

FASTSIMCOAL2 INPUT FILES

Vitor Sousa

vmsousa@fc.ul.pt

Cesky Krumlov 2020

https://github.com/vsousa/EG_cE3c

Examples of observed SFS

1PopExpInst20Mb_DAFpop0.obs

```
1 observations
d0_0    d0_1    d0_2    d0_3    d0_4    d0_5    d0_6    d0_7    d0_8    d0_9    d0_10
19973842 24630    810     173     145     111     88      84      61      56      0
```

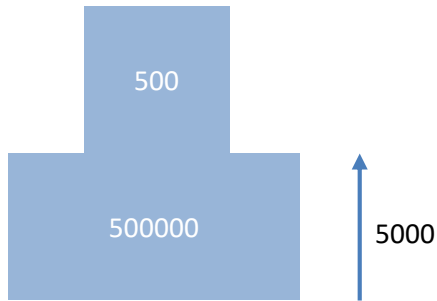
2PopDivMigr20Mb_jointDAFpop1_0.obs

```
1 observations
      d0_0    d0_1    d0_2    d0_3    d0_4    d0_5
d1_0 19985747 8350   1628   360    62     8
d1_1  9660     0      0      0      0
d1_2  4790     0      0      0      0
d1_3  3280     0      0      0      0
d1_4  2490     0      0      0      0
d1_5  1760     13     18     13     19     0
```

2PopDiv20Mb_jointDAFpop1_0.obs

```
1 observations
      d0_0    d0_1    d0_2    d0_3    d0_4    d0_5
d1_0 19985547 8211   1415   316    55    10
d1_1  1266     101    37     16     5     1
d1_2  61142    20     8      2      0
d1_3  48631    12     5      0      0
d1_4  47915     9      2      3      1
d1_5  1189     46     22     19     18     0
```

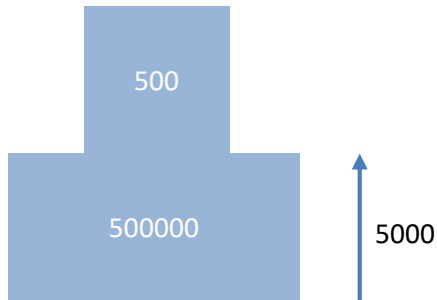
Parameter estimation settings files



1PopExplnst20Mb

- Additional files necessary to estimate parameters:
- Template file (TPL) defining the model
 - Estimation file (EST) with search range for parameters

Parameter estimation settings files



1PopExpInst20Mb

- Additional files necessary to estimate parameters:
- **Template file (TPL) defining the model**
 - Estimation file (EST) with search range for parameters

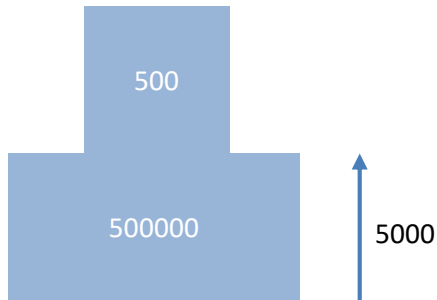
Tags for parameter we want to estimate:
\$NPOP\$, \$TEXP\$, \$RESIZE\$

Template file (filename.tpl)

```
1PopExpInst20Mb/1PopExpInst20Mb.tpl
```

```
//Parameters for the coalescence simulation program : fsimcoal2.exe
1 samples to simulate :
//Population effective sizes (number of genes)
$NPOP$
//Samples sizes and samples age
10
//Growth rates: negative growth implies population expansion
0
//Number of migration matrices : 0 implies no migration between demes
0
//historical event: time, source, sink, migrants, new deme size, new growth rate,
migration matrix index
1 historical event
$TEXP$ 0 0 0 $RESIZE$ 0 0
//Number of independent loci [chromosome]
1 0
//Per chromosome: Number of contiguous linkage Block: a block is a set of contiguous loci
1
//per Block:data type, number of loci, per generation recombination and mutation rates
and optional parameters
FREQ 1 0 2.5e-8 OUTEXP
```

Parameter estimation settings files



1PopExpInst20Mb

- Additional files necessary to estimate parameters:
- Template file (TPL) defining the model
 - **Estimation file (EST) with search range for parameters**

Tags for parameter we want to estimate:
\$NPOP\$, \$TEXP\$, \$RESIZE\$

Estimation file (filename.est)

1PopExpInst20Mb/1PopExpInst20Mb.est

```
// Search ranges and rules file
// *****

[PARAMETERS]
//#isInt? #name #search #min #max
//all Ns are in number of haploid individuals
1 $NPOP$ logunif 1000 1e7 output
1 $NANC$ logunif 10 1e5 output
1 $TEXP$ unif 10 1e5 output

[RULES]

[COMPLEX PARAMETERS]
0 $RESIZE$ = NANC/NPOP hide
```

INPUT files for fastsimcoal2:

Defining an evolutionary model with TPL file

2PopDivMig.tpl

```
//Parameters for the coalescence simulation program : fsimcoal2.exe
2 samples to simulate :
//Population effective sizes (number of genes)
$NPOP1$
$NPOP2$
//Samples sizes and samples age
5
5
//Growth rates: negative growth implies population expansion
0
0
//Number of migration matrices : 0 implies no migration between demes
2
//Migration matrix 0
0 0
$MIG10$ 0
//Migration matrix 1: No migration
0 0
0 0
//historical event: time, source, sink, migrants, new deme size, new growth rate, migration matrix
index
2 historical event
$TMIG_STOP$ 0 0 0 1 0 1
$TDIV_POP01$ 1 0 1 $RESIZE$ 0 1
//Number of independent loci [chromosome]
1 0
//Per chromosome: Number of contiguous linkage Block: a block is a set of contiguous loci
1
//per Block:data type, number of loci, per generation recomb. and mut. rates and optional parameters
FREQ 1 0 2.5e-8 OUTEXP
```

