

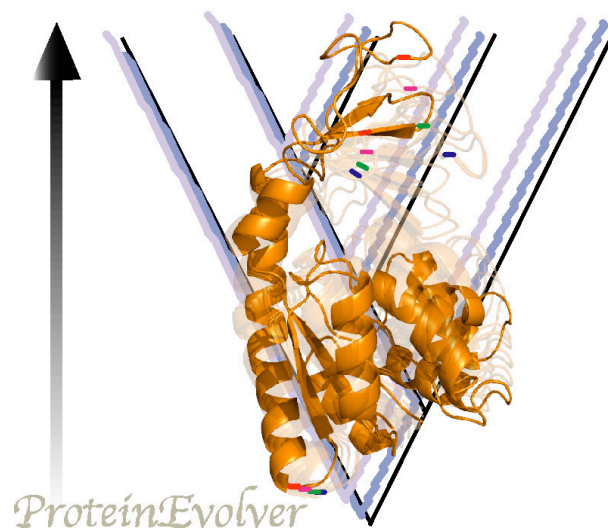
## Documentation for **ProteinEvolver**

*Simulation of protein evolution along phylogenies accounting for structural protein stability*

Current version is 1.2.0

© 2012-2013 Miguel Arenas ([marenas@cbm.uam.es](mailto:marenas@cbm.uam.es) - [miguelmmmab@gmail.com](mailto:miguelmmmab@gmail.com)),  
David Posada ([dposada@uvigo.es](mailto:dposada@uvigo.es)) and Ugo Bastolla ([ubastolla@cbm.uam.es](mailto:ubastolla@cbm.uam.es)).

April 9, 2013



# Contents

<b>Disclaimer</b>	<b>3</b>
<b>Credits</b>	<b>3</b>
<b>1. Purpose</b>	<b>3</b>
<b>2. Executables and compilation</b>	<b>4</b>
<b>3. ProteinEvolver Usage</b>	<b>4</b>
<b>3.1. Command line</b>	<b>5</b>
<b>3.2. Parameters file</b>	<b>6</b>
3.2.1. Arguments for the parameters file	6
<b>3.3. Input files for the SCS models</b>	<b>11</b>
3.3.1. Arguments for the file of SCS models	11
<b>3.4. Additional input files</b>	<b>13</b>
3.4.1. User-specified tree. No coalescent simulation	13
3.4.2. Recombination hotspots	13
3.4.3. MRCA/GMRCA sequence user-specified	14
3.4.4. Empirical user-specified amino acid matrix	15
3.4.5. Site by site variable substitution rate	15
<b>3.5. Default settings</b>	<b>16</b>
<b>3.6. Output Files</b>	<b>19</b>
3.6.1. Output files for the SCS models	19
<b>3.7. Solving message errors</b>	<b>20</b>
<b>4. ProteinEvolver model</b>	<b>21</b>
<b>4.1. Markov substitution models</b>	<b>22</b>
4.1.1. DNA models	22
4.1.2. Empirical amino acid models	22
<b>4.2. SCS models</b>	<b>22</b>
4.2.5 Crossing the SCS models with classic Markov substitution models	25
4.2.6 Evolution of sequences along trees under the SCS models	26
<b>4.3. Coalescent simulations</b>	<b>26</b>
4.3.1. Recombination	27
4.3.2 Migration	32
4.3.3 Convergence of demes	35
4.3.4 Demography	35
4.3.5 Tip Dates	36
<b>5. Program benchmarking</b>	<b>37</b>
<b>6. History</b>	<b>37</b>
<b>7. Acknowledgments</b>	<b>39</b>
<b>8. References</b>	<b>39</b>

## Disclaimer

This program is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation; either version (at your option) of the License. This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details. You should have received a copy of the GNU General Public License along with this program; if not, write to the Free Software Foundation, Inc., 59 Temple Place – Suite 330, Boston, MA 02111-1307, USA.

## Credits

This program was developed at the Bioinformatics Unit, Centre for Molecular Biology “Severo Ochoa”, Consejo Superior de Investigaciones Científicas (CSIC). Madrid, Spain.

## 1. Purpose

*ProteinEvolver* generates samples of coding and amino acid sequences evolved along phylogenies under structurally constrained substitution (SCS) models. These models consider the stability of the protein structure to evaluate candidate mutations. Thus, the mutations can be fixed (substitutions) or rejected depending on the energy of the protein structure with the mutated sequence.

The simulation of molecular evolution occurs along phylogenetic histories, which can be either user-specified or simulated by the coalescent modified with recombination (including recombination hotspots and coldspots), migration, demographics and longitudinal sampling (among other options).

### ***What is exactly implemented in ProteinEvolver?***

*ProteinEvolver* implements protein stability substitution models that consider contact matrices, configurational entropy per residue in unfolded and misfolded proteins, configurational entropy offset (misfolded) and energy functions [1, 2]. In addition, these structural models can be crossed with parametric DNA substitution models such as *JC* [3], *K80* [4], *F81* [5], *HKY* [6], *SYM* [7], *GTR* [8] and even *GTnR* [extended from, 8] for the simulation of DNA data or, with empirical amino acid substitution models such as *Blosum62* [9], *CpRev* [10], *Dayhoff* [11], *DayhoffDCMUT* [12], *HIVb* [13], *HIVw* [13], *JTT* [14], *JonesDCMUT* [12], *LG* [15], *Mtart* [16], *Mtmam* [17], *Mtrev24* [18], *RtRev* [19], *VT* [20], *WAG* [21] or any user-specified matrix, for the simulation of proteins. Indeed, heterogeneous substitution rates among sites by a gamma distribution (+G) and proportion of invariable sites (+I) are also implemented. Furthermore, the user can modify the substitution rate for each site, for example allowing fix particular sites (for example, catalytic sites).

The molecular evolution is simulated forward in time along the phylogeny. The user can either specify a particular phylogenetic tree or, simulate a coalescent history [22, 23]. In the latter case, *ProteinEvolver* implements the coalescent with recombination [22] (which can be homogeneous or heterogeneous along the sequence [following, 24]), variable population size (by a growth rate and demographic periods), a variety of migration models (such as island [25], stepping-stone [26] and continent-island [27])

with temporal variation of migration rates and convergence of demes or subpopulations, the simulation of haploid or diploid data, and longitudinal sampling [see, 28].

The user has to specify a PDB file (which can be downloaded from the Protein Data Bank, <http://www.rcsb.org/PDB/home/home.do>) and a sequence, both are assigned to the root of the phylogeny. Then, a variety of input evolutionary parameters should be specified (see next sections).

The output allows multiple options. Alignments can be printed in *phylip*, *fasta* and *nexus* formats. The ancestral sequence (MRCA or GMRCA, most recent common ancestor and grand most recent common ancestor, respectively [see, 29, 30]) can be also printed. Energies for the native structural protein can be printed as a function of the temperature and energies for the simulated proteins can be also printed to study how the incorporation of substitution events influences the structural protein stability. When coalescent histories are simulated, the simulated tree or ancestral recombination graph (ARG) [31] can be also printed.

We recommend check the attached folder “example\_input\_files”, it contains a variety of examples of evolutionary scenarios that can be simulated using *ProteinEvolver*.

## 2. Executables and compilation

Executable files are provided for Linux Debian and MacOS X (Intel and G4 processors), and a Makefile is provided for compilation in any OS with a C compiler. This makefile can be optimized for different users, for example using the optimization option `-fast` instead of `-O3`, for Mac processors. To compile the program type (it may take a few minutes): *make all*

It should print something like:

*Building ProteinEvolver version 1.2.0*

*gcc -c -O3 -Wall ProteinEvolver1.2.0.c*

*gcc -lm -O3 -Wall -o ProteinEvolver1.2.0 ProteinEvolver1.2.0.o*

*Finished compiling.*

A second makefile “Makefile\_MPI” is provided to compile a MPI version (which could be convenient for diverse simulation experiments [e.g., 32, 33]). This Makefile might need some modifications for particular OS. To compile the program type: *make -f Makefile\_MPI*. MPI libraries have to be installed in the host for running *ProteinEvolver* in parallel. The minimum number of processors is two. An example of execution for 3 processors is the next: *mpirun -np 3 ProteinEvolver1.2.0*

## 3. ProteinEvolver Usage

The input of the program consists of a series of arguments and parameter values (Table 1) that can be written in the command line or, more conveniently, specified in a text file called “parameters” that should be located at the same directory of the executable. These arguments include the parameter values used in the simulations and several printing options that control the amount of information that is sent to the console or the type of output files.

In addition, other input files are required for diverse specifications (see sections 3.3 and 3.4):

- For protein stability substitution models: particular settings for the structural model, PDB file and amino acid contacts matrix.
- A user- specified empirical amino acid model.

- A vector to modify the heterogeneous substitution rates per site.
- User-specified tree/s.
- A user-specified sequence for the root of the phylogeny.
- Settings for recombination hotspots.

### 3.1. Command line

In Mac systems, the command line is provided by the Terminal, while in Windows, this is provided by the windows console or command prompt. If the user specifies any argument in the command line, *ProteinEvolver* will use the values specified for those parameters, and default values for the parameters not included in the command line. It is highly recommended check the information shown in the screen to be sure that the simulations have been parameterized as intended. Some example command lines are the following:

Amino acid data by the empirical JTT substitution model applied along a coalescent tree,

```
./ProteinEvolver1.2.0 -n2 -s8 150 -e200 1 -@JTT
```

Amino acid data by the structural protein stability (SCS) × JTT substitution model applied along a simulated coalescent tree. Note that it also requires the input files *Pop\_evol.in* (and derived files) and, *seqGMRCa*.

```
./ProteinEvolver1.2.0 -n2 -s8 255 -e200 1 -@JTT -f20 0.04 0.06 0.05 0.05 0.08 0.02 0.05 0.05 0.03 0.07 0.04 0.06 0.05 0.05 0.05 0.05 0.05 0.05 0.04 0.06 -zPop_evol.in -xseqGMRCa
```

DNA data by the classic JC substitution model applied along a user-specified tree in the input file *treefile*,

```
./ProteinEvolver1.2.0 -n2 -ptreefile -v1 0.5 -a0.7 -i0.5
```

Coding DNA data by the structural protein stability (SCS) × JC substitution model applied along a user-specified tree in the input file. Note that it also requires the input files *treefile*, *Pop\_evol.in* (and derived files) and, *seqGMRCa*.

```
./ProteinEvolver1.2.0 -n2 -ptreefile -v1 0.5 -f4 0.25 0.20 0.30 0.25 -mPop_evol.in -xseqGMRCa -a1 0.7 -i0.5
```

DNA data according to the classic JC substitution model applied along an ARG simulated by the coalescent with recombination (including recombination hotspots), longitudinal sampling, demographic periods, migration by an island model, variable migration rate according to temporal periods, convergence of demes, vector of rates per site and multiple output files. Note that it also requires the input files *UserHetRec*, *HetRatesVector* and *seqGMRCa*.

```
./ProteinEvolver1.2.0 -n2 -s8 765 -e1000 2 -=4 1995 1 1 2003 4 6 1997 2 3 2001 7 8 -/1200 -g1 3 1000 1250 1000 1300 1550 2000 1560 1000 3000 -q1 4 2 2 3 1 -t3 100 800 0.002 0.001 0.003 -%1 1 2 10000 -r2.3e-6 -hUserHetRec -o0.1 -u4.1e-5 -f4 0.25 0.20 0.30 0.25 -a0.7 -i0.52 -_HetRatesVector -xseqGMRCa -bsequences -c1 1 0 -jtrees -ktimes -dbreakpoints -*NetworkFile -y2 -#245
```

**We strongly recommend the use of the parameters file (instead of the command line, see next section) in order to avoid text errors and a better checking of the input settings.**

### 3.2. Parameters file

If no argument is specified, the program will only read the text file “*parameters*” which should be placed at the same directory of the executable. If no arguments are specified in the command line, and there is no a “*parameters*” file, the program will stop and throw an error.

