# Data simulation (including genomics) QMSim software

Zulma Vitezica[†]

[†] INRA, UMR 1289, Castanet-Tolosan 31326 France

zulma.vitezica@ensat.fr

**ALIMENTATION**
**AGRICULTURE**
**ENVIRONNEMENT**

# QMSim: why to use it ?

✓ It was design to simulate large-scale genotyping data in multiple and complex livestock pedigrees

✓ A wide variety of genome architectures from infinitesimal model to single-locus model

✓ It is a user-friendly tool for simulating data

✓ Computationally efficient in termes of both time and memory

# QMSim[†]: where to find it ?

The code is written in C++ language

Executable files are freely available for Windows and Linux platforms at:

http://www.aps.uoguelph.ca/~msargol/qmsim/

[†]Sargolzaei & Schenkel (2009), Bioinformatics 25:680-681.

# How the simulation is carried out ?

In 2 steps:

- ✓ *First step:* A **historical population** is simulated
  - – in order to create initial LD and
  - – to establish mutation-drift equilibrium
  - – expansion and contraction of the population

- ✓ *Second step:* One or multiple **recent population structures** are generated

# Parameter file

- ✓ It must be in ASCII format

- ✓ It consists of five main sections

- ✓ The order of commands within each section is not important

- ✓ All commands end with a semicolon

- ✓ No semicolon → error message and program exits.

```
/*******************************
**     Global parameters      **
*******************************/
title = "Example 1 - 10k SNP panel
...;

/*******************************
**    Historical population   **
*******************************/
begin_hp;
      ....;
end_hp;

/*******************************
**        Populations         **
*******************************/
begin_pop = "p1";
      ....;
end_pop;

/*******************************
**           Genome           **
*******************************/
begin_genome;
      ....;
end_genome;

/*******************************
**       Output options       **
*******************************/
begin_output;
      ....;
end_output;
```

# 1. Global parameters section

```
/*********************************
 **      Global parameters      **
 *********************************/
title = "Example 1 - 10k SNP panel";
seed = "seed.txt";
```

An arbitrary title

The random number generator (RNG*) requires a seed file.

If it is not specified → RNG will be seeded from the system clock

For each run the initial seed numbers will be backed up in output folder
→ This allows to repeat the run !

**Parameter file:** ex01.prm
**Output folder:** r_ex01/

```
.----------------------------------------.
| Example 1 - 10k SNP panel |
`----------------------------------------'
```

Output

Initial seed is backed up in [r_ex01/seed].
parameter file is backed up in [r_ex01/ex01.prm].

* Mersenne Twister algorithm (Matsumoto & Nishimura, 1998)

# 1. Global parameters section

```
/*********************************
 **      Global parameters      **
 *********************************/
title = "Example 1 - 10k SNP panel";
nrep  = 1;        //Number of replicates
h2    = 0.2;      //Heritability
qtlh2 = 0.2;      //QTL heritability
phvar = 1.0;      //Phenotypic variance
```

Range: 0 - 10,000

Overall heritability (Polygenic + QTL

QTL effect is simulated

```
title = "Example 8
nrep  = 1;
h2    = 0.2;
qtlh2 = 0.0;
phvar = 1.0;
```

Only polygenic effect is simulated

```
title = "Example 11
nrep  = 1;
h2    = 0.2;
qtlh2 = 0.05;
phvar = 1.0;
```

Both, polygenic and QTL effects are simulated

# 1. Global parameters section

```
/*******************************
 **      Global parameters     **
 ******************************/
title = "Example 1 - 10k SNP panel";
nrep  = 1;        //Number of replicates
h2    = 0.2;      //Heritability
qtlh2 = 0.2;      //QTL heritability
phvar = 1.0;      //Phenotypic variance
no_male_rec;      // No record for males
```

A sex limited trait like milk yield

When males do not have records, but selection or culling are based on

**EBVs** → Ok

**Phenotypes** → Males will be randomly selected or culled

INP ENSAT

INRA

# 2. Historical population section

```
/*******************************
**    Historical population    **
*******************************
begin_hp;
    hg_size = 420 [0]            //
              420 [200];
    nmlhg   = 20;                //
end_hp;
```