INPUT files for fastsimcoal2: Defining an evolutionary model with TPL file

Number of samples
to simulate

2PopDivMig.tpl

```
//Parameters for the coalescence simulation program : fsimcoal2.exe  
2 samples to simulate :  
//Population effective sizes (number of genes)
```

\$NPOP1\$

\$NPOP2\$

```
//Samples sizes and samples age
```

5

5

```
//Growth rates: negative growth implies population expansion
```

0

0

```
//Number of migration matrices : 0 implies no migration between demes
```

2

```
//Migration matrix 0
```

0 0

\$MIG10\$ 0

```
//Migration matrix 1: No migration
```

0 0

0 0

```
//historical event: time, source, sink, migrants, new deme size, new growth rate, migration matrix  
index
```

```
2 historical event
```

\$TMIG_STOP\$ 0 0 0 1 0 1

\$TDIV_POP01\$ 1 0 1 **\$RESIZES\$** 0 1

```
//Number of independent loci [chromosome]
```

1 0

```
//Per chromosome: Number of contiguous linkage Block: a block is a set of contiguous loci
```

1

```
//per Block:data type, number of loci, per generation recomb. and mut. rates and optional parameters
```

FREQ 1 0 2.5e-8 OUTEXP

Deme sizes (2N)

Sample sizes

Growth rates

Migration
matrices

Historical events

Keep these as default
(not used for SFS)

Definition of genetic data type to simulate.

For SFS inference use:

FREQ 1 0 fixedMutationRate OUTFREQ

NOTE: for the SFS you cannot jointly infer the effective sizes and mutation rates! You need to **give a fixed mutation rate** if you have the number of monomorphic sites.

Otherwise, with "-0" option the mutation rate is ignored.

FREQ indicates you will use the SFS

OUTFREQ means the expected SFS will be output

TPL files

These files are very important! Check carefully all the definitions. Errors in the TPL file are difficult to detect and imply the model specification is incorrect! This means that all inferences will be wrong, and also that all parameter estimates will be incorrect!

Defining population sizes and sample sizes

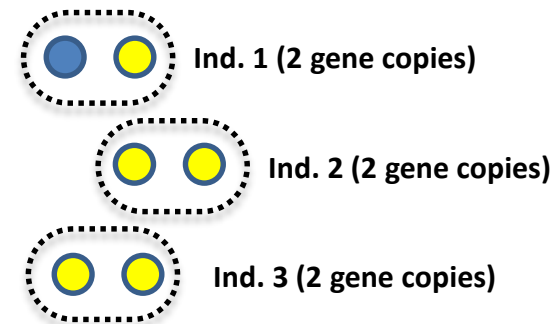
2PopDivMigr10Loci.par

```
//Parameters for the coalescence simulation program : fsimcoal2.exe
2 samples to simulate :
//Population effective sizes (number of genes)
$NPOP1$
$NPOP2$
//Samples sizes and samples age
6
6
//Growth rates: negative growth implies population expansion
0
0
```

Parameter tags

Population effective sizes are given in number of gene copies. For a diploid species with $N=500$ individuals, this corresponds to a $2N=1000$ gene copies, as each individual carries two gene copies at any given site.

The sample size is also given in gene copies. The value of 6 means that we sampled 3 diploid individuals.



TPL files

MIGRATION

```
//Number of migration matrices : 0 implies no migration between demes
1
//migration matrix
0.000 $MIG_01$
$MIG_10$ 0.000
```

Parameter tags

The migration matrix can be asymmetric, and in the case the entry m_{ij} list the **migration rates backward in time** from population in row i to population in column j . The above-mentioned matrix states that, for each generation (backward in time), any gene from population 0 has probability MIG_01 to be sent to population 1, and that a gene from population 1 has a probability MIG_10 to move to population 0.

If no migration matrix is defined, no migration is assumed between populations.