In this file anything within brackets will be ignored. Examples of the “*parameters*” file are included in the distribution. Then, only type:

*./ProteinEvolver1.2.0*

#### 3.2.1. Arguments for the parameters file

Possible arguments for the parameters file are the following (# means number, *NAME* means a word):

##### General settings

*-n# : Number of replicates*

The number of samples to be generated. Each sample is an independent realization of the evolutionary process. This specification is mandatory. Example: *-n10*

##### Coalescent settings

*-s# #: Sample size; Number of sites*

The number of sequences to be generated for each sample and the total sequence length (in nucleotides or amino acids). Example: *-s6 255*

*-e# #: Effective population size; Haploid / Diploid*

The effective size (*N*) of the population from which the sample was theoretically drawn. If there are several demes, this argument is the effective population size for each deme. The second number means that the data set can be simulated as haploid (1) or diploid (2). Example: *-e100 1*

*-=# : Tip dates*

The time of the tips can be different with this option. For example, 4 samples: 1995:sequence 1; 2003: sequences 4 and 6; 1997: sequences 2 and 3; 2001: sequences 7 and 8. This option does not work if there is any convergence of demes at younger times. Example (of above): *- =4 1995 1 1 2003 4 6 1997 2 3 2001 7 8*

*-/# : Generation time*

The time per generation. Example: *-/300*

*-g0 # or -g1 # (# # #) : Demographics settings. Exponential growth rate or Demographic periods*

The first number specifies the model, exponential growth rate (0) or demographic periods (1). These parameters are looking back in time, so it is not a good idea to

specify a negative growth rate for the last period, as the coalescent time could become infinite in the past.

*Rate of exponential growth per individual per generation*, after “0” the growth rate must be specified. Example: `-g0 1e-5` (= `-g0 0.00001`)

*Demographic periods*, after “1” the user has to specify the number of periods (from the present to the past) and  $N$  during those periods. The first number here specifies the number of periods. For each period should be three consecutive numbers indicating the size  $N$  at the beginning and at the end of the period, and the duration of the period in generations. Example: `-g1 3 1000 1250 1000 1300 1550 2000 1560 1000 3000`

The exponential growth rate during the period (positive or negative) will be deduced from the specified  $N$  at the beginning and at the end of the period. The growth rate derived for the last period will continue into the indefinite past. This implementation is borrowed from [34]. This option is incompatible with the exponential growth rate option (`-g`). Again, these parameters are looking back in time, so it is not a good idea to specify a negative growth rate for the last period, as the coalescent time could become infinite in the past.

*-q# # (#): Migration model and population structure*

The first number specifies the migration model (island model=1, stepping-stone model=2, continent-island model=3). The second number specifies the total number of demes or subpopulations sampled. The next  $n$  numbers specify the number of individuals (or sequences) per deme (note that the specified sample size (`-s`) must be equal to the sum of these). For the island-continent model, deme #1 will be the continent while the other demes will be islands (see details in section 4.3.2). Example: `-q2 2 3 3` (a stepping-stone model, two demes with three samples each).

*-t# (#)(#): Migration rate*

This parameter introduces the migration rate, which can be constant or variable with time according to temporal periods. The first number specifies the number of temporal periods, then:

For only 1 period, the second number is the migration rate (constant). Example: `-t1 0.001` (only 1 period with migration rate = 0.001).

For more than 1 period, the second number/s are the time/s for the beginning of a new migration rate and the third/s numbers are the corresponding migration rate/s for each period. Example: `-t2 100 0.001 0.005` (2 periods, the first period occurs from  $t = 0$  to  $t = 100$  with a migration rate = 0.001, the second period occurs from  $t = 100$  to the end of the simulation with a migration rate = 0.005). Example: `t3 100 800 0.002 0.001 0.003` (3 periods: from  $t = 0$  to  $t = 100$  with migration rate = 0.002, from  $t = 100$  to  $t = 800$  with migration rate = 0.001, from  $t = 800$  to the end of the simulation with a migration rate = 0.003).

*-%# (# # #) : Events of convergence of demes*

The first number specifies the total number of convergent events. For each convergence event should be three consecutive numbers. The first number and the second number are the numbers of the demes to converge. The third number is the time to that convergence. With this option the user can build the demes evolutionary tree but it is only available when the migration model is activated (despite the migration rate could be zero). Examples: `-%1 1 2 2000` (for 2 initial demes (`-q2`), convergence of deme 1 with deme 2 at time 2000 to create a new deme 3). `-%3 1 2 400 3 4 1900 5 6 2000` (for 4 initial demes (`-q4`) convergence of deme 1 with deme 2 at time 400 to create a new

-v1 #, -v6 #####, -v12 #####: relative substitution rates



The first number specifies the total number of relative substitution rates that can be 1 (for transition/transversion rate ratio), 6 (for relative symmetric substitution rates) or 12 (relative asymmetric substitution rates).

If the first number = 1, the second number is the transition/transversion rate ratio. When set to 0.5, the probability of transitions and transversions are the same, like in models “*JC*” [3] or “*F81*” [5] models. Example: *-v1 2.1*

If the first number = 6, the following 6 numbers are the relative symmetric substitution rates  $A \leftrightarrow C$ ,  $A \leftrightarrow G$ ,  $A \leftrightarrow T$ ,  $C \leftrightarrow G$ ,  $C \leftrightarrow T$  and  $G \leftrightarrow T$ . Example: *-v6 1.0 5.0 1.0 1.0 5.0 1.0*.

If the first number = 12, the following 12 numbers are the relative asymmetric substitution rates AC, CA, AG, GA, AT, TA, CG, GC, CT, TC, GT and TG. Note that some calculations under this option can be problematic (complex eigenvalues/vectors) if rates are too asymmetric, especially with codon models. Example: *-v12 1.1 1.2 1.3 1.4 1.5 1.6 1.7 1.8 1.9 2.0 1.0 1.0*

*-mNAME : SCS for coding DNA simulation*

This argument should be followed by the name of the file with settings for the SCS models (see section 3.3). When this argument is specified the program will simulate coding DNA data (which of course can be translated to amino acid data). In particular, under this model the mutations and substitutions (accepted mutations) occur at DNA level but their energies are tested at amino acid level taking into account the stability of the protein structure (see section 4.2). Importantly, do not activate this argument and the argument “-z” (which directly simulates amino acid data considering the stability of the mutations in the protein structure) at the same time. Example: *-mPop\_evol.in*

*-@NAME : empirical amino acid model*

This option allows for the simulation of proteins. The empirical amino acid models implemented in *ProteinEvolver* are the following: *Blosum62* [9], *CpRev* [10], *Dayhoff* [11], *DayhoffDCMUT* [12], *HIVb* [13], *HIVw* [13], *JTT* [14], *JonesDCMUT* [12], *LG* [15], *Mtart* [16], *Mtmam* [17], *Mtrev24* [18], *RtRev* [19], *VT* [20], *WAG* [21]; in this case just specify the name of the empirical amino acid model.

But in addition, it is possible to specify a user-defined empirical amino acid model by an input file (see section 3.4); in this case just specify the name of the input file with the empirical amino acid model. Example: *-@JTT*. Example: *-@UserEAAM*.

*-zNAME : SCS model for coding protein simulation*

This argument should be followed by the name of the file with settings for the SCS models (see section 3.3). Importantly, when this argument is specified the program will simulate protein data. In particular, under this model the mutations and substitutions (accepted mutations) occur at the amino acid level and their energies take into account the stability of the protein structure. Importantly, do not activate this argument and the argument “-m” at the same time. Example: *-zPop\_evol.in*

*-a# : alpha shape of the gamma distribution*

A gamma distribution (+G) can simulate substitution rate variation among sites [35]. Alpha is the shape of this distribution. Smaller alphas imply stronger rate variation. Example: *-a0.7*.

*-\_NAME : factor for variable substitution rate site by site*

After this argument the user should specify the name of an input file, which contains a vector with values for each site (see details in section 3.4). Example: `-_HetRatesVector`

*-i# : proportion of invariable sites*

A proportion of sites can be set to be invariable (+I). Example: `-i0.3`

*-xNAME : GMRCA input file*

The user can specify its own MRCA/GMRCA sequence in a text file (see section 3.4). This file must contain only a DNA or amino acid sequence (depending on the substitution model applied), and the length has to be equal to the number of sites specified in the arguments “-s” or in the tree file “-p”. By default, the *MRCA/GMRCA sequence* is simulated from the nucleotide or amino acid frequencies (see above, `-f` argument). Example: `-xseqGMRCA`

## Output settings

*-bNAME : print sequences*

When this argument is invoked, aligned sequences are printed to the specified text file in the *Results* folder. Example: `-bsequences`

*-c# # #: format for printing sequences*

This option specifies the format of the output alignments. The first argument indicates *Phylip* sequential, *Fasta* and *Nexus* formats (1-3, respectively). The second argument indicates that alignments for each replicate are printed into a single file (0) or different files (1). The third argument prints the sequence of all internal nodes (1). Example: `-c1 1 0` (which means: *phylip* format, a file for each replicate, sequences corresponding to internal nodes are not printed).

*-\$: print catMRCA/GMRCA*

This option prints the ancestral catMRCA/GMRCA sequences in output files. These output files will be incorporated to the *Results* folder. Example: `-$`

*-jNAME : print genealogies*

When this option is specified, genealogies for each recombinant fragment are printed, in *Newick* format, to the specified text file, in the *Results* folder. Example: `-jtrees`

*-kNAME : print times*

When this argument is specified the coalescent times for each genealogy will be printed to the specified text file, in the *Results* folder. This option will slow down the simulations. Example: `-ktimes`

*-\*NAME : print ARGs*

When this argument is specified the ancestral recombination graph (ARG) will be printed to the specified text file in the *Results* folder. Then, this file can be directly introduced into the *NetTest* web server (<http://darwin.uvigo.es/software/nettest/>) [36] in order to visualize the ARG. Example: `-*NetworkFile`

*-dNAME : print breakpoints*

When this argument is specified breakpoint positions are printed into the specified text file, in the *Results* folder. This option only works for coalescent simulations. Example: `-dbreakpoints`

### Other settings

*-y# : noisy level*

This option controls the level of information that will be printed to the screen.

0 : does not print run information, just the simulation progress.

1 : run settings and run information summarizing the simulations.

2 : calculation status, initial demes and event information, run settings for each replicate

3 : print ancestral status for each sequence at each event + *MRCA* status, tip dates insertion, demes, molecular evolution.

4 : potential recombining locations (*g* and *G* vectors) and information about recombinant fragments evolution.

Note that higher levels of noisy will slow down the simulations. The default level is 1

Example: *-y1*

*-## : Seed*

Seed for the random number generator. If no seed is specified, the computer clock will be used. When a seed is fixed the process can be always reproduced if same settings are specified. Example: *-#386658297*

## 3.3. Input files for the SCS models

The main input file to specify SCS models must be specified in the parameters file by the argument “*-m*” (for the simulation of coding DNA data) or “*-z*” (for the simulation of proteins). Examples of the main input file for these models can be found in the “examples” folder, with name “Pop\_evol.in”. Theory is described in the section 4.2.

Importantly, this simulation requires that the MRCA/GMRCA sequence must codify the amino acid sequence of the PDB file (for coding data simulation) or be equal to the amino acid sequence of the PDB file (for protein simulation), see details in the theory section. **As a consequence, for these models it is highly recommended to fix the MRCA/GMRCA sequence by an input file (see section 3.2, argument “*-x*” in the parameters file) because a MRCA/GMRCA sequence computed using the nucleotide or amino acid frequencies could not satisfy such a condition.**

### 3.3.1. Arguments for the file of SCS models

Possible arguments are the following (*#* means number, *NAME* means a word):

#### Structural settings

*PDB= NAME: File from the Protein Data Bank*

A file from the Protein Data Bank (<http://www.rcsb.org/PDB/home/home.do>) (PDB) must be specified. Example: *PDB= ITRE.PDB*

*CHAIN= NAME: Chain of the PDB file*

The particular chain of the PDB file to be considered must be specified. Example: *CHAIN= A*

*FILE\_STR= NAME: List of contact matrices*

A file with list of contact matrices must be specified. The folder with example input files already contains a file “*structures.in*” with a huge list of contact matrices downloaded from the PDB and which can be applied to compute energies from any PDB protein structure. However, the authors could provide other files with contact matrices (considering particular user-specified details) upon request. Example: *FILE\_STR= structures.in*

*TEMP= #: Temperature*

The temperature used to compute protein energies must be specified (see section 4.2). We recommend a range between 1.25 and 2.0. Example: *TEMP= 1.5*

*S0= #: Configurational entropy per residue (unfolded)*

The configurational entropy per residue for the unfolded protein must be specified. We recommend a range between 0.025 and 0.075. Example: *S0= 0.05*

*SC1= #: Configurational entropy per residue (misfolded)*

The configurational entropy per residue for the misfolded protein must be specified (see section 4.2). We recommend a range between 0.025 and 0.075. Example: *SC1= 0.05*

*SC0= #: Configurational entropy offset (misfolded)*

The configurational entropy offset for the misfolded protein must be specified (see section 4.2). By default this it is 0. Example: *SC1= 0.0*

*REM3= #: Third cumulant in REM calculation*

The third cumulant in REM calculation (to compute the structural protein energy) must be specified (see section 4.2). By default this it is 0. Example: *REM3= 0*

*NEUTRAL= #: If 1, Neutral landscape (**neutral SCS model**), otherwise population size dependent selection (**fitness SCS model**)*

If this parameter is set to 1, the neutral substitution models of structural protein stability will be applied (see section 4.2). Note that this model does not use the population size (next parameter), by default we recommend this neutral model. Example: *NEUTRAL= 1*

*NPOP= #: Population size for protein structures*

The population size to simulate molecular evolution under fitness substitution models of structural protein stability. Importantly, note that this option requires that the setting “*NEUTRAL=*” must be set to 0, so a fitness substitution model is specified. Example: *NPOP= 10*

### Other settings

*TYPE\_BL= #: Type of branch lengths (by mutations or substitutions)*

This argument specifies if branch lengths are either considered by events of mutation (1) or substitution (2). Note that the expected number of events = branch length × number of sites. Example: *TYPE\_BL= 2*

*OUTPUT\_LEVEL= #: Amount of output files*

This argument specifies the amount of output files printed using protein stability substitution models (see section 3.6.1). “2”, all output files are printed. “1”, only the final output files are printed. “0”, the output files are not printed. Importantly, note that a total of 4 output files related with these substitution models are printed per branch and

therefore, it can slow down the simulations and may need a lot of space in the hard disk.  
Example: `OUTPUT_LEVEL=1`

### 3.4. Additional input files

The following input files are optional and if required, must be introduced in the same directory of the executable of *ProteinEvolver*. Indeed, note that the settings contained in these input files cannot be specified in the command line. The name of these files must be introduced in the main settings (parameters file or command line).

#### 3.4.1. User-specified tree. No coalescent simulation

Alternatively to the coalescent simulation, the user can specify a tree in order to evolve sequences along its branches. This possibility is convenient when the user already has a tree (e.g., inferred from real data) or in order to avoid coalescent assumptions such as the molecular clock. This procedure is similar to other software such as *SeqGen* [37] or *Evolver* [38] (note that models implemented in these programs do not consider protein structures).

The tree must be incorporated into an input file, which should be specified from the main settings by the argument “-p” (e.g., `-ptreefile`; see section 3.2).

Then, the input file consists of a text line with a range of sites and a tree in Newick format. The range specifies the first and last site (nucleotide or amino acid position) where the following tree should work. Importantly, all sites should be covered by a phylogenetic tree, so it always should start by 1 and finish by the total number of sites. Only one tree is allowed because the whole protein must exist in all nodes. The file does not allow empty lines. Trees should be rooted.

An example is shown below,

```
1 765 ((taxonA:0.1,taxonB:0.1):0.4,(taxonC:0.1,taxonD:0.1):0.5,taxonE:1.0);
```

Note that by this procedure the total number of sites, number of taxa and name of taxa are specified. In that example, a total of 765 sites with 5 taxa (taxonA, taxonB, taxonC, taxonD, taxonE). When this option is applied all coalescent input settings are ignored (see section 3.2).

#### 3.4.2. Recombination hotspots

The user can optionally simulate recombination hotspots following the algorithm implemented in SNPsim [24] when using coalescent simulations. This file must be specified in the main settings (`-hNAME`; see section 3.2.1) and should be placed in the directory of the executable. An example file “*UserHetRec*” is included in the folder with examples.

Possible arguments for *recombination hotspots file* are (# means number):

`-k#` : *Hotspot recombination rate*

Expected recombination rate at the hotspots sites. If the hotspots are homogeneous (option `-t#` is not invoked) all the hotspots have the same rate. Example: `-k1e-4` (= -

*k0.0001)*

*-h# : Expected number of hotspots*

This is the expected number of hotspots for a given sample in the absence of interference. This parameter corresponds to the intensity parameter for a Poisson distribution from which the actual number of hotspots is drawn. For a given sample, the actual number of hotspots will change around this value. It does not have to be an integer. NOTE: When interference is specified we need to divide this number by the interference interval (*-z#*) to obtain the expected number of hotspots. It does not have to be an integer. Example: *-h1.1*

*-q# : Fixed number of hotspots*

This option fixes the number of hotspots inside the region of interest, so every sample will have the same number. In this case the hotspot locations are chosen from a uniform distribution. If the hotspots overlap, they will be displaced to the closest available location. Note that in this case no recombination events will originate from a hotspot located outside the region of interest. Example: *-q3*

*-v# : Hotspot imprecision*

The hotspot imprecision corresponds to the variance of a Normal distribution for the specific site to recombine around the hotspot center (chosen by a Poisson process). The bigger the imprecision, the wider is the hotspot. If the imprecision is 0, all the recombination events happen exactly at the hotspot center. See figures 2 and 3. Example: *-v0*

*-m# : Hotspot width*

This option specifies the width of the hotspots. In this case any site in the hotspot has the same probability of recombination. If the width is 1 all the recombination events happen exactly at the hotspot center. This parameter has to be bigger than 0. See figures 2 and 3. Example: *-m1*

*-t# : Hotspot heterogeneity*

This parameter indicates that there is hotspot heterogeneity, that is, hotspots may have different recombination rates. This heterogeneity is accomplished through the use of the continuous gamma distribution. The shape parameter of this distribution (%) will control the strength of this heterogeneity. The smaller the shape the stronger the heterogeneity. This is similar to the application of Yang [35]. Example: *-t0.5*

*-z# : Hotspot interference*

This parameter indicates whether the location of the hotspots is not independent of each other. If this parameter is 1 there is no interference, if it is between 0 and 1 hotspots tend to cluster, and if it is bigger than 1 hotspots will tend to be pushed away from each other. Example: *-z1*

### 3.4.3. MRCA/GMRCA sequence user-specified

By default, the GMRCA or MRCA sequence is simulated according to the nucleotide or amino acid frequencies. However, the user can optionally specify its own root sequence by a text file (the name of this file must be specified in the main settings (*-xNAME*),

which should be located in the directory of the executable. The file just contains a single sequence. An example file “*seqGMRC*” is included in the folder with examples.

The option is recommended for protein stability substitution models where the sequence assigned to the root node must be equal (for the simulation of proteins) or codify (for the simulation of coding DNA) for the amino acid sequence of the PDB file.

**Consequently, under SCS models, it is highly recommended introduce directly (by using this MRCA/GMRCA file) the amino acid sequence of the PDB file or a coding sequence that codifies for such a PDB sequence.**

### 3.4.4. Empirical user-specified amino acid matrix

The user can optionally specify a particular empirical amino acid matrix from a text file, which has to be located in the directory of the executable. An example file “*userEAAM*” is included in the folder with examples.

This file consists in two sets of parameter values: First, the substitution rates for each amino acid (values must start after the letter for the corresponding amino acid). Second, 20 amino acid frequencies can be introduced (values must start after a “z”). The order of amino acids must be the following: A R N D C Q E G H I L K M F P S T W Y V.

### 3.4.5. Site by site variable substitution rate

The user can optionally modify the substitution rate site by site. This file must be specified in the main settings by “- *NAME*” and should be located in the directory of the executable. An example file is included with the package (file named “*HetRatesVector*”) and another example is shown below.

The file consists on two sets of numbers specified after an “r”. First, the sequence length (255 in the example). Second, values for each site separated by a single space. Note that the number of values must be equal to the number of sites.

[illegible]

The value for each site consists in a factor that multiplies the original substitution rate for that site and should be between 1 and 0. Thus, for example a value of 1 means that the original substitution rate remains invariant while a value of 0 means that the substitution rate is 0 and therefore that site will never mutate. This option can be useful when the user already know the functional importance of the amino acids in the protein, so for example catalytic sites could have a value of 0 in order to keep the protein activity.

### **3.5. Default settings**

By default *ProteinEvolver* simulates 10 samples of 6 individuals with 201 sites, a constant effective size of 1000, constant size, no recombination, no migration and a mutation rate of  $1e-7$  under the *JC* DNA substitution model with nucleotide frequencies equal to 0.25. Noisy level is 1. The sequences will be printed to the *sequences* output file in the *Results* folder. To run the program with the equivalent arguments we should type:

```
./ProteinEvolver1.2.0 -n10
```

```
./ProteinEvolver1.2.0 -n10 -s6 201 -e1000 2 -u1.0e-07 -f4 0.25 0.25 0.25 0.25 -  
bsequences -y1
```



**Table 1. Key arguments for ProteinEvolver.** The user can specify several parameters to define different simulation scenarios. These arguments can be entered in the command line or read from text file.

Parameter	Arg	Example value	Application
Number of replicates	n	200	All
Sample size	s	8	All
Number of sites (bp or amino acids)	s	765	All
Effective population size	e	1000	All
Haploid / Diploid	e	1	All
Tip dates	=	2 1995 1 3 2003 4 8	All <sup>1</sup>
Generation Time	/	400	All
Exponential growth rate	g	$2.1 \times 10^{-5}$	Demography
Demographic periods <sup>2</sup>	g	1000 5000 200	Demography
Migration model	q	1	Migration
Number of demes	q	“4” 2 2 3 1	Demes/Migration
Population structure	q	4 “2 2 3 1”	Demes/Migration
Migration rate (constant or variable with time)	t	$1.2 \times 10^{-4}$ or 100 “0.001 0.005”	Migration
Convergence of demes <sup>3</sup>	%	1 2 5000	Demes/Migration
Homogeneous recombination rate	r	$5 \times 10^{-6}$	Homogeneous recombination/Recombination hotspots
Fixed number of recombination events	w	3	Homogeneous recombination/Recombination hotspots
Recombination hotspots	h	UserHetRec	Recombination hotspots
Substitution rate	u	$5.1 \times 10^{-4}$	All
Outgroup branch length	o	0.1	All
User-specified tree/s	p	Treefile	Genetic data simulation
Nucleotide frequencies	f	0.4 0.3 0.1 0.2	DNA / DNA – structural protein stability substitution models
Transition / transversion ratio	v	2.1	DNA / DNA – structural protein stability substitution models
Relative substitution rates	symmetrical v	1.0 2.3 2.1 3.0 4.2 1.0	DNA / DNA – structural protein stability substitution models
Relative	asymmetrical v	0.1 0.2 0.3 0.4 0.9	DNA / DNA – structural

substitution rates		0.6 0.9 0.8 0.9 0.1	protein stability substitution models
		0.2 0.3	
File for the coding DNA – SCS <sup>5</sup> models	m	Pop_evol.in	Coding DNA – SCS models
Amino acid frequencies	f	0.04 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.06	Amino acid / Protein – structural protein stability substitution models
Empirical amino acid model	@	JTT	Amino acid / Protein – structural protein stability substitution models
File for the Protein – SCS <sup>5</sup> models	z	Pop_evol.in	Protein – SCS models
Rate variation among sites <sup>4</sup>	a	0.4	All
Variable substitution rates by site	–	HetRatesVector	All
Proportion of invariable sites	i	0.2	All
User-defined sequence for GMRCA	x	seqGMRCA	All
Print sequences using diverse options	b, c	sequences, 1 1 0	All
Print GMRCA sequence	\$		All
Print simulated omega per site and/or per branch	+		Codon models with variable omega per site and/or branch
Print simulated trees	j	trees	All
Print simulated ARG	*	NetworkFile	All
Print times for genealogies	k	times	All
Print breakpoints	d	breakpoints	All
Apply a seed	#	3444556	All
Level of output information printed in the screen	y	2	All

<sup>1</sup>In presence of convergence of demes, the time of the tip dates must be younger than the time of the convergence of demes.

<sup>2</sup>from 1000 to 5000 effective size during 200 generations.

<sup>3</sup>deme 1 is converging with deme 2 at time 5000 to make a new ancestral deme 3.

<sup>4</sup>shape of the gamma distribution.

<sup>5</sup>SCS models can be neutral or fitness-based landscape.

### 3.6. Output Files

All output files produced by the program are written to the folder “*Results*”:

- *Sequences file*: contains the aligned sequences in *Phylip* sequential, *Fasta* or *Nexus* formats.
- *Breakpoints file*: contains the breakpoints ordered by event time. Disabled when input tree/s are user-specified.
- *Trees file*: contains the trees in newick format for each recombinant fragment.
- *Network file*: contains the simulated ARG by a branches format (each line contains two connected nodes) [36]. Disabled when input tree/s are user-specified. See section 4.3.1.2 for details.
- *Times file*: contains information about time and branch length for each tree.
- *CatMRCA/GMRCA*: contains the corresponding sequences catMRCA (concatenated MRCA fragments) and GMRCA (sequence of the GMRCA node) in the ARG.
- *Settings file*: contains a summary of applied settings showed in the screen at the end of the simulation. This option must be activated in the source code.

#### 3.6.1. Output files for the SCS models

The substitution models that account for the structural protein stability can produce additional output files that are printed within an output folder “*ProteinStability*”, located in the output folder “*Results*”.

Note that these output files are printed for each branch so be careful when printing all output files in extensive simulations, it could require a lot of space in the hard disk and might slow down the simulations. See the argument *OUTPUT\_LEVEL* in the section 3.3.1 to control the amount of printed output files.

The name for each output file consist in the prefix “*ProteinStability\_*” followed by the number of replicate “*Replicate#\_*”, the branch number “*Branch#\_*”, and details about the particular settings specified: PDB name “*NAME*”, temperature “*T#*”, entropy “*S#*”, protein population size “*N#*”, and optionally GC bias for coding DNA simulation “*GC#*” (note that “*#*” means a number).

Example from the simulation of coding DNA evolution,

```
ProteinStability_Replicate1_Branch0_1TREA_DeltaG_2.dat
ProteinStability_Replicate1_Branch0_1TREA_T1.80_S00.05_N10_GC0.40_ave.dat
ProteinStability_Replicate1_Branch0_1TREA_T1.80_S00.05_N10_GC0.40_dna.dat
ProteinStability_Replicate1_Branch0_1TREA_T1.80_S00.05_N10_GC0.40_stab.dat
ProteinStability_Replicate1_Branch0_1TREA_T1.80_S00.05_N10_GC0.40_final.dat
```

Example from the simulation of protein evolution,

```
ProteinStability_Replicate5_Branch12_1TREA_DeltaG_2.dat
ProteinStability_Replicate5_Branch12_1TREA_T1.80_S00.05_N10_ave.dat
ProteinStability_Replicate5_Branch12_1TREA_T1.80_S00.05_N10_stab.dat
ProteinStability_Replicate5_Branch12_1TREA_T1.80_S00.05_N10_final.dat
```

Files “*DELTA\_G\_2.dat*” show the protein energy values at amino acid level ( $\Delta G/L$ ) as a function of the temperature (T) computed by the REM2 (Random Energy Model) energy function (see section 4.2).

Files “*ave.dat*” show information about the user-specified settings, computed energy, fitness and, transition and mutation loads.

Files “*dna.dat*” show information about the user-specified settings and GC content.

Files “*stab.dat*” show, for each accepted mutation (substitution), the structural energy of the mutated and native proteins, derived fitness and number of mutation attempts to reach such a substitution. It is recommended to explore the influence of substitution events in the protein structure.

Files “*final.dat*” show the final results derived from the substitutions introduced on a branch. It indicates the mean of fitness, energy, entropy, transition and mutation loads and, rejected mutations.

### 3.7. Solving message errors

- Parameter values for the SCS models. Entropy and temperature must be chosen carefully, specially for fitness SCS models. See above recommended settings.

- There is a common message of error due to not enough memory in scenarios with very high recombination rates. Note that recombination increases the number of nodes of the ARG and the number of recombinant fragments within nodes (see section 4.3.1). Then, when  $\rho$  ( $= 4N\mu$ ) is very high, the ARG must save in memory a high amount of information, and with time such amount exponentially increases. Thus, at a given time the machine could not have enough available memory and the program stops with an error message like the following:

```
malloc: *** mmap(size=1584291840) failed (error code=12)
```

```
*** error: can't allocate region
```

```
Could not reallocate segments (1584144000 bytes)
```

In this situation we recommend to reduce the value of the following parameters: recombination rate, population size, length of the sequences and sample size. Of course, the other option is just to use another machine with more memory.

- Using a fixed number of recombination events keep in mind that the recombination rate introduced should generate an accordingly number of recombination events. Otherwise the simulation could never finish (because the assumption of the fixed number of recombination events is never successful) leading to crash the execution.

- Using growth rates or demographic periods you could find an error like:

```
ERROR: Coalescent time (nan) is infinite
```

```
This might suggest that the growth rate is too negative  
and the coalescent time is therefore infinite.
```

```
Try a smaller value
```

This is because the population increases too much going back in time and thus, the coalescent time gets infinite. The best option to solve this error is to reduce the growth rate. In the case of demographic periods this could be done for example by using longer times for such period.

- Asymmetric substitution rates (for example, from user-specified amino acid matrices) could be problematic. In these cases the program writes an output message about “complex roots”.

- Input tree from user. The file with the input tree must follow: first site should be 1 and the last site should be the total number of sites, no empty lines, the tree should be rooted.

If you experience any unexpected error or there is any doubt, please do not hesitate to contact us: miguelmmmab@gmail.com. Thanks for your contribution!

## 4. ProteinEvolver model

*ProteinEvolver* implements site-dependent substitution models that consider the stability of the protein structure. These models compute the structural energy of mutated proteins (including misfolded and unfolded configurations) given the entropy and the temperature. Then, mutations can be accepted (substitutions) or rejected according to the Moran’s model. Two substitution models have been implemented, the neutral model (which does not consider population size) and the fitness model (which needs a user-specified population size).

Indeed, these models can be crossed with classic site-independent substitution models (see below) to compute substitution rates per site and state. The implemented classical DNA Markov models of substitution are all that currently exist (e.g., *JC* [3], *K80* [4], *F81* [5], *HKY* [6], *SYM* [7], *GTR* [8] and even *GTnR* [extended from, 8]) for the simulation of DNA data. The implemented empirical amino acid substitution models are *Blosum62* [9], *CpRev* [10], *Dayhoff* [11], *DayhoffDCMUT* [12], *HIVb* [13], *HIVw* [13], *JTT* [14], *JonesDCMUT* [12], *LG* [15], *Mtart* [16], *Mtmam* [17], *Mtrev24* [18], *RtRev* [19], *VT* [20], *WAG* [21] and any user-specified matrix. In addition, heterogeneous substitution rates among sites by a gamma distribution +G and proportion of invariable sites +I are implemented. Furthermore, the user can alter the substitution rate for each site, this allows to fix particular sites. For example, it would make sense to fix sites related with the activity of the protein such as catalytic positions.

Using these substitution models, sequences can be evolved from the root to the tip nodes of a given phylogeny. Such a phylogeny can be user-specified or simulated by the coalescent. The program implements an extension of the coalescent with recombination, demographics and migration based on the neutral Wright-Fisher model [22, 25, 34] following [30, 39]. Given the specified recombination and migration rates (and other parameters like the effective population size (N) and growth rate) random genealogies are produced. Recombination hotspots are also implemented following *SNPsim* [24]. Several migration models (island [27], stepping-stone [26] and continent-island [27]) are allowed and migration rate can change with time. The evolution of the demes (or species tree) can be fixed. Complex demographic histories can be implemented by defining demographic periods in which population sizes augment, reduce, or remain constant. Samples can be collected at same or different times [see, 28] and simulated data can be haploid or diploid.

The result is a random sample of aligned coding DNA or amino acid sequences.

## 4.1. Markov substitution models

Once the sample genealogy has been constructed, nucleotide or amino acid data can be simulated along its branches. *ProteinEvolver* implements multiple nucleotide and amino acid Markov models through the specification of different parameters (including base frequencies, relative substitution rates, transition/transversion ratio, a proportion of invariable sites and rate variation among sites) and empirical matrices. Conveniently, the sequence of the root node can be specified at random, according to the nucleotide or amino acid frequencies, or the user can specify its own sequence (highly recommended for structural protein stability substitution models).

Variation in the substitution rate among sites can be user-established by the shape of a gamma distribution (+G). Indeed, *ProteinEvolver* also implements variation of substitution rates by a user-specified site-by-site vector. Finally, a proportion of invariable sites (+I) can be also simulated.

### 4.1.1. DNA models

*ProteinEvolver* implements the general time reversible model for nucleotide substitution (GTR) [8], a non reversible version (GTnR) [8], and models nested therein, like JC [3], K80 [4], F81 [5] or HKY [6].

### 4.1.2. Empirical amino acid models

*ProteinEvolver* implements a variety of empirical amino acid models: *Blosum62* [9], *CpRev* [10], *Dayhoff* [11], *DayhoffDCMUT* [12], *HIVb* [13], *HIVw* [13], *JTT* [14], *JonesDCMUT* [12], *LG* [15], *Mtart* [16], *Mtmam* [17], *Mtrev24* [18], *RtRev* [19], *VT* [20], *WAG* [21]. In addition, it is possible to specify a user-defined empirical amino acid model by an input file (see section 3.4.4).

## 4.2. SCS models

The SCS models take into account the stability of the protein structure. These models compute the structural energy of mutated proteins given the entropy, contact matrices, the protein structure and the temperature. Then, mutations are evaluated and can be accepted (substitutions) or rejected.

Our site-dependent SCS models estimates the stability of the mutated sequence folded into the target structure at the simulation temperature by means of a contact free energy function. The contact matrix  $C_{ij}$  takes the value 1 if residues  $i$  and  $j$  are close in space and 0 otherwise (by a threshold distance of 4.5Å), and it is sufficient to reconstruct the three-dimensional structure of the protein up to good accuracy [40]. We assume that the free energy of a protein with sequence  $A$  folded into the contact matrix  $C$  is given by the sum of its pairwise contact interactions,

$$E(A, C) = \sum_{ij} C_{ij} U(A_i, A_j)$$

where  $U(a, b)$  is the contact interaction matrix that expresses the free energy gained when amino acids  $a$  and  $b$  are brought in contact. We adopt the contact interaction matrix determined in Bastolla et al. [41]. For proteins that fold with two-states

thermodynamics, i.e. for which only the native structure and the unfolded structure are thermodynamically important, stability against unfolding is defined as the free energy difference between the folded and the unfolded states. The free energy of the unfolded state is estimated as  $sL$ , where  $L$  is protein length and  $s$  is an entropic parameter, is estimated as  $\Delta G \sim E(A, C_{nat}) + sL$ , where  $C_{nat}$  is the native structure and  $s = 0.074$  was determined fitting the above equation to a set of 20 experimentally measured unfolding free energy, yielding a correlation coefficient  $r = 0.92$  (U. Bastolla, unpublished data). The accuracy of this method for predicting the stability effect of mutations is comparable to state-of-the-art atomistic methods such as Fold-X [42], and its computational simplicity allows to use it for simulating protein evolution for long evolutionary trajectories and complex phylogenetic histories.

Stability against unfolding is however not sufficient to characterize protein stability. Consequently, the model has also to check the stability against compact, incorrectly folded conformations of low energy that can act as kinetic traps in the folding process and, in many cases, give raise to pathological aggregation. The term positive design indicates sequence features that favor protein stability by decreasing the free energy of the native structure. On the other hand, stability against misfolding is realized by increasing the energy of key contacts that are frequently found in alternative structures, which is termed negative design [43-45]. Therefore, it is also influenced by mutations at positions that are distant in the native structure.

Stability against misfolded structures is difficult to estimate, and several models of protein evolution do not consider it, despite its importance is being more and more recognized. Here we consider the set of alternative compact matrices of  $L$  residues that can be obtained from non-redundant structures in the Protein Data Bank. This procedure, called threading, guarantees that the contact matrices fulfill physical constraints on chain connectivity, atomic repulsion, and hydrogen bonding (secondary structure), which are not enforced in the contact energy function.

The free energy of this misfolded ensemble can be estimated in analogy with the Random Energy Model [REM, [46]] as,

$$G_{misfold} \approx \langle E(A, C) \rangle - \frac{\sigma^2}{2k_B T} - k_B T s_c L$$

where  $\langle E(A, C) \rangle$  is the mean and  $\sigma^2$  is the variance of the energy of alternative structures [47]. This formula holds for temperatures above the freezing temperature at which the entropy of the misfolding ensemble vanishes. At lower temperatures the free energy maintains the same frozen value [46]. A precise computation showed that the third moment of the energy can not be neglected (Minning et al. 2013). We therefore consider the equation,

$$\begin{aligned} G_{misfold} &\approx \langle E \rangle - \frac{\langle (E - \langle E \rangle)^2 \rangle}{2k_B T} + \frac{\langle (E - \langle E \rangle)^3 \rangle}{6(k_B T)^2} - k_B T s_c L \\ &= \sum_{i < j} \langle C_{ij} \rangle U_{ij} - \sum_{i < j < k < l} D_{ijkl} \frac{U_{ij} U_{kl}}{2T} + \sum_{i < j, k < l, m < n} F_{ijklmn} \frac{U_{ij} U_{kl} U_{mn}}{6T^2} - k_B T s_c L \end{aligned}$$

where the tensors  $\langle C_{ij} \rangle$ ,  $D_{ijkl}$  and  $F_{ijklmn}$  are the averages over the set of contact matrices of  $L$  residues respectively of single contacts, contact correlations and triples of contacts and,  $U_{ij} = U(A_i, A_j)$  only depends on the protein sequence [45]. The computing time is considerably reduced by approximating the above formula with one that only depends on pairs of residues,

$$G_{\text{misfold}} \approx A^{(1)} \langle U \rangle - \frac{A^{(2)} \langle U \rangle^2 + \sum_{ij} B_{ij}^{(1)} U_{ij}^2}{2k_B T} + \frac{A^{(3)} \langle U \rangle^3 + \langle U \rangle \sum_{ij} B_{ij}^{(2)} U_{ij}^2 + \sum_{ij} B_{ij}^{(3)} U_{ij}^3}{6(k_B T)^2} - k_B T s_c L$$

The quantities  $A^{(1)}$   $A^{(2)}$   $A^{(3)}$   $B_{ij}^{(1)}$   $B_{ij}^{(2)}$   $B_{ij}^{(3)}$  only depend on the set of alternative contact matrices and on protein length  $L$ , and they are pre-computed before the simulation starts. In this way, we can evaluate how the misfolded free energy changes upon mutation only performing order  $L$  operations for computing  $U(A_i = b, A_j) - U(A_i = a, A_j)$  when the residue at the mutated site  $i$  changes from state  $a$  to  $b$ . The stability of the native state is finally evaluated as the difference in free energy between the native, the unfolded and the misfolded states,  $\Delta G = E(A, C_{\text{nat}}) - G_{\text{misfold}} - k_B T s_u L$ .

Note that, even if the two configurational entropies per residue  $s_u$  (unfolded ensemble)  $s_c$  act additively, the free energy is not simply a function of their sum, since it is only  $s_c$  that determines the freezing temperature of the misfolded ensemble.

For modeling protein evolution, we still have to define how protein stability influences fitness. The simplest possibility is a neutral fitness landscape where the fitness is a binary variable and all proteins with stability above a given threshold, i.e.  $\Delta G < \Delta G_{\text{thr}}$  are considered viable and equally fit, whereas all proteins below threshold are considered lethal and therefore, discarded. The threshold is  $\Delta G_{\text{thr}} = \Delta G(A_0, C_{\text{nat}})$  where  $A_0$  is the protein sequence in the Protein Data Bank, which means that the neutral SCS model is not sensible to variations of the temperature or configurational entropies.

The neutral fitness landscape can be generalized to a landscape in which fitness is an increasing function of stability, and in particular it is proportional to the fraction of protein that is in the native state [47],

$$f(A) = \frac{1}{1 + e^{\Delta G(A, C_{\text{nat}}) / kT}}$$

Note that the fitness landscape can be reduced to the neutral landscape in the limit of very small temperature, since in this limit the fitness tends to 1 if  $\Delta G < 0$  and to zero if  $\Delta G > 0$ .

We then assume that the mutation rate is very small so that the population is monomorphic, and model selection through the Moran's birth-death process [48], which yields the fixation probability,

$$\Pi(ij) = \frac{1 - \left( \frac{f_i}{f_j} \right)^a}{1 - \left( \frac{f_i}{f_j} \right)^{2N}}$$

where  $f_i$  is the fitness of the wild-type,  $f_j$  is the fitness of the mutant,  $N$  is the effective



population size and  $a = 2$  or  $1$  for a haploid or diploid population, respectively. Given the probability of fixation, the succession of mutant fixations can be depicted as a Markov process, in which the genotype of the population moves from one sequence to another one according to the mutation and fixation probabilities.

#### **4.2.5 Crossing the SCS models with classic Markov substitution models**

The SCS models implemented in *ProteinEvolver* can be crossed with any parametric DNA substitution model (such as JC, HKY or GTR) or any empirical amino acid substitution model (such as JTT, WAG or LG). The matrices of change and the nucleotide or amino acid frequencies are considered to compute the rates of change per site and state.

Two types of material can be simulated using protein stability substitution models, coding DNA and amino acid sequences.

##### **4.2.5.1. Simulation of coding DNA data with the SCS models**

The simulation of coding DNA data consists basically in the following algorithm (which can be specified with argument “-m” in the parameters file).

I) mutation rates among DNA states and per site are computed according to the rates of change provided by the matrix of the classic Markov model (JC, HKY or GTR, parameters specified with the command “-v” in the parameters file) and the nucleotide frequencies “-f4”.

II) According to the mutation rates, a mutation at DNA level is introduced.

III) The DNA sequence is translated to an amino acid sequence and the mutation is classified as synonymous or non-synonymous.

IV) If the mutation is synonymous, it is accepted since it does not alter the protein sequence and structure.

V) If the mutation leads to a nonsynonymous change, it is evaluated by the computation of the structural protein energy for the mutated amino acid protein (section 4.2.2), the energy is used to compute a fitness value (section 4.2.3) and the mutation can be accepted or rejected according to the Moran’s birth-death process [48]. The mutation can be accepted or rejected and the DNA and protein sequences are updated according to such a result.

The number of mutations or substitutions (accepted mutations) is computed according to the branch length and the user can specify if such branch length should be considered as mutations or substitutions (see section 3.3.1).

##### **4.2.5.2. Simulation of protein data with the SCS models**

The simulation of protein data consists in the following algorithm (which can be specified with argument “-z” in the parameters file).

I) mutation rates among amino acid states and per site are computed according to the rates of change provided by an empirical amino acid matrix (WAG, JTT or HIVb, parameters specified with the command “-@” in the parameters file) and the nucleotide frequencies “-f20”.

II) According to the mutation rates, a mutation at amino acid level is introduced. Note that the mutation will be always classified as non-synonymous.

III) The mutation is evaluated by the computation of the structural protein energy for the mutated amino acid protein (section 4.2.2), the energy is used to compute a fitness value (section 4.2.3) and the mutation can be accepted or rejected [48]. The protein sequence is updated according to such a result.

The number of mutations or substitutions (accepted mutations) is computed according to the branch length and here the user can specify if such branch length should be again considered as mutations or substitutions (see section 3.3.1).

#### 4.2.6 Evolution of sequences along trees under the SCS models

While under classic Markov substitution models the evolution of the sequence is preformed site by site along the branches of the tree, because the site-independent aspect of these models, the simulation under SCS models requires the evolution of the whole molecule along the tree because the site-dependent aspect.

This leads to a problem when dealing with recombination because half a protein cannot be evolved independently. Consequently, the user-specified tree cannot be reticulated and recombination is only allowed in *ProteinEvolver*, for these structural substitution models, by a special algorithm described in 4.3.1.3. Here, the evolution starts from a molecule assigned to the GMRCA node of the ARG [31]. By this procedure all the material involved in the recombination events is considered.

#### 4.3. Coalescent simulations

The coalescent proceeds backwards starting from the sample of  $s$  gametes. Time is scaled in units of  $2N$  generations, where  $N$  is the effective population size. For a given site, without recombination or migration, and under constant population size, the time to the most recent common ancestor (TMRCA) is:

$$E(TMRCA) = 2 \left( 1 - \frac{1}{s} \right)$$

$$Var(TMRCA) = \sum_{i=2}^s \frac{4}{i^2(i-1)^2}$$

The times to a coalescence (CA), recombination (RE) or migration (MI) event are exponentially distributed, with intensity equal to their respective rates (see below). The next event will be the one that would occur before according to these expectations.

$$Time\ to\ CA \propto Exp[rateCA] \cdot 2N$$

$$Time\ to\ RE \propto Exp[rateRE] \cdot 2N$$

$$Time\ to\ MI \propto Exp[rateMI] \cdot 2N$$

The rate of coalescence depends only on the number of lineages ( $k$ ):

$$\text{rateCA} = \frac{k(k-1)}{2}$$

### 4.3.1. Recombination

The rate of recombination depends on the population size ( $N$ ) and on the total recombination rate at all valid recombination sites ( $G$ ). A *valid* recombination site has to have at both sides ancestral material that has not found its MRCA yet.

$$G = \sum_{j=1}^k \sum_{i=1}^L r_{Gi}$$

$$\rho = \text{rateRE} = 2NG$$

Given that a recombination event occurs, a gamete is chosen according to the total rate at potential recombining sites in that gamete. Breakpoint sites are chosen according to the recombination probabilities per site ( $r_{Gi}$ ). The expected number of recombination events in a panmictic population with constant size is:

$$E(\text{number of recombination events}) = \rho \sum_{i=1}^{s-1} \frac{1}{i}$$

Importantly, in the presence of recombination, different regions of the alignment might evolve under different genealogies (Figure 1), which together conform to the ancestral recombination graph. The number of genealogies will be the number of breakpoints + 1.

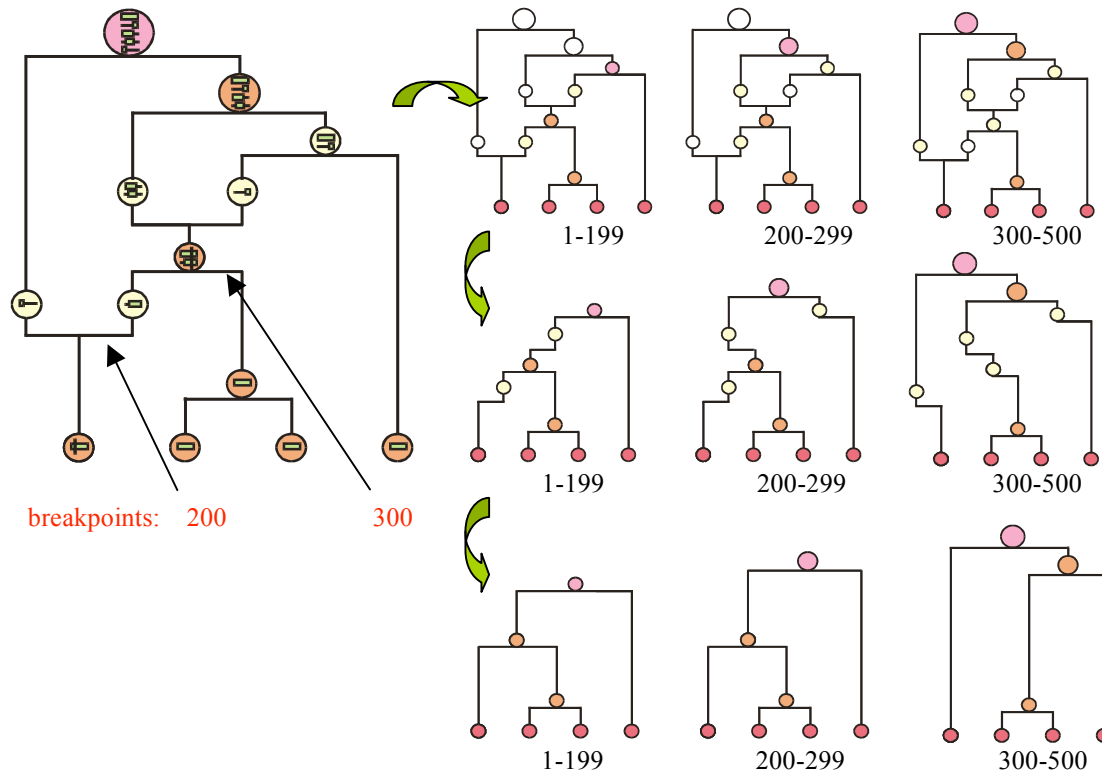


Figure 1. Representation of the ancestral recombination graph and the binary tree embedded.

#### 4.3.1.1 Recombination hotspots

This implementation is based on *SNPsim* [24], the population genetic model of recombination hotspots developed by Wiuf and Posada, to result in heterogeneous recombination rates along the chromosomes.

Given an expected number of recombination hotspots, a background homogenous recombination rate and a hotspot recombination rate, the algorithm starts by choosing the position and number of recombination hotspots for a particular sample. Adding recombination events around the hotspot center results in the specification of a probability distribution for the recombination rate along the region of interest. Other recombination events coming from recombinational hotspots centered outside the region of interest are also considered. This fast simulation results in different recombination rates for different sites along the region (hotspots and coldspots). Next model for hotspot recombination is described following the algorithm [24].

#### Hotspot recombination model

The hotspots recombination model aims to represent the idea that some sites in a chromosome are more likely to recombine than others (“recombination hotspots”). This general model is composed of two basic recombination rates and a set of hotspots sites. The background recombination rate ( $R_B$ ) is the per generation recombination rate at any site in the chromosome of length  $L$ , while the hotspot recombination rate ( $R_H$ ) is the additional per generation recombination rate at the recombination hotspot.  $X$  is the number of hotspots.

$$\text{Background recombination rate } (R_B) = 4Nr_B L$$

$$\text{Hotspot recombination rate } (R_H) = 4Nr_H X$$

$$\text{Global recombination rate } (R_G) = R_B + R_H$$

In real life we do not expect the recombination hotspot to be always restricted to the same single site, to have always the same intensity, or to occur independently from other hotspots. The model described above can be generalized to include these relevant biological features. When the hotspot is not restricted to a single site, and under constant population size, the expected number of recombination events is

$$E(\text{number of recombination events}) = R_G \sum_{i=1}^{s-1} \frac{1}{i}$$

#### Hotspot location: interference

We will assume that the distance between hotspots is Gamma distributed,  $\Gamma(m, \lambda)$ ,  $m > 0$ . If  $m=1$ , we have a Poisson process with intensity  $\lambda$ . Allowing  $m \neq 1$  introduces interference. If  $m > 1$  hotspots are pushed away from each other; if  $0 < m < 1$  they tend to be clustered (although the algorithm will only accept integer values for  $m$ ). The average number of hotspots in the gene is  $\lambda/m$  (see Figure 2).

*Hotspot imprecision*

We will consider that the hotspot location actually represents the center of the hotspot, and where recombination is more likely, but that there are also some sites around where recombination occurs with some frequency that decreases with increasing distance from the hotspot center. This can be represented by a Normal or a Uniform distribution for the location of the recombination. When the Normal distribution is used, this has a mean equal to the location of the hotspot center and variance called hotspot imprecision ( $\sigma^2$ ) (see Figure 2). When the variance is small the hotspot tends to be narrow. If the variance is 0, the hotspot is 1 bp wide. When the Uniform distribution is used, recombination events occur with the same probability along a given width for the hotspot. In addition there is the possibility of recombination events coming from hotspots located outside the region of interest of length  $L$ . To implement this idea we can extend the region of interest by a number of sites  $K$  at each end. In the Normal distribution an arbitrary, but seemingly reasonable value for  $K$  that assures that wide hotspots outside  $L$  are taking into account is

$$K = 10\sqrt{\sigma^2}$$

If the uniform distribution is used,  $K$  equals half the width of the hotspot.

*Hotspot heterogeneity*

In addition, not all hotspots have to be equally “hot”. We can model this heterogeneity of recombination rates using a gamma distribution with a mean of 1. The shape of this distribution ( $\alpha$ ) will determine the strength of this hotspot heterogeneity. The smaller  $\alpha$ , the bigger the heterogeneity of recombination rates at the hotspots.

*Implementation*

A nice feature of the hotspot model is that it allows for the construction of a distribution of recombination rate along the chromosome ( $\mathfrak{R}$ )(Figure 2) for every sample.

The algorithm starts by building  $\mathfrak{R}$ . The first step is to set up the number of hotspot centers ( $X$ ) along the extended region  $L + 2K$ . The average number of hotspots in the gene is  $\lambda/m$ , and when  $m > 1$  we use a thinning algorithm to locate these hotspots (Figure 2). If  $m=1$  we distributed the hotspots according to a Poisson distribution with intensity  $\lambda$ .

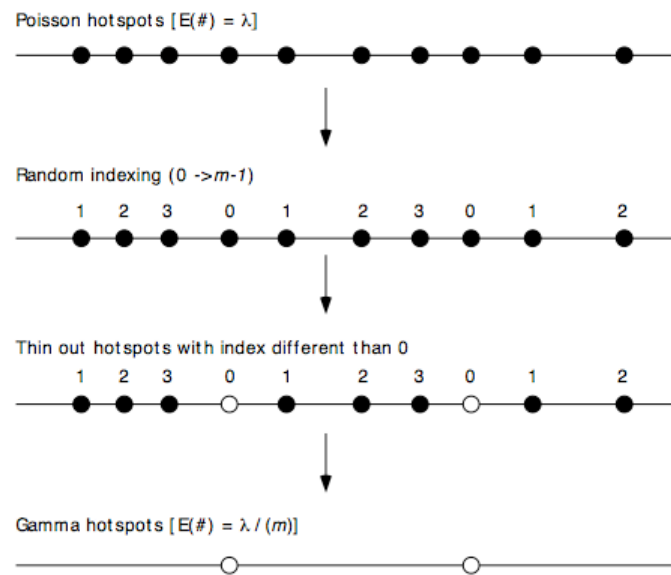


Figure 2. Thinning algorithm for the location of hotspots with interference ( $m=4$ ). The index for the first hotspot site is chosen uniformly between 0 and  $m-1$ . This plot was taken from SNPsim manual [24].

Alternatively the user can specify a fixed number of hotspots, which will be located uniformly (see option -q in recombination hotspots file).

The distribution  $\mathfrak{R}$  is can be constructed according to a Normal ( $x_i, \sigma^2$ ), or a Uniform (hotspot width) distribution. If there is hotspot recombination, a random gamma variable will scale the recombination events at each hotspot. Figures 3 and 4 represent a realization of this process.

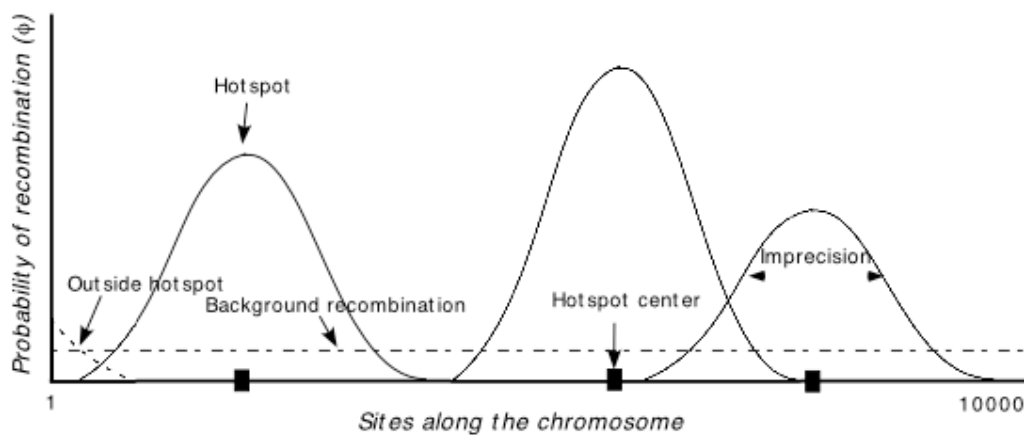


Figure 3. Schematic representation of  $\mathfrak{R}$  in the case of three Normal hotspots for the region of interest ( $L=10000$ ). In this case the hotspot imprecision is quite big. Note that some recombination probability is contributed by a hotspot outside the region of interest. The background recombination is the same for all sites. In this case there is hotspot heterogeneity. This plot was taken from SNPsim manual [24].

We can use now  $\mathfrak{R}$  now to set the global recombination rate per each site  $i$ ,

$$r_{Gi} = r_{Bi} + r_{Hi}\varphi_i$$

where  $\varphi_i$  is the probability of block recombination at site  $i$ .

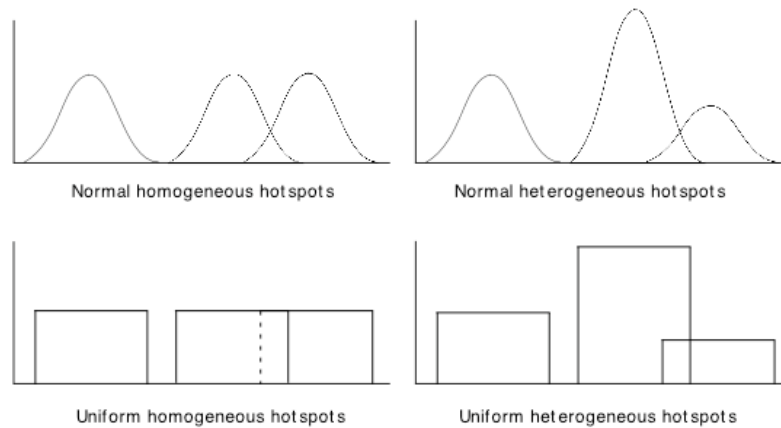


Figure 4. Schematic representation of different  $\mathfrak{R}$ . This plot was taken from *SNPsim manual* [24].

Finally, the total recombination rate at all valid recombination sites is considered to compute times for recombination events (as described before). Breakpoint sites are chosen according to the recombination probabilities per site ( $r_{Gi}$ ).

#### 4.3.1.2 Visualization and characterization of the simulated ARG

The simulated ARG can be directly exported to the *NetTest* web server (<http://darwin.uvigo.es/software/nettest/>) [36] in order to visualize and characterize the ARG [36, 49]. The ARG can be printed by *ProteinEvolver* in branches format [36] (argument “\*” in the parameters file) to be introduced in *NetTest*.

#### 4.3.1.3 Algorithm to evolve molecules with the SCS models along the ARG

In the presence of recombination, classic Markov substitution models of evolution independently evolve each site along the evolutionary history of the corresponding recombinant fragment [see for example 30, 39, 50, 51]. Nevertheless, the SCS models require the evolution of the whole molecule because the dependence among sites. This implied the development of a new algorithm to evolve whole molecules under recombination.

After building the ARG, the molecular evolution starts by the assignation of a sequence to the GMRCa node. *ProteinEvolver* incorporates an algorithm for the evolution of the whole molecule along the ARG (see Figure 5), which is an adaptation of the algorithm developed to by Arenas and Posada to simulate intracodon recombination [30]. The process starts from the sequence assigned to the GMRCa node. This sequence can suffer substitutions along the branches and new whole sequences are assigned to the descendant nodes (see steps 1-2 in the Figure 5). Later the evolution process can reach a

recombinant node (shown in grey, see step 3) and a sequence is assigned to such a recombinant node. However, at that point the evolution stops and comes back to continue by other path. This is forced to occur since there is no information about the sequence at the parental recombinant node. Later, the process reaches the parental recombinant node (step 5), at this point note that there are whole sequences in both recombinant nodes, and there is a combination of the material according to the recombinant breakpoint. Consequently, a new whole sequence is generated (step 6) and this new sequence continues the evolutionary process. By this algorithm the molecular evolution considers the exchanges of material introduced by the recombination and that are described in the ARG.

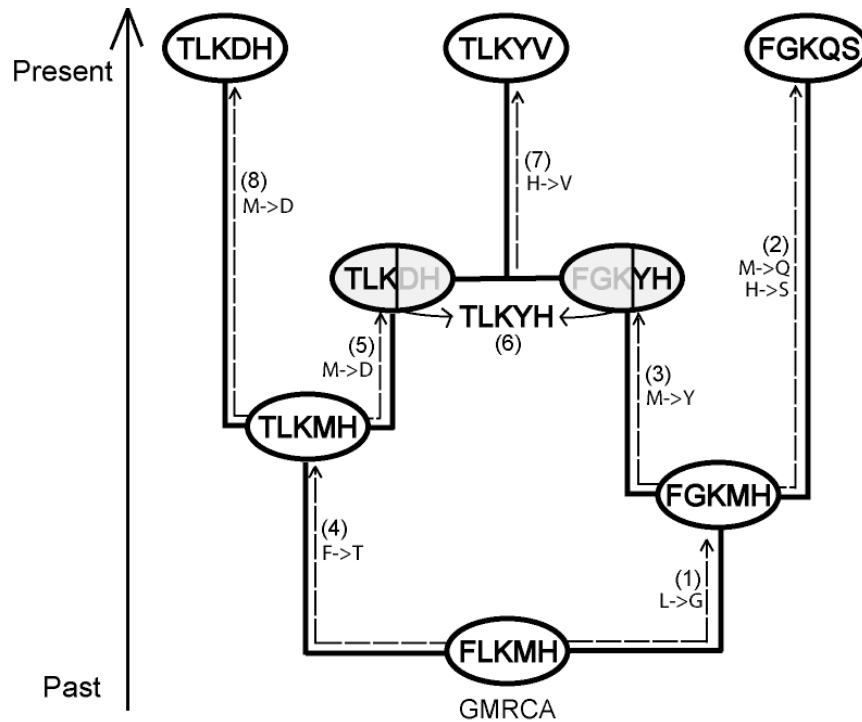


Figure 5. Example of the evolution of a whole molecule along an ARG.

### 4.3.2 Migration

The simplest and most widely used model of population structure is the “finite island model” [25, 27, 52-54] however the stepping-stone [26] and an approximation to the continent-island [27] are also implemented. The rate of migration depends on the population size ( $N$ ), the migration rate per deme ( $m$ ) and the number of available demes ( $q$ ).

$$rateMI = 2Nm q$$

Importantly, in the coalescent with migration, lineages can only coalesce with other lineages in the same deme. When a migration event occurs, a given lineage changes from one deme to another (Figure 6).



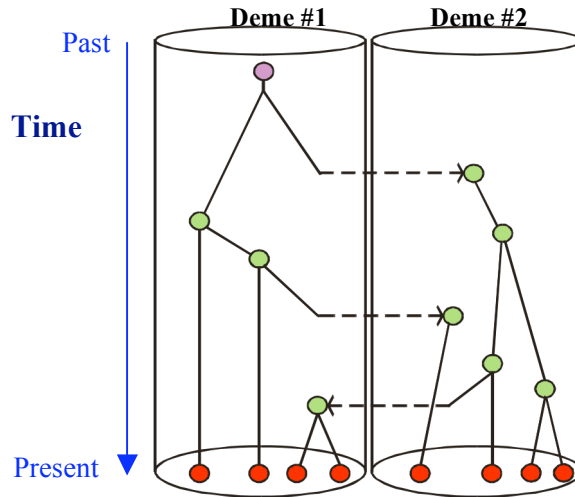


Figure 6. Coalescent with migration for two demes.

#### 4.3.2.1 Island model

The Island Model describes an array of  $q$  demes, each of constant size (Figure 7). Migration events may occur between any two demes.

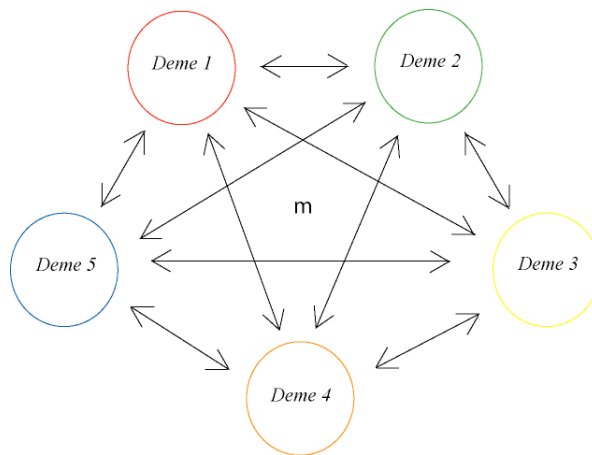


Figure 7. Island model. All demes have the same population size and the same migrant rate per deme.

#### 4.3.2.2 Stepping-stone model

Under the stepping-stone migration events occur between neighboring demes (Figure 8).



Figure 8. 1D stepping-stone model implemented in ProteinEvolver.

#### 4.3.2.3 Continent-island model

The continent-island model implemented in *ProteinEvolver* allows symmetrical migration between a deme (continent) and the other demes (islands), see Figure 9. Note that there is not possible to directly migrate individuals among islands.

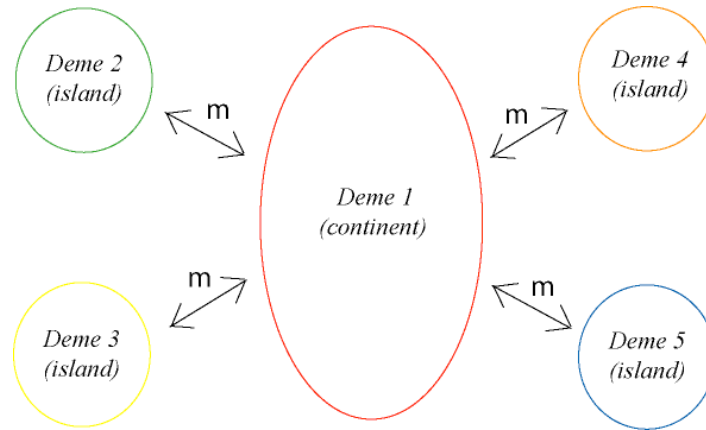


Figure 9. Continent-island model implemented in *ProteinEvolver*.

#### 4.3.2.4 Temporal variation of the migration rate

*ProteinEvolver* allows for temporal variation where a migration rate can be applied at different time periods (Figure 10).

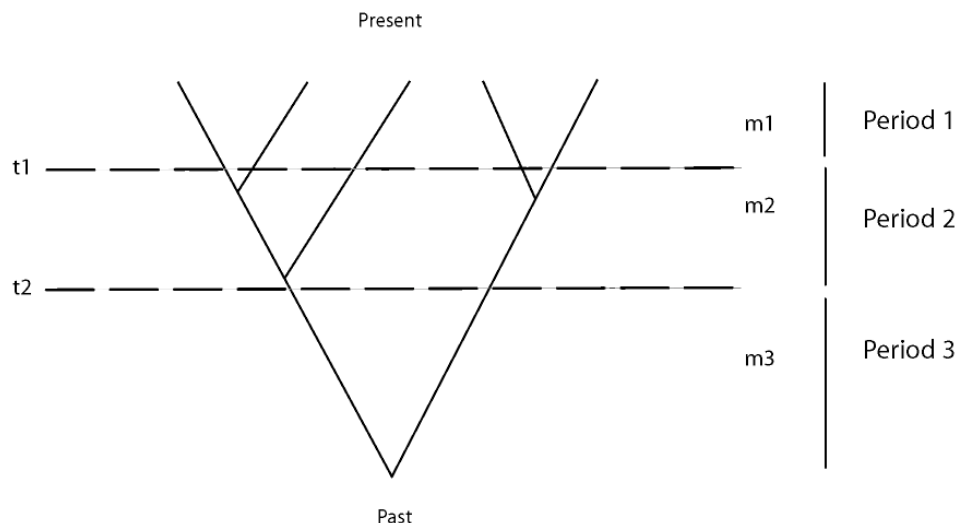


Figure 10. Temporal variation of the migration rate implemented in *ProteinEvolver*. The figure shows 3 temporal periods defined by two times ( $t1$  and  $t2$ ) and where each period implements a given migration rate ( $m1$ ,  $m2$  or  $m3$ ).

### 4.3.3 Convergence of demes

The evolution of the demes can be user-specified. Two demes can be converged at a given time in order to build a new big deme with the lineages of the previous demes. The convergence of demes always occurs under a migration model. An example is shown in the Figure 11. It is not possible use convergence of demes in presence of tip dates when the time of the convergence of demes is younger that the time of the tip dates.

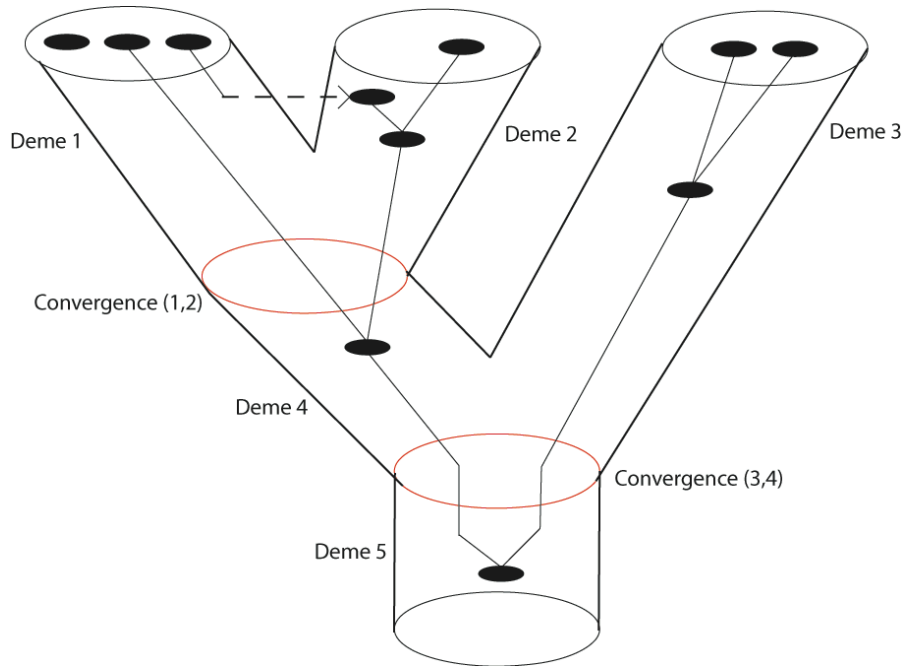


Figure 11. Demes evolution. Two demes can evolve convergences to generate an ancestral big deme, these convergences are introduced by the user.

### 4.3.4 Demography

*ProteinEvolver* allows for the specification of exponential population growth rate. The only modification concerns to the expected time to a coalescence, where  $t$  is the current time:

$$time\ to\ CA \sim \frac{\log \left[ e^{\beta t} + \beta \text{Exp} \left[ \frac{k(k-1)}{2} \right] 2N \right]}{\beta} - t$$

If the growth rate is negative, the coalescence time may be infinite (i.e., coalescence does not happen), and *ProteinEvolver* will stop and issue an error message. Alternatively, the user can define any number of demographic periods (Figure 12).  $\beta_i$  is the growth rate inferred for a demographic period  $i$  that goes from size  $N_{Bi}$  in the past to size  $N_{Ei}$  in  $I_i$  generations:

$$\beta_i = \frac{-\log\left(\frac{N_{Bi}}{N_{Ei}}\right)}{l_i}$$

The time to coalescence will be:

$$\text{Time to CA} \sim \frac{\log\left[\text{Exp}\left(\frac{k(k-1)}{2}\right)\beta_i 2N_{Ei}e^{-\beta(t-t_{i-1})} + 1\right]}{\beta_i}$$

where  $t$  is the current time and  $t_i$  is the cumulative time from the present:

$$t_i = \sum_{j=0}^{j=i} l_j$$

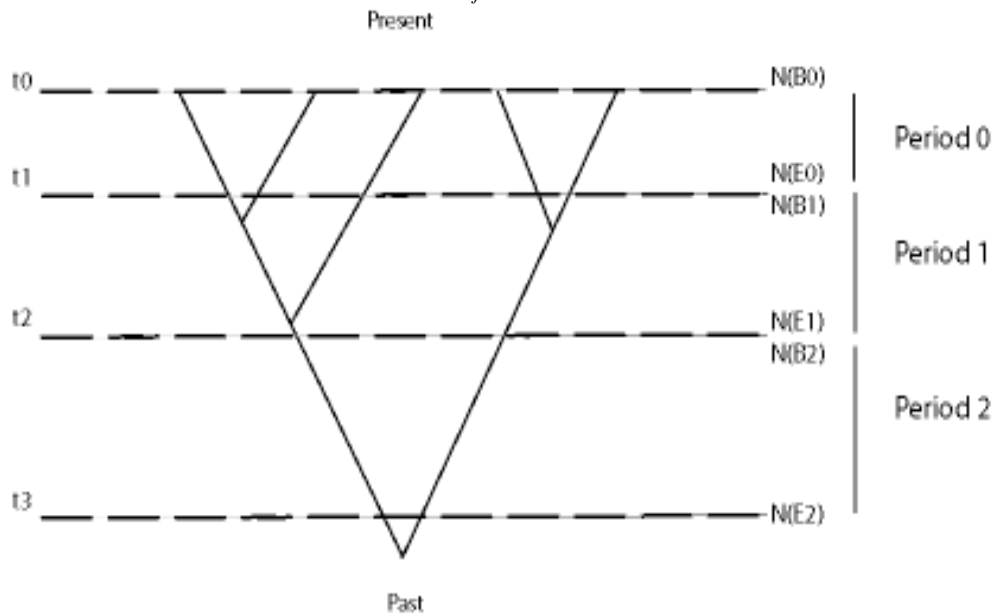


Figure 12. Demographic periods. The growth rate after the last period will be the same as the one implied by the last period.

#### 4.3.5 Tip Dates

Tip dates can be specified by the user to simulate sequences with different times, such as samples taken at different times. For this option, the user must introduce a generation time. Figure 13 shows an example of four samples taken at different times, the oldest sample contains only one sequence (seq00001), and the younger sample contains three sequences (seq00004, seq00006 and seq00005).

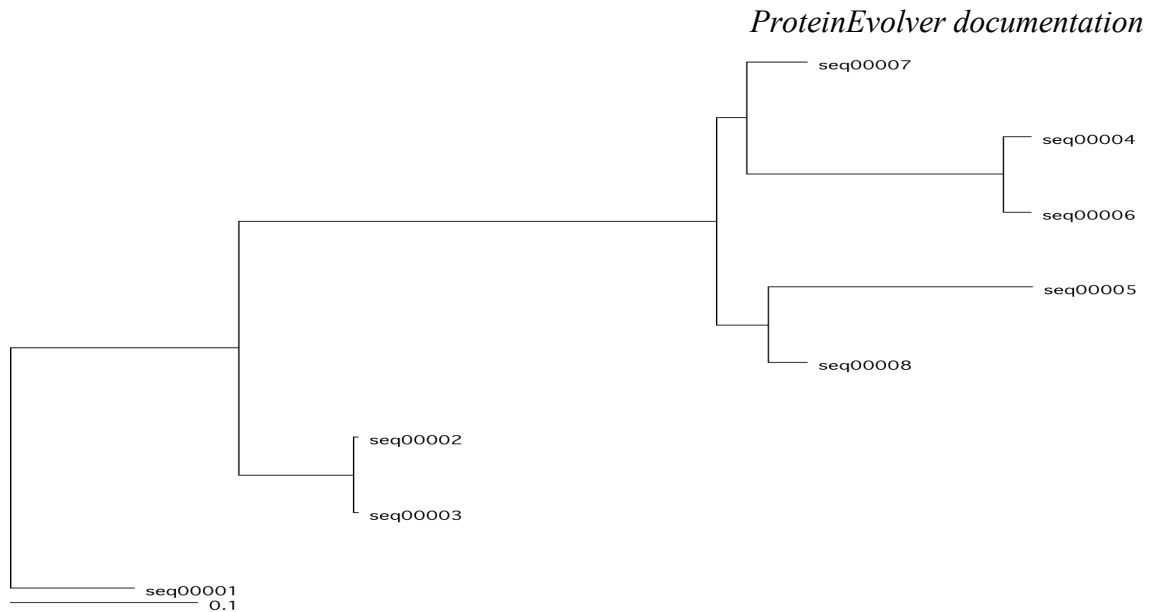


Figure 13. Tip dates simulation examples. Three examples of four samples taken at different times with 1, 2, 2 and 3 sequences per sample (from present to past).

## 5. Program benchmarking

The models on which *ProteinEvolver* is based have been separately validated. The method to estimate protein folding stability against unfolding and misfolding has been developed in [45, 55, 56] and yields good correlations with experimentally measured free energies. It has been widely used with some variations in many simulations of protein evolution (Bastolla 2001; etc.), since its computational simplicity allows its application over long evolutionary trajectories.

The substitution processes implemented in the SCS models were verified using HYPHY [57] and agreed with the branch lengths. The implemented site-independent DNA substitution models were previously validated [30, 39]. The empirical amino acid substitution models were accurately estimated using *ProtTest* [58]. Coalescent simulations were contrasted with the theoretical expectations for the mean and variances for different values, like the number of recombination and migration events, or the times to the most recent common ancestor [39, 59]. Recombination hotspots agree with theoretical expectations and breakpoints were likely to occur in recombination hotspots regions [24].

## 6. History

*Version 1.0.0 (June 2012)*

- Implementation of *Pop\_evol* (code developed by Ugo Bastolla) in *ProteinEvolver*, a mutational model to evolve DNA sequences accounting for structural protein stability (DNA → AA → DNA) and the HKY substitution model.

*Version 1.0.1 (July-August 2012)*

- PDB from input file.
- Seed for structural protein stability substitution model comes now from the main input file (parameters).

*Version 1.0.2 (August 2012)*

- Structural protein stability accounting for rates of change among all nucleotide states (GTR and GTnR).

*Version 1.0.3 (September 2012)*

- Solved bug in the "Compute\_load" function. Other warnings from Pop\_Evol were solved.

- Solved bug in GTR and GTnR models for PopEvol section (in "Mutate\_nuc" and "Compute\_load" functions).

- Structural protein stability substitution model crossed with any empirical substitution model (Blosum62, CpRev, Dayhoff, DayhoffDCMUT, HIVb, HIVw, JTT, JonesDCMUT, LG, Mtart, Mtmam, Mtrev24, RtRev, VT, WAG, UserEAAM).

- Validation. User-specified tree was inferred by Hyphy from the simulated alignmets (for both DNA and amino acid data) under the structural protein stability substitution model.

*Version 1.0.4 (September 2012)*

- Coalescent with recombination is allowed for the simulation under structural protein stability substitution models.

When a recombination occurs, the evolution stops at the first parental recombinant node and goes in other direction. Later, when the other parental recombinant node is reached, there is a combination of the material, according to the recombination breakpoint, to generate the material for the descendant node, which continues the evolution. Recombination can be homogeneous (recombination rate is constant among all sites) or heterogeneous (recombination hotspots, functions developed by David Posada). However, recombination cannot be applied using user-specified input trees and structural protein stability substitution models because part of a protein cannot be evolved under such a substitution model.

*Version 1.0.5 (September 2012)*

- Included options to print the energy output files per branch produced by the structural protein stability substitution models.

- Print information about the attempted mutations in the structural protein stability substitution models.

*Version 1.2.0 (September-October 2012)*

- Proportion of invariable sites +I and variable substitution rates per site (heterogeneity) according to a gamma distribution +G, for the structural protein stability substitution models.

- Heterogeneity by a user-defined vector. Each site can evolve under a particular rate user-defined. This is compatible with the specification of +G.

- Print number and mean of the introduced substitution events and, those nonsynonymous for DNA structural protein evolution models.

- Maximun number of characters for the name of taxas is 10 (Phylip format).

- More robust generation of the random seed.

- Work with bigger input trees (MAX\_LINE).

*Version 1.2.0 (September-December 2012)*

- Implementation of the neutral evolution model to evolve sequences accounting for the protein structure. Here population size is not required, follow this with the variable "NEUTRAL".

## 7. Acknowledgments

This work was supported by the “Juan de la Cierva” fellowship JCI-2011-10452 to MA from the Spanish Government. Several functions were taken from code provided by R. Nielsen and Z. Yang, other functions were provided from *SNPsim* [24] (by David Posada and Carsten Wiuf) and *Mosaic* (by David Posada). We want to thank J. Carlos Mouriño at the Supercomputing Center of Galicia (CESGA) for extensive help with code parallelization in first versions of the coalescent software.

## 8. References

1. Mendez R, Fritsche M, Porto M, Bastolla U: **Mutation bias favors protein folding stability in the evolution of small populations.** *PLoS Comput Biol* 2010, **6**(5):e1000767.
2. Bastolla U, Porto M, Roman HE, Vendruscolo M: **Structural approaches to sequence evolution.** Berlin, Heidelberg: Springer; 2007.
3. Jukes TH, Cantor CR: **Evolution of protein molecules.** In: *Mammalian Protein Metabolism*. Edited by Munro HM. New York, NY: Academic Press; 1969: 21-132.
4. Kimura M: **A simple method for estimating evolutionary rate of base substitutions through comparative studies of nucleotide sequences.** *J Mol Evol* 1980, **16**:111-120.
5. Felsenstein J: **Evolutionary trees from DNA sequences: A maximum likelihood approach.** *J Mol Evol* 1981, **17**:368-376.
6. Hasegawa M, Kishino K, Yano T: **Dating the human-ape splitting by a molecular clock of mitochondrial DNA.** *J Mol Evol* 1985, **22**:160-174.
7. Zharkikh A: **Estimation of evolutionary distances between nucleotide sequences.** *J Mol Evol* 1994, **39**(3):315-329.
8. Tavaré S: **Some probabilistic and statistical problems in the analysis of DNA sequences.** In: *Some mathematical questions in biology - DNA sequence analysis*. Edited by Miura RM, vol. 17. Providence, RI: Amer. Math. Soc.; 1986: 57-86.
9. Henikoff S, Henikoff JG: **Amino acid substitution matrices from protein blocks.** *Proc Natl Acad Sci U S A* 1992, **89**(22):10915-10919.
10. Adachi J, Waddell PJ, Martin W, Hasegawa M: **Plastid genome phylogeny and a model of amino acid substitution for proteins encoded by chloroplast DNA.** *J Mol Evol* 2000, **50**(4):348-358.
11. Dayhoff MO, Schwartz RM, Orcutt BC: **A model of evolutionary change in proteins.** In: *Atlas of protein sequence and structure*. Edited by Dayhoff MO, vol. 5, Suppl. 3. Washington D. C.; 1978: 345-352.
12. Kosiol C, Goldman N: **Different versions of the Dayhoff rate matrix.** *Mol Biol Evol* 2005, **22**(2):193-199.
13. Nickle DC, Heath L, Jensen MA, Gilbert PB, Mullins JI, Kosakovsky Pond SL: **HIV-specific probabilistic models of protein evolution.** *PLoS One* 2007, **2**(6):e503.
14. Jones DT, Taylor WR, Thornton JM: **The rapid generation of mutation data matrices from protein sequences.** *Comput Appl Biosci* 1992, **8**(3):275-282.

15. Le SQ, Gascuel O: **An improved general amino acid replacement matrix.** *Mol Biol Evol* 2008, **25**(7):1307-1320.
16. Abascal F, Posada D, Zardoya R: **MtArt: A New Model of Amino Acid Replