINE_MODEL SIMLA INDEX interactions.simla 0.05 2 1 5 MIN 0.26 0.0
```

**DEFINE\_MODEL SIMLA LABEL simla-cfg float int snp-label MIN|MAJ float float [...]**

```
DEFINE_MODEL SIMLA LABEL interactions.simla 0.05 2 RL5 MIN 0.26 0.0
```

Both lines do the same thing with the exception of how they specify which loci are to be associated with the model.

The filename specified by simla-cfg represents the list of interactions (see below). If no interactions are required, you may use the keyword, NO\_INTERACTIONS- otherwise, the file must exist. The next number (0.05 in the example above) is the target prevalence. The last parameter before the loci is the maximum interaction size. This is just maximum number of loci that will be interacting with one another.

The MIN/MAJ value determines whether the disease is associated with the minor or major allele. The next parameter specifies the beta value associated with that locus. Finally, the last value required for each locus is the type. A 0.0 represents a recessive trait. The locus becomes more dominant as it approaches 1.0. Each locus to be considered must have each of these parameters (locus specification, MIN|MAJ, beta and type).

### SIMLA configuration file

The simla-cfg is just a file that specifies the beta values associated with each of the interactions desired. For each interaction, specify them in the following way:

```
1x2x3 0.26
```

This tells genomeSIMLA that the 1st, 2nd and 3rd locus (in the order they are encountered on the DEFINE\_MODEL line) interact with a beta value of 0.26. Users can add as many or as few interactions as they wish.