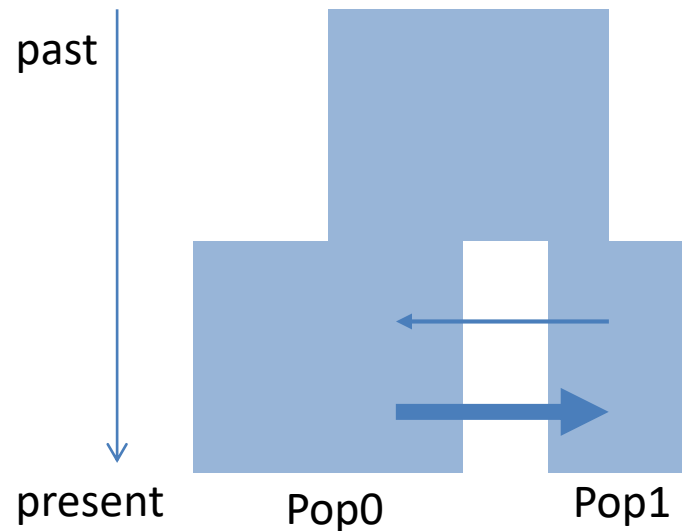
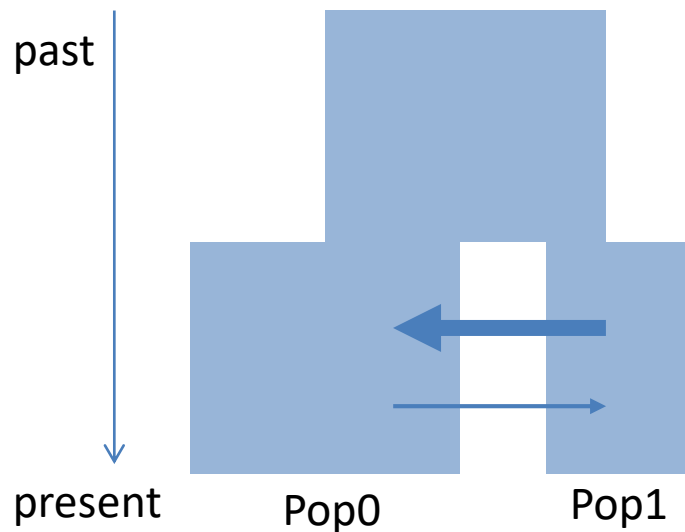
```
1PopStationary10Loci.par
```

```
//Number of migration matrices : 0 implies no migration between demes
0
```

A note on looking backward in time

Assuming that we look **forward in time** and that the size of the arrows are proportional to the migration rate, to what model does the following migration matrix corresponds to?

```
//Number of migration matrices : 0 implies no migration between demes  
1  
//migration matrix  
0.000 0.005  
0.001 0.000
```



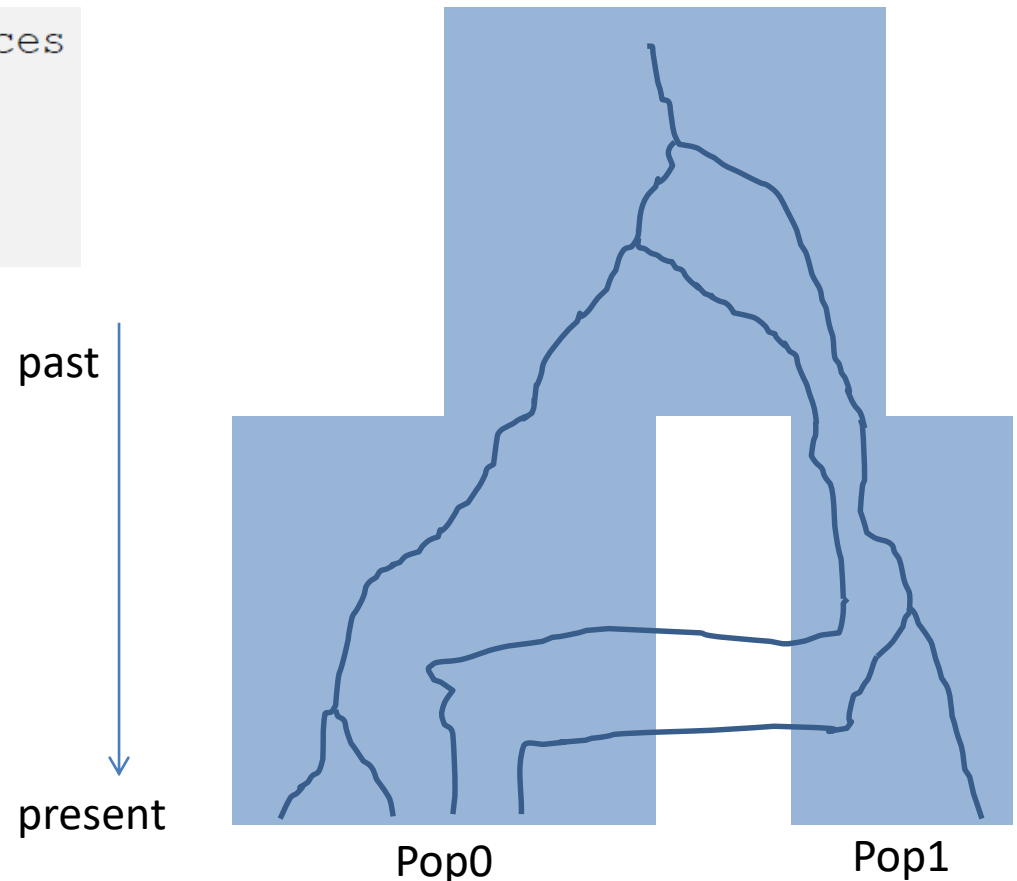
A note on looking backward in time

Assuming that we look **forward in time** and that the size of the arrows are proportional to the migration rate, to what model does the following migration matrix corresponds to?

```
//Number of migration matrices  
1  
//migration matrix  
0.000 0.005  
0.001 0.000
```

This means that there are more lineages migrating ("jumping") from pop0 to pop1 backward in time.

Thus, in Pop0 there are many individuals whose ancestors were migrants from Pop1 into Pop0.

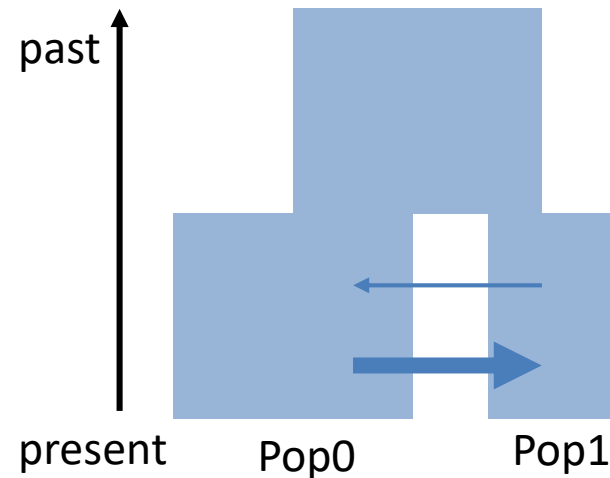
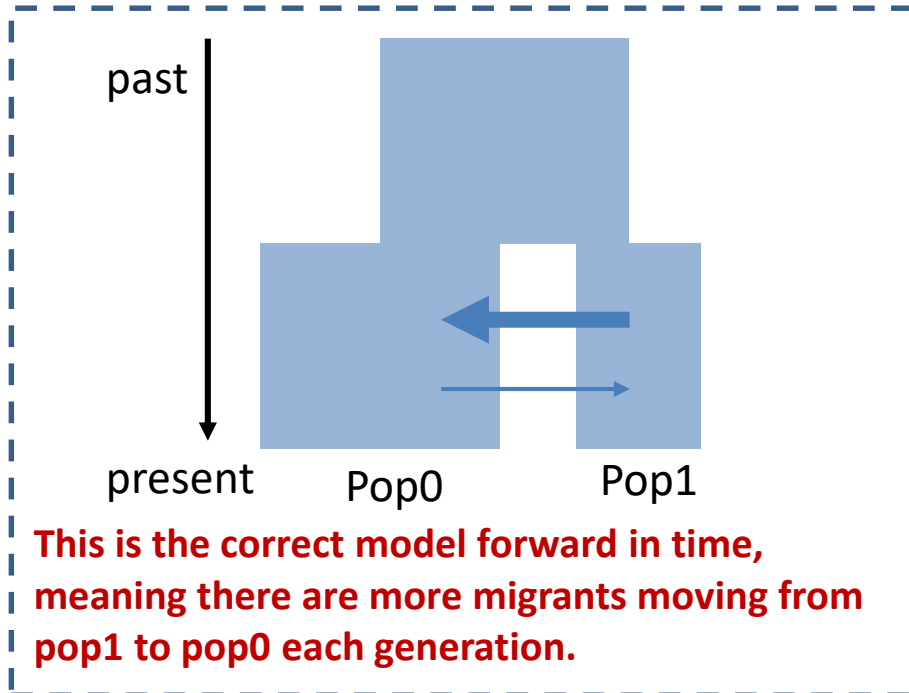


A note on looking backward in time

Assuming that we look forward in time and that the size of the arrows are proportional to the migration rate, to what model does the following migration matrix corresponds to?

```
//Number of migration matrices : 0 implies no migration between demes  
1  
//migration matrix  
0.000 0.005  
0.001 0.000
```

Note that in the PAR and TPL files everything is backward in time!!



Historical events in fastsimcoal2

Historical events can be used to:

- Change the size of a given population
- Change the growth rate of a given population
- Change the migration matrix to be used between populations
- Move a fraction of the genes of a given population to another population. This amounts to implementing a (stochastic) admixture or introgression event.
- Move all genes from a population to another population. This amounts to fusing two populations into one looking backward in time.
- One or more of these events at the same time

Defining the historical events is crucial to have a correct model!

Historical events (backward in time)

Each historical event is coded with a line with the following arguments

time, **source**, **sink**, **migrants**, **new deme size**, **new growth rate**, migration matrix index

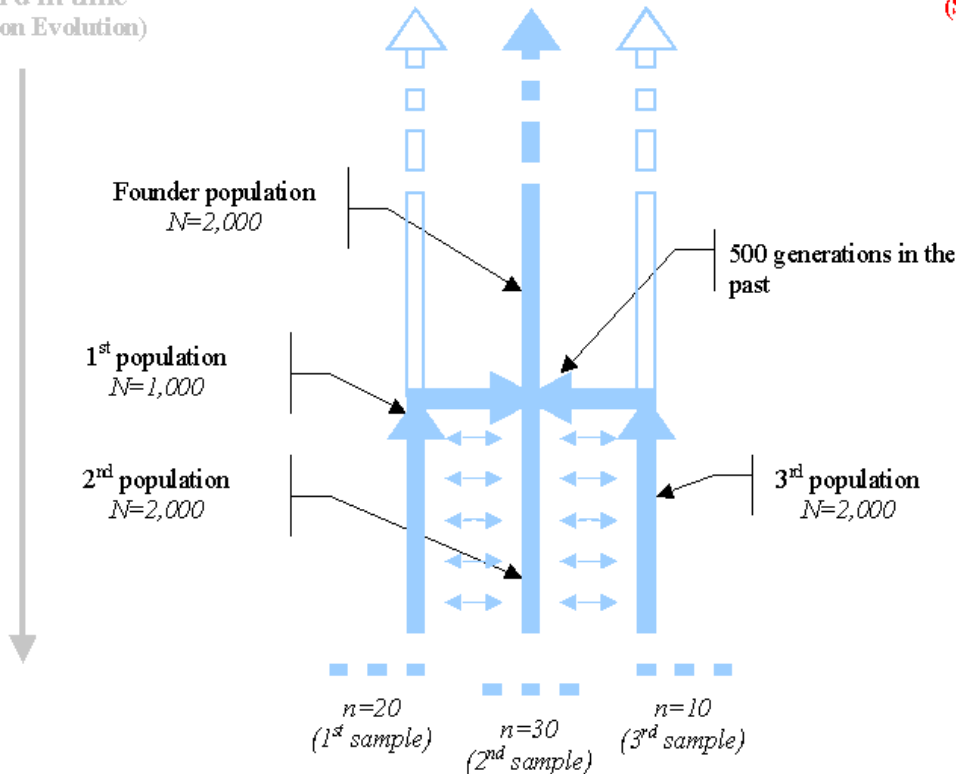
500 0 1 1 1 0 1

500 2 1 1 1 0 1

500 generations ago, 100% (migrants=1.0) of lineages in pop0 (source =0) migrated to pop1 (sink=1). The size of the sink (pop1) remained the same (new deme size=1.0, i.e. N2=2000). The new growth rate is zero. The migration rate that is active after the event is given in the migration matrix 1.

Forward in time
(Population Evolution)

Backward in time
(Simcoal2)



Historical events (backward in time)

Each historical event is coded with a line with the following arguments

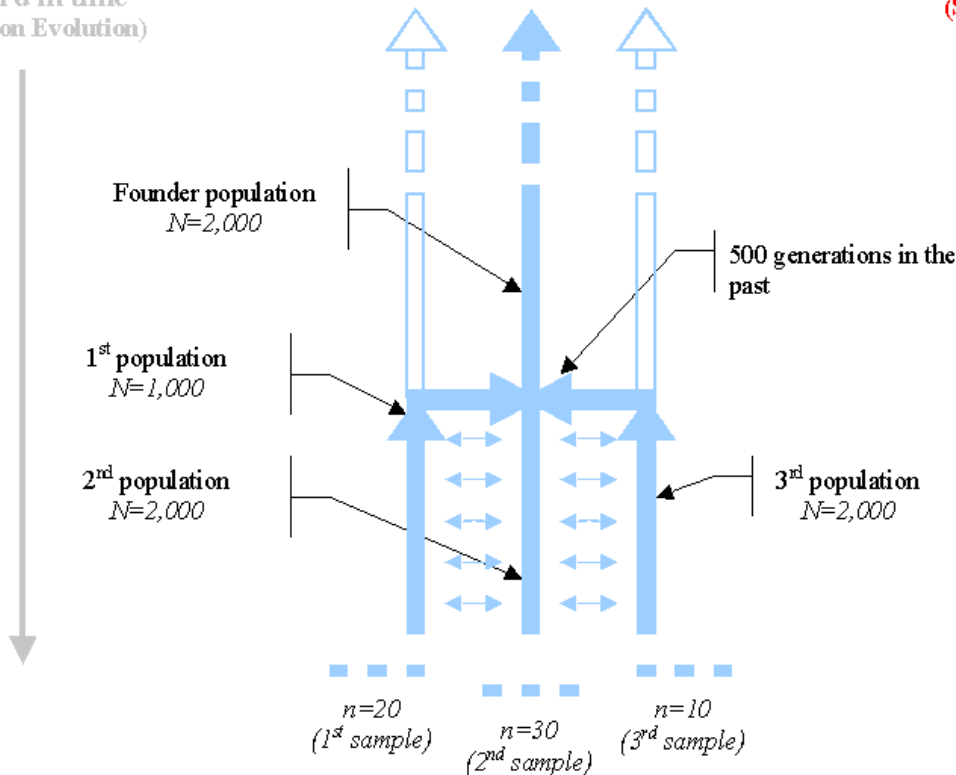
time, **source**, **sink**, **migrants**, **new deme size**, **new growth rate**, migration matrix index

500 0 1 1 1 0 1

500 2 1 1 1 0 1

Forward in time
(Population Evolution)

Backward in time
(Simcoal2)



500 generations ago, 100% of lineages (**migrants=1.0**) in **pop2 (source =2)** migrated to **pop1 (sink=1)**. The size of the sink (pop1) remained the same (**new deme size=1.0**, i.e. $N_2=2000$). The new growth rate is zero. The migration rate that is active after the event is given in the migration matrix 1.

Historical events in fastsimcoal2

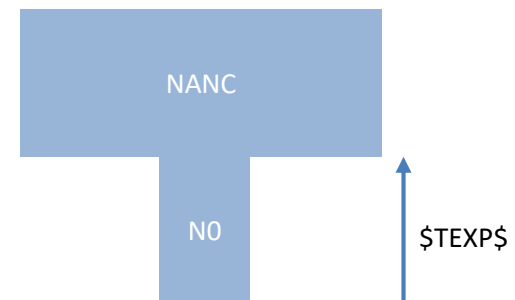
Change the size of a given population

1PopContrInst10Loci.par

```
//Parameters for the coalescence simulation program : fsmcoal2.exe
1 samples to simulate :
//Population effective sizes (number of genes)
1000
//Samples sizes and samples age
10
//Growth rates: negative growth implies population expansion
0
//Number of migration matrices : 0 implies no migration between demes
0
//historical event: time, source, sink, migrants, new deme size, new growth rate, migration matrix index
1 historical event
$TEXP$ 0 0 0 $RESIZE$ 0 0
```

- **\$TEXP\$** generations ago, 0% (migrants=0) of lineages in pop0 (source) migrated to pop1 (sink). This means that 100% of lineages remained in pop0.
- The sink population (pop0) has a size **\$RESIZE\$** times larger after the event (**\$RESIZE\$=\$NANC\$/\$NO\$**). Given NO diploids at time zero, it implies that $NANC=NO*RESIZE$ diploids.
- The migration matrix valid after the event is the migration rate 0. Since it is not defined it implies no migration.

Recent instantaneous
demographic contraction



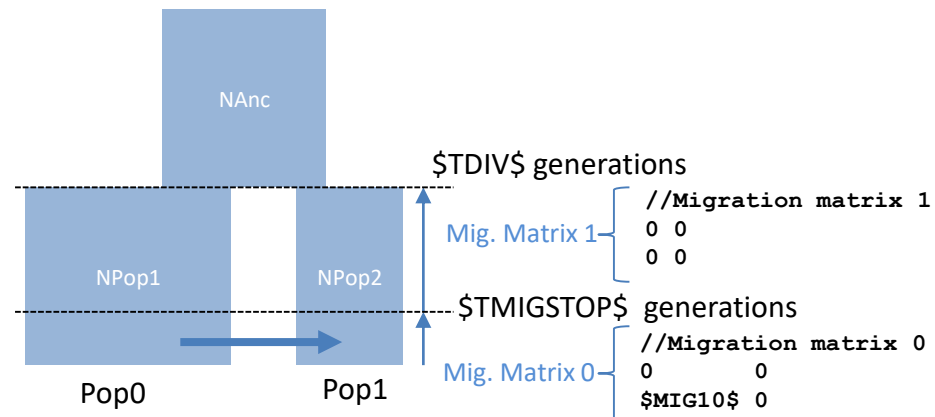
1PopContrInst10loci.par

Historical events in fastsimcoal2

Models with changes in the migration matrix to be used between populations

2PopDivMigr10Loci.par

```
//Number of migration matrices : 0 implies no migration between demes
2
//Migration matrix 0
0 0
$MIG10$ 0
//Migration matrix 1: No migration
0 0
0 0
//historical event: time, source, sink, migrants, new deme size, new growth rate, migration matrix
index
2 historical event
$TMIGSTOP$ 0 0 0 1 0 1
$TDIV$ 1 0 1 $RESIZE$ 0 1
```



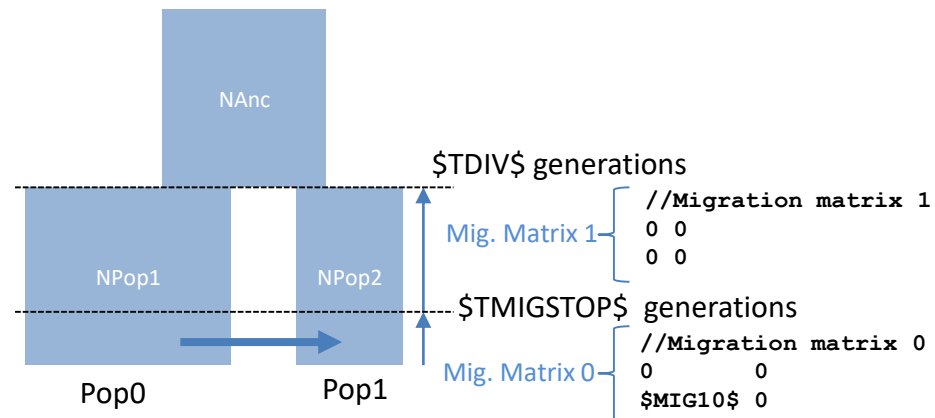
Historical events in fastsimcoal2

Migration matrix can change through time

```

2PopDivMigr10Loci.par
//Number of migration matrices : 0 implies no migration between demes
2
//Migration matrix 0
0 0
$MIG10$ 0
//Migration matrix 1: No migration
0 0
0 0
//historical event: time, source, sink, migrants, new deme size, new growth rate, migration matrix
index
2 historical event
$TMIGSTOP$ 0 0 0 1 0 1
$TDIVS$ 1 0 1 $RESIZES$ 0 1
    
```

- At generation \$TMIGSTOP\$ in the past, 0% (migrants=0) of lineages migrated from pop0 (source=0) to pop1 (sink=0).
- After the historical event, the deme size of the sink population (pop1) remained the same (new deme size=1).
- After the historical event the growth rate was set to zero.
- After the historical event the migration rate matrix was set to matrix 1, i.e. no migration between populations.



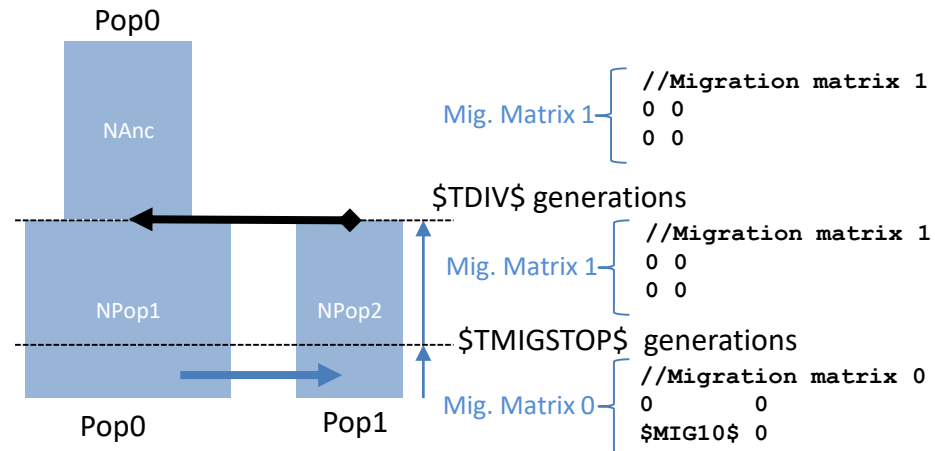
Historical events in fastsimcoal2

Population split (merge populations going backwards in time)

```

2PopDivMigr10Loci.par
//Number of migration matrices : 0 implies no migration between demes
2
//Migration matrix 0
0 0
$MIG10$ 0
//Migration matrix 1: No migration
0 0
0 0
//historical event: time, source, sink, migrants, new deme size, new growth rate, migration matrix
index
2 historical event
$STMIGSTOP$ 0 0 0 1 0 1
$STDIV$ 1 0 1 $RESIZES$ 0 1
    
```

- At generation \$STDIV\$ in the past, 100% (migrants=1) of lineages migrated from pop1 (source=1) to pop0 (sink=0).
- After the population split, the deme size of the sink population (pop0) is \$NANC\$, and hence \$RESIZES\$=\$NANC\$/\$NPOPO\$).
- After the historical event the growth rate of the sink population pop0 is zero.
- After the historical event the migration rate matrix was set to matrix 1, i.e. no migration between populations.



Estimation file (.est)

Estimation file ("NoMigSan_Maya.est")

```
NoMigSan_Maya.est
// Search ranges and rules file
// *****

[PARAMETERS]
//#isInt? #name #search.#min #max
//all Ns are in number of haploid individuals
1 $NPOP1$ unif 10 1e5 output
1 $NPOP2$ unif 10 1e5 output
1 $NANC$ unif 10 1e5 output
1 $NBOTP1$ unif 1 1e3 output bounded

1 $TDIV$ unif 100 1e4 output
0 $RELTBOT$ unif 1e-5 1 hide bounded

[RULES]

[COMPLEX PARAMETERS]

0 $RES_BOT_START$ = $NBOTP1$/$NPOP1$ hide
0 $RES_BOT_END$ = $NPOP1$/$NBOTP1$ hide

1 $TBOT_START$ = $TDIV$ * $RELTBOT$ output
1 $TBOT_END$ = $TBOT_START$ + 10 hide

0 $RESIZE0$ = $NANC$/$NPOP1$ hide
```

Each line of [PARAMETERS] section must contain the following:

#isInt? 0 for continuous, 1 for integers

#name Parameter tag name

#search "unif" for uniform scale
"logunif" for log10 scale

#min minimum search range (lower bound)

#max maximum search range. If the keyword bounded is not used, then if likelihood is higher near maximum value, fastsimcoal2 will keep increasing the maximum value. The **bounded** keyword prevents this.

Complex parameters depend on the values of other parameters. Only one operation per line can be done. Thus, you cannot have something with many operations in a single line:

$\$BLA\$ = (\$BL\$ * \$A\$) + (\$BLA\$ / \$LA\$) - \text{WRONG!}$

Estimation file (.est)

Estimation file ("NoMigSan_Maya.est")

NoMigSan_Maya.est

```
// Search ranges and rules file
// *****

[PARAMETERS]
//#isInt? #name #search.#min #max
//all Ns are in number of haploid individuals
1 $NPOP1$ unif 10 1e5 output
1 $NPOP2$ unif 10 1e5 output
1 $NANC$ unif 10 1e5 output
1 $NBOTP1$ unif 1 1e3 output bounded

1 $TDIV$ unif 100 1e4 output
0 $RELTBOT$ unif 1e-5 1 hide bounded

[RULES]

[COMPLEX PARAMETERS]

0 $RES_BOT_START$ = $NBOTP1$/$NPOP1$ hide
0 $RES_BOT_END$ = $NPOP1$/$NBOTP1$ hide

1 $TBOT_START$ = $TDIV$ * $RELTBOT$ output
1 $TBOT_END$ = $TBOT_START$ + 10 hide

0 $RESIZE0$ = $NANC$/$NPOP1$ hide
```

Note that complex parameters can be used to define the order of events.

By using a $\$RELTBOT\%$ between $1e-5$ and 1 , and then specifying that

$$\$TBOT_START\$ = \$TDIV\$ * \$RELTBOT\%$$

we define that the $TBOT_START$ is always more recent than the time of divergence.

If this is not well specified you can get errors, because events need to happen in a specific order.

Estimation file (.est)

Estimation file ("NoMigSan_Maya.est")

NoMigSan_Maya.est

```
// Search ranges and rules file
// *****

[PARAMETERS]
//#isInt? #name #search.#min #max
//all Ns are in number of haploid individuals
1 $NPOP1$ unif 10 1e5 output
1 $NPOP2$ unif 10 1e5 output
1 $NANC$ unif 10 1e5 output
1 $NBOTP1$ unif 1 1e3 output bounded

1 $TBOT_END$ unif 100 1e4 output
0 $TDIV_TBOT_INT$ unif 10 1e3 hide

[RULES]

[COMPLEX PARAMETERS]

0 $RES_BOT_START$ = $NBOTP1$/$NPOP1$ hide
0 $RES_BOT_END$ = $NPOP1$/$NBOTP1$ hide

1 $TBOT_START$ = $TBOT_END$ - 10 output
1 $TDIV$ = $TBOT_END$ + $TDIV_TBOT_INT$ output

0 $RESIZE0$ = $NANC$/$NPOP1$ hide
```

Another solution is to actually estimate the time between time events, as shown on the left.

In this case, we would estimate the parameter $\$TDIV_TBOT_INT\$$

And then in complex parameters:

$\$TDIV\$ = \$TBOT_END\$ + \$TDIV_TBOT_INT\$$

Estimation file (.est)

Estimation file ("NoMigSan_Maya.est")

```
NoMigSan_Maya.est
// Search ranges and rules file
// *****

[PARAMETERS]
//#isInt? #name      #search.#min  #max
//all Ns are in number of haploid individuals
1 $NPOP1$          unif  10   1e5   output
1 $NPOP2$          unif  10   1e5   output
1 $NANC$           unif  10   1e5   output
1 $NBOTP1$         unif  1    1e3   output  bounded

1 $TBOT_END$       unif  100  1e4   output
0 $TDIV_TBOT_INT$ unif  10   1e3   hide

[RULES]

[COMPLEX PARAMETERS]

0 $RES_BOT_START$ = $NBOTP1$/$NPOP1$      hide
0 $RES_BOT_END$   = $NPOP1$/$NBOTP1$      hide

1 $TBOT_START$   = $TBOT_END$ - 10        output
1 $TDIV$         = $TBOT_END$ + $TDIV_TBOT_INT$ output

0 $RESIZE0$      = $NANC$/$NPOP1$         hide
```

Finally, a note about inferring bottlenecks associated with founder events.

It is difficult to jointly infer the duration and Effective population size of a bottleneck.

Instead, we can infer the bottleneck intensity, which is given by

$$I_B = \frac{\text{Time Duration Bottleneck (generations)}}{(2 * \text{Effective size during bottleneck})}$$

Thus, we usually fix the duration of the bottleneck and infer the effective size.

In this case, we fix the duration of the bottleneck to be 10 generations.

If \$NBOTP1\$ is larger than 500, then actually we would infer no bottleneck, as $I_B < 0.01$ ($10/(2*500)$).

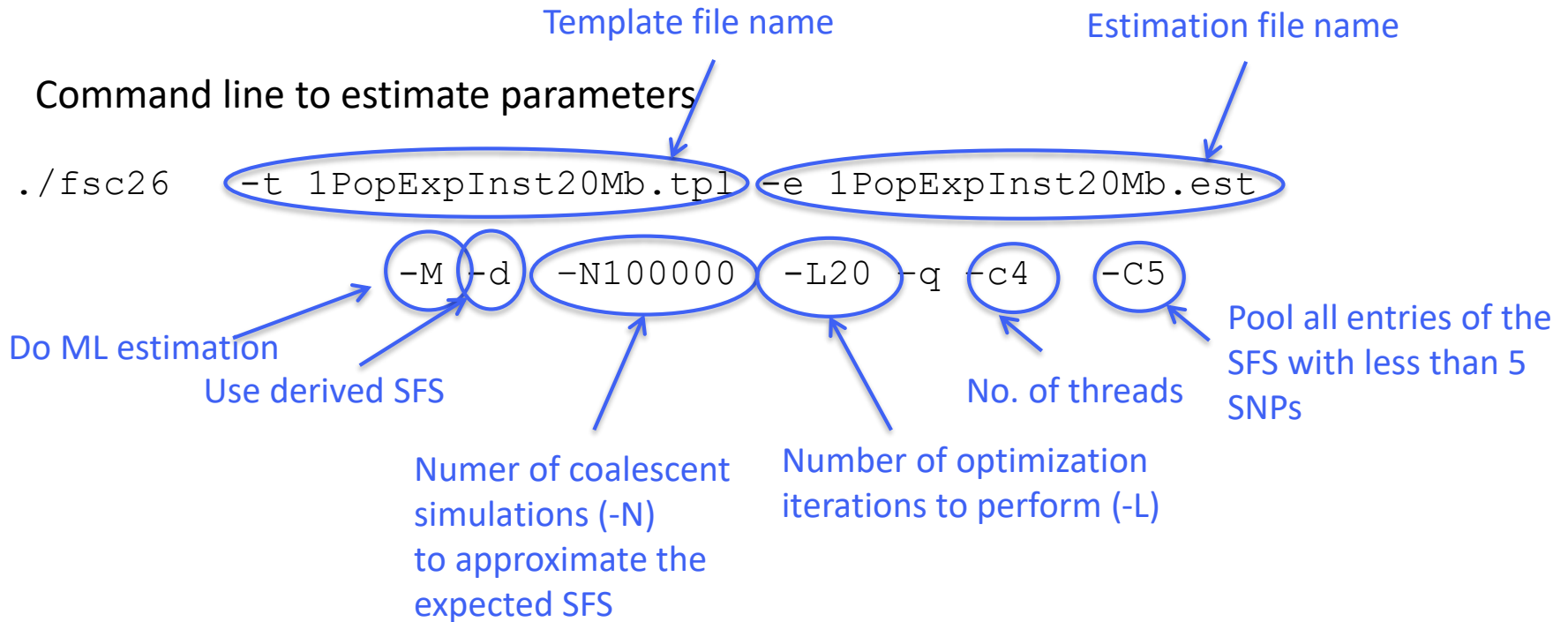
Launching parameter estimations

Command line to estimate parameters

```
./fsc26 -t 1PopExpInst20Mb.tpl -e 1PopExpInst20Mb.est  
-M -d -N100000 -L20 -q -c4 -C5
```

Observed SFS file must have the same name as template file and extension
_DAFpop0.obs. e.g. `1PopExpInst20Mb_DAFpop0.obs`

Launching parameter estimations



Observed SFS file must have the same name as template file and extension
_DAFpop0.obs. e.g. 1PopExpInst20Mb_DAFpop0.obs