➡ To create initial LD

➡ Evolutionary foces: mutation and drift (no selection, no migration)

➡ Random mating: union of gametes randomly sampled from the male and female gametic pools

➡ Discrete generations

➡ Only a single historical population

# 2. Historical population section

Historical generation sizes

```
/*******************************
**    Historical population    **
*******************************
begin_hp;
    hg_size = 420 [0]
             420 [200];
    nmlhg    = 20;                //
end_hp;
```

Constant size of 420

hg_size = **v1** [**v2**]

**v1** the historical generation **size**
  *Range:* 2 – 100,000

**v2** the historical generation **number**
  *Range:* 0 – 150,000

# 2. Historical population section

Historical **bottleneck** or **expansion** can be simulated

```
/********************************
 **    Historical population    **
 *******************************/
begin_hp;
   hg_size = 2000 [0]
            200 [1000];
   nmlhg   = 40;
end_hp;
```

```
/**********************************
 **    Historical population    **
 *********************************/
begin_hp;
   hg_size = 100 [0]
            100 [950]
            3000 [1000];
   nmlhg   = 200;
end_hp;
```

Gradual decrease in size from 2000 to 200

Expansion in the last historical generation from 100 to 3000

**LD** in livestock extends over longer distances than in humans

# 2. Historical population section

```
/********************************
 **    Historical population    **
 ********************************/
begin_hp;
    hg_size = 2000 [0]
              200 [1000];
    nmfhg    = 40;
end_hp;
```

Number of males

Default : equal number of males and females

nm**f**hg → **first** historical generation

Sex ratio will be constant across historical generations.
It can be changed in the last generation

```
/********************************
 **    Historical population    **
 ********************************/
begin_hp;
    hg_size = 2000 [0]
              200 [1000];
    nmlhg    = 40;
end_hp;
```

nm**l**hg → **last** historical generation

# 3. Population section

```
/*******************************
 **         Populations          **
 ********************************/
begin_pop = "p1";
        ....;
end_pop;

begin_pop = "p2";
        ....;
end_pop;
```

One or multiple recent populations

For the *first defined recent population* (i.e. p1), **founders** must come from the last historical population

For *subsequent populations* (i.e. p2), **founders** can be chosen from one or more (up to 10) previously defined populations (i.e. p1)

Multiple recent populations can be analyzed
    separately (one pedigree for each population) or
    jointly (by creating one pedigree for all populations) for inbreeding and EBV

# 3. Population section

**Choosing founders for a population**

```
/*******************************
 **        Populations         **
 *******************************/
begin_pop = "line1";
    begin_founder;
        male   [n = 20,  pop = "hp", select = tbv /h];
        female [n = 400, pop = "hp", select = tbv /h];
    end founder;
```

Parameters for choosing founders

Number of male/female to be selected

It indicates from which population the base animals must be selected

Type of selectionn

**select:** rnd (default), phen, tbv and ebv
**/l :** to select low values
**/h :** to select high values

**hp:** historical population (last historical generation)

```
/*********************************
 **        Populations          **
 *********************************/
begin_pop = "line1";
    begin_founder;
        male    [n = 20,  pop = "hp", select = tbv /h];
        female [n = 400, pop = "hp", select = tbv /h];
    end_founder;
    ng   = 20;          //Number of generations
end_pop;


begin_pop = "line2";
    begin_founder;
        male    [n = 20,  pop = "hp", select = tbv /l];
        female [n = 400, pop = "hp", select = tbv /l];
    end_founder;
    ng   = 20;          //Number of generations
end_pop;


//Cross between line1 and line 2 to generate F2
begin_pop = "cross";
    begin_founder;
        male    [n = 20, pop = "line1", gen = 20];
        female [n = 400, pop = "line2", gen = 20];
    end_founder;
    ng   = 2;          //Number of generations
```

Choosing founders for a population for F2 design

Crossing between populations/lines is allowed

```
/********************************
 **        Populations        **
 ********************************/
begin_pop = "line1";
    begin_founder;
        male    [n = 20,  pop = "hp", select = tbv /h];
        female [n = 400, pop = "hp", select = tbv /h];
    end_founder;
    ng  = 20;          //Number of generations
end_pop;

begin_pop = "line2";
    begin_founder;
        male    [n = 20,  pop = "hp", select = tbv /l];
        female [n = 400, pop = "hp", select = tbv /l];
    end_founder;
    ng  = 20;          //Number of generations
end_pop;

//2 males and 10 females from line 2 immigrate to line 1
begin_pop = "line1_c";
    begin_founder;
        male    [n =  8, pop = "line1", gen = 10];
        male    [n =  2, pop = "line2", gen = 10]; //2 male immigrants
        female [n = 90, pop = "line1", gen = 10];
        female [n = 10, pop = "line2", gen = 10]; //10 female immigrants
    end_founder;
    ng  = 5;           //Number of generations
```

Migration can be simulated

# 3. Population section

```
/********************************
 **          Populations          **
 ********************************/
begin_pop = "p1";
    begin_founder;
        male    [n =
        female [n = 2500
    end_founder;
    ls   = 1 2 [0.2];          //Litter size
    pmp  = 0.5;                //Proportion of male progeny
    ng   = 10;                 //Number of generations
    md   = p_assort/ebv;       //Mating design
    sr   = 0.4;                //Replacement ratio for sires
    dr   = 0.2;                //Replacement ratio for dams
    sd   = ebv /h;             //Selection design
    cd   = phen/l;             //Culling design
    ebv_est = blup;
```

**ls:** number of progeny per dam

**ls:** Probability of the litter sizes

# 3. Population section

```
/*******************************
 **          Populations          **
 *******************************
begin_pop = "p1";
    begin_founder;
        male    [n =    50,
        female [n = 2500, pop       ];
    end_founder;
    ls   = 1 2 [0.2];        //Litter size
    pmp = 0.5;               //Proportion of male progeny
    ng   = 10;               //Number of generations
    md   = p_assort/ebv;     //Mating design
    sr   = 0.4;              //Replacement ratio for sires
    dr   = 0.2;              //Replacement ratio for dams
    sd   = ebv /h;           //Selection design
    cd   = phen/l;           //Culling design
    ebv_est = blup;
```

**pmp:** range 0-1, default is equal to 0.5

**pmp:** 0.5 /fix_litter
Sex ratio will be fixed within litters (progeny of a dam)

# 3. Population section

> **rnd** (default)**, rnd_ug** (a dam can mate with more than one sire in each generation), **p_assort** (similarity), **n_assort** (dissimilarity), **minf** and **maxf** (inbreeding is minimized in the next generation)

```
*******************
                ulations        **
******************/
                1";
                er;
           n =    50, pop = "hp"];
           n = 2500, pop = "hp"];
                er;
    ls     2 [0.2];        //Litter size
    pmp =   ;              //Proportion of male progeny
    ng  = 10;              //Number of generations
    md  = p_assort/ebv;    //Mating design
    sr  = 0.4;             //Replacement ratio for sires
    dr  = 0.2;             //Replacement ratio for dams
    sd  = ebv /h;          //   tion design
    cd  = phen/l;
    ebv_est = blup;
```

> Assortative mating base on **phen**, **ebv** or **tbv**

# 3. Population section

```
/*******************************
**          Populations          **
*******************************/
begin_pop = "p1";
      _founder;
      le    [n =    50, pop = "hp"];
      male [n = 2500, pop = "hp"];
      _founder;
   = 1 2 [0.2];          //Litter size
p  = 0.5;                //Proportion of male progeny
ng = 10;                 //Number of generations
md = p_assort/ebv;       //Mating design
sr = 0.4;                //Replacement ratio for sires
dr = 0.2;                //Replacement ratio for dams
sd = ebv /h;             //Selection design
cd = phen/l;             //Culling design
ebv_est = blup;
```

**sr :** 40% of sires will be replaced in all generations

**sr :** 0.4 [1] 0.5 [5]
40% of sires will be culled for generation 1 to 5, and 50% from generation 5 to last generation

**sr :** 1, discrete generations (default)

# 3. Population section

```
/*********************************
 **         Populations          **
 *********************************/
begin_pop = "p1";
    begin_founder;
        male   [n =    50, pop = "hp"];
        female [n = 2500, pop = "hp"];
    end_founder;
    ls  = 1 2 [0.2];          //Litter size
    pmp = 0.5;                //Proportion of male progeny
    ng  = 10;                 //Number of generations
    md  = p_assort/ebv;       //Mating design
    sr  = 0.4;                //Replacement ratio for sires
    dr  = 0.2;                //Replacement ratio for dams
    sd  = ebv /h;             //Selection design
    cd  = phen/l;             //Culling design
    ebv_est = blup;
```
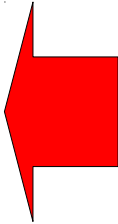
**rnd, phen, tbv ebv** and **age** (only for culling)

Breeding value estimation method

**/l** or **/h** to select low or high values

```
/********************************
 **          Populations          **
 ********************************/
begin_pop = "p1";
    begin_founder;
        male    [n =  20, pop = "hp"];
        female [n = 400, pop = "hp"];
    end_founder;
    ls  = 2;
    pmp = 0.5 /fix;
    ng  = 10;
    begin_popoutput;
        data;
        stat;
        genotype /snp_code /gen 8 9 10;
    end_popoutput;
end_pop;
```

**Population specific parameters for saving outputs**

p1_**mrk**_007.txt

p1_**qtl**_007.txt

**data:** save individual's data except their genopype
*(File name:* 'population name'_**data**_'replicate number'.txt

**stat:** save brief statistic on simulated data

**genotype:** save genotype data

# 4. Genome section

**Example – 10k SNP panel**

```
/*************************
**          Genome
*************************
begin_genome;
    begin_chr = 10;
        chrlen = 100;      //Chromosome length
        nmloci = 1000;     //Number of markers
        mpos   = rnd;      //Marker positions
        nma    = all 2;    //Number of marker alleles
        maf    = eql;      //Marker allele frequencies
        nqloci = 25;       //Number of QTL
        qpos   = rnd;      //QTL positions
        nqa    = all 2;    //Number of QTL alleles
        qaf    = eql;          //QTL allele frequencies
        qae    = rndg 0.4;     //QTL allele effects
    end_chr;
    mmutr      = 2.5e-5 /recurrent; //Marker mutation rate
    qmutr      = 2.5e-5;            //QTL mutation rate
    r_mpos_g;      // Randomize marker positions across genome
    r_qpos_g;      // Randomize QTL positions across genome
end_genome;
```

Number of chromosomes: 10
**chrlen :** range 1-5,000 cM

Samples from uniform distribution in each replicate

All marker loci will have 2 alleles

In the first historical generation, then drift and mutation

# 4. Genome section

**Example – 10k SNP panel**

```
/*****************************
**          Genome          **
                         ****
        chrlen = 10        //Chromosome    gth
        nmloci = 10        //Number    arkers
        mpos   = rn        //Marke   positions
        nma    = al  2;    //Nu   er of marker al
        maf    = eq ;      //  arker allele freque
        nqloci = 25;       //Number of QTL
        qpos   = rnd;      //QTL positions
        nqa    = all 2;    //Number of QTL alleles
        qaf    = eql;        //QTL allele frequencies
        qae    = rndg 0.4;   //QTL allele effects
end_chr;
        mmutr  = 2.5e-5  /                        ate
        qmutr  = 2.5e-5;
        r_mpos_g;    // Ra                          home
        r_qpos_g;    // Rand                           ne
end_genome;
```

nqloci: range 1-50,000 on the chromosome

Samples from uniform distribution in each replicate

Nb of QTL alleles in the first historical generation (all: same number)

Equal allele frequencies in the first historical generation

It will be sampled from gamma distribution with shape 0.4

# Example – 10k SNP panel

**More genome information**

```
/*********************************
 **           Genome           **
 *********************************/

begin_genome;
    begin_chr = 10;
                        ...omosome length
                        ...er of markers
                        ...er positions
    ...       ...        Number of ...
    maf      = eql;      //Marker al...
    nqloci   = 25;       //Number of ...
    qpos     = rnd;      //QTL posit...
    nqa      = all 2;    //Number of QTL ...
    qaf      = eql;      //QTL allele f...les
    qae      = rndg 0.4; //QTL allele eff...
    end_chr;
    mmutr    = 2.5e-5 /recurrent; //Marker mutation rate
    qmutr    = 2.5e-5;            //QTL mutation rate
    r_mpos_g;    // Randomize marker positions across genome
    r_qpos_g;    // Randomize QTL positions across genome
end_genome;
```

In recurrent mutation, no new allele is generated.
Default: infinite-allele model

SNP recurrent mutations are generally very rare and no evidence that mutation contributes to erosion of LD between SNP ( Ardlie et al., 2002)

In historical population

Other possibilities :
  Missing marker/QTL genotypes
  Genotyping errors can be simulated (marker/QTL)

# 5. Output section

```
/*******************************
 **        Output options       **
  *******************************/
begin_output;
    linkage_map;
    hp_stat;
end_output;
```

Marker and QTL linkage map

Save brief statistics on historical population

```
/*******************************
  **        Output options       **
   *******************************/
begin_output;
    linkage_map;
    allele_effect;
end_output;
```

Save allele effects

# QMSim outputs

```
/********************************
**          Populations          **
********************************/
begin_pop = "p1";
    begin_founder;
        male    [n =  20, pop = "hp"];
        female  [n = 400, pop = "hp"];
    end_founder;
    ls   = 2;
    pmp  = 0.5  /fix;
    ng   = 10;
    begin_popoutput;
        data;
        stat;
        genotype /snp_code /gen 8 9 10;
    end_popoutput;
end_pop;
```

p1_**data**_001.txt

```
.----------.
| Example 1 |
'----------'
Progeny Sire    Dam     G     Sex NMPrg NFPrg F          Homo      Phen      Res       Polygene   QTL
1       0       0       0     M   33    27    0.000000 0.696797 +1.323314 +0.331291 -0.000000 +0.992023
2       0       0       0     M   21    19    0.000000 0.695996 +0.933861 +1.323803 -0.000000 -0.389942
3       0       0       0     M   9     11    0.000000 0.673574 +0.903691 -0.106867 -0.000000 +1.010557
4       0       0       0     M   20    20    0.000000 0.685385 +0.502346 +0.068033 -0.000000 +0.434313
5       0       0       0     M   18    22    0.000000 0.696096 -0.038755 +0.870122 +0.000000 -0.908877
6       0       0       0     M   11    9     0.000000 0.692092 +2.246078 +1.202401 +0.000000 +1.043677
7       0       0       0     M   34    26    0.000000 0.704304 +1.312932 +1.393522 +0.000000 -0.080591
8       0       0       0     M   22    18    0.000000 0.692793 +1.375544 +1.060612 +0.000000 +0.314932
```

```
  .............
|  Example 1  |
 `............'
------------------ Inbreeding -------------------
                      Inbred              All
Gen.        No.    Mean      SD      Mean      SD
0             0   0.0000  0.0000   0.0000  0.0000
1             0   0.0000  0.0000   0.0000  0.0000


----------------- Homozygosity ------------------
Gen.                Mean                SD
0                0.68254159         0.01207245
1                0.68200626         0.01103250


------------------- Phenotype -------------------
Gen.                Mean                SD
0                0.08440969         1.01093563
1                0.04504056         1.02152016


--------------------- QTL -----------------------
Gen.                Mean                SD
0                0.04889285         0.56092140
1               -0.00533798         0.55392545
```

```
begin_popoutput;
    data;
    stat;
    genotype /gen 8 9 10;
end_popoutput;
end_pop;
```

| Gen. | Progeny | Male% | Male | Selected | Female | Selected | Sire | Culled | Dam | Culled |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 420 | 0.047619 | 20 | 0 | 400 | 0 | 0 | 0 | 0 | 0 |
| 1 | 400 | 0.500000 | 200 | 8 | 200 | 80 | 20 | 8 | 400 | 80 |
| 2 | 400 | 0.500000 | 200 | 8 | 200 | 80 | 20 | 8 | 400 | 80 |
| 3 | 400 | 0.500000 | 200 | 8 | 200 | 80 | 20 | 8 | 400 | 80 |
| 4 | 400 | 0.500000 | 200 | 8 | 200 | 80 | 20 | 8 | 400 | 80 |
| 5 | 400 | 0.500000 | 200 | 0 | 200 | 0 | 20 | 0 | 400 | 0 |
| Overall | 2420 | 0.421488 | 1020 | 32 | 1400 | 320 | 100 | 32 | 2000 | 320 |

p1_**stat**_001.txt

p1_**mrk**_001.txt

```
                begin_popoutput;
                     data;
                     stat;
                     genotype /snp_code /gen 4 5;
                end_popoutput;
           end_pop;
```

.-----------.
|  Example 1  |
`-----------'

Progeny Sire    Dam        Genotypes (0 = a1,a1; 2 = a2,a2; 3 = a1,a2; 4 = a2,a1; 5 = missing; The first
allele is paternal and the second allele is maternal) ...
1621     1408    365        0 2 2 2 2 2 0 2 2 2 2 4 4 3 0 3 3 3 0 0 0 2 0 3 2 2 4 0 2 0 0 4 3 2 3 3 3 3
0 2 2 0 0 3 4 3 4 4 0 2 0 4 0 0 3 2 2 4 2 0 4 2 3 2 3 0 0 4 0 4 2 2 3 0 4 2 0 3 0 0 4 3 3 0 2 0 2 2 0
2 2 2 0 0 0 0 0 2 0 0 2 0 0 2 0 4 2 2 3 4 0 4 2 2 0 0 0 0 4 4 2 2 3 4 0 0 2 3 2 4 0 3 4 0 2 2 4 0 0 2
0 2 0 3 0 0 3 2 2 0 0 0 4 0 4 2 2 2 2 0 3 2 2 2 0 4 2 2 2 0 2 2 0 2 0 0 2 2 4 3 2 0 2 0 0 2 0 2 2 0 2
2 2 2 2 1623    1408    1096      3 4 4 4 2 2 4 3 4 2 4 4 2 2 0 0 0 0 3 0 0 0 2 0 3 2 2 4 0 2 0 0 4 3 2
3 3 3 3 0 2 2 0 0 3 4 3 4 4 0 2 0 4 0 0 3 2 2 4 2 0 4 2 3 2 3 0 0 4 0 4 2 2 3 0 4 2 0 3 0 0 4 0 3 3 4
0 4 2 0 4 4 2 3 3 3 3 0 2 0 0 2 0 0 2 0 2 2 2 0 2 0 4 2 2 0 0 0 0 2 2 2 0 2 0 0 2 0 4 4 0 3 4 0 2 2
4 0 3 2 3 2 0 0 0 0 0 2 4 3 0 3 2 3 4 4 4 2 2 3 0 2 2 2 2 0 2 2 2 2 0 2 2 0 2 2 0 0 2 2 2 0 2 0 2 0 0 2 0
4 2 3 4 2 2 2 2 0 0 2 2 0 0 0 2 2 2 0 2 0 3 4 0 4 4 2 0 2 2 2 2 0 0 2 0 2 2 2 3 4 2 3 4 2 4 4 4 3 4 2
3 2 0 0 2 0 0 0 0
```

.-----------.
|  Example 1  |
`-----------'

Progeny Sire    Dam       Genotypes (paternal allele, maternal allele) ...
1621     1139    285       2 2 1 1 2 1 1 1 1 2 2 1 2 1 1 2 1 1 1 1 2 1 2 1 1 1 1 2 1 1 2 2 1 1 1 1 1 1 2
2 1 1 1 1 2 2 2 2 2 2 2 2 1 1 2 2 2 1 1 1 1 2 1 2 1 1 1 2 1 1 1 2 1 1 2 1 1 2 1 2 1 1 2 1 2 1 2 1 1 2
2 2 2 1 2 1 2 2 1 2 1 1 2 2 2 2 2 2 2 2 1 1 1 2 1 1 2 2 2 2 1 2 1 1 2 2 2 1 1 2 1 2 1 2 2 2 2 2 2 2 1 1
1 2 2 2 2 2 1 1 1 1 1 2 1 2 2 2 1 2 1 2 2 1 1 1 2 1 2 2 1 1 2 1 1 2 2 1 2 1 2 1 2 1 1 2 2 1 1 1 1 1 2
2 2 2 1 2 1 1 1 2 1 1 1 1 1 2 2 2 2 2 1 2 1 1 1 1 2 1 2 2 2 2 1 1 1 2 1 2 1 2 2 2 1 2 1 2 1 2 2 2
1 2 1 1 1 2 1 2 1 2 2 1 2 2 1 1 1 1 1 1 1 1 1 2 1 1 2 1 2 1 2 2 2 1 1 1 2 1 2 1 2 1 1 1 1 1 1 2 1 2 1 1 2 1 2
1 1 2 1 1 1 1 1 1 1 1 1 2 1 2 1 2 1 2 1 2 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 1 2 2 1 1 2 2 1 1 1 1 1 2 1 1
2 2 1 2 1 2 1 2 1 2 1 2 2 1 1 2 2 2 1 1 2 1 2 1 1 2 2 2 2 2 2 2 1 1 2 2 2 2 2 2 2 2 2 2 2 2 2 1 1
2 2 2 2 2 2 2 2
```

```
/*********************************
 **       Output options        **
 *********************************/
begin_output;
    linkage_map;
    hp_stat;
end_output;
```

Marker and QTL linkage map

```
            .-----------.
           |  Example 1  |
            `-----------'

     ----------- QTL linkage map ---
    ID          Chr      Position
    ------------------------------------
    Q1           1         8.88876
    Q2           1        13.35024
    Q3           1        17.76099
    Q4           1        22.12918
    Q5           1        29.68482
    Q6           1        37.76335
    Q7           1        43.84122
    Q8           1        46.93041
    Q9           1        47.16755
    Q10          1        48.56634
```

```
/********************************
 **       Output options
 ********************************
```

```
begin_output;
    linkage_map;
    hp_stat;
end_output;


 .----------.
|  Example 1  |
 `----------'

--------------- Statistics for the last historical generation ----------------
Scaled QTL mean (var)      :      -0.00000 (0.3)
QTL scale factor           :       1.64027
Polygenic mean (var)       :       0.00000 (0)


                                  Overall       Chr1       Chr2       Chr3       Chr4       Chr5
Chromosome length (male)   :     3000.000    100.000    100.000    100.000    100.000    100.000
Chromosome length (female) :     3000.000    100.000    100.000    100.000    100.000    100.000

No. loci                   :        10740        358        358        358        358        358
No. markers                :         9990        333        333        333        333        333
No. QTL                    :          750         25         25         25         25         25
No. non-seg. loci          :         1241         47         30         45         39         43
No. non-seg. markers       :           89          2          3          4          2          5
No. non-seg. QTL           :         1152         45         27         41         37         38


                                     male      female
Density in total           :      3.58000    3.58000 per 1 cM
Density for markers        :      3.33000    3.33000 per 1 cM
Density for QTL            :      0.25000    0.25000 per 1 cM
```

```
/*******************************
 **       Output options       **
 *******************************/
begin_output;
    linkage_map;
    allele_effect;      ←  Save allele effects
end_output;
```

```
      .-----------.
      |  Example 1  |
      `-----------'

ID        Chr      Allele:Effect ...
------------------------------------------------
Q1        1        1: 0.066403    2:-0.001068
Q2        1        1:-0.050405    2: 0.031267
Q3        1        1:-0.006917    2: 0.009631
Q4        1        1:-0.000543    2: 0.000171
Q5        1        1:-0.001498    2: 0.004858
Q6        1        1: 0.001299    2:-0.000535
Q7        1        2: 0.000000
Q8        1        1:-0.004849    2: 0.003374
Q9        1        1:-0.014103    2: 0.018606
Q10       1        1: 0.048198    2:-0.006161
Q11       1        1: 0.000189    2:-0.001423
```

# Conclusion

**+**

To create LD

Population expansion or bottleneck

QTL + polygenic

Dense marker map

Sex limited traits

Multiple recent populations / lines

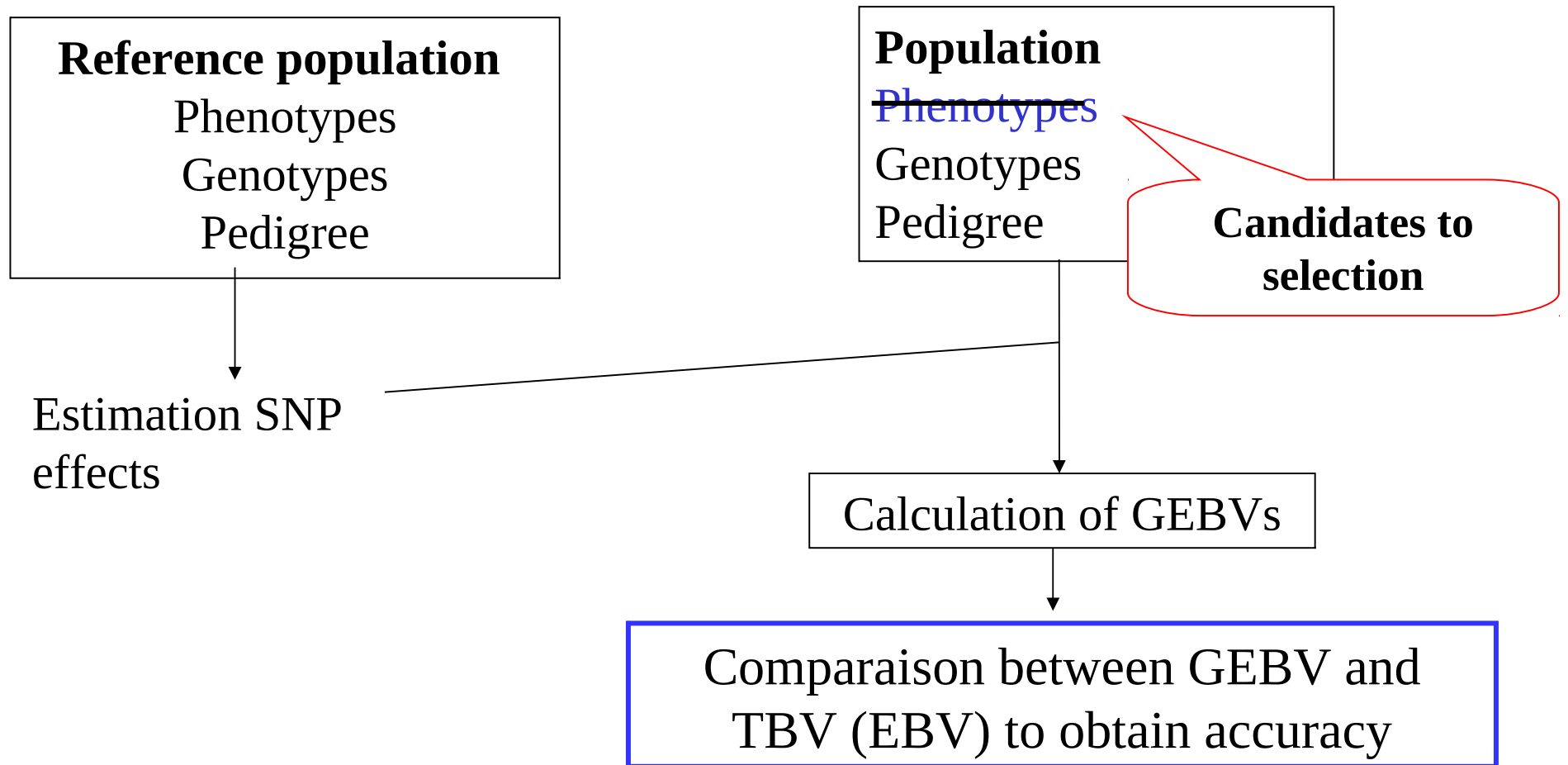Crossing between populations / lines

**QMSim**

**-**

A single historical population

No fixed effects

Only additive effects

# Genomic selection : validation

# Example of simulation

Generation -1000 to -5 — Random mating (N=100)

Generation -5 to -1 — Expanded to N=3000

*Pedigree recording and genotyping start*

Generation 1 to 9 ← Training data — 200 ♂ x 2,000 ♀ / generation

Generation 10 ← Validation data: candidates to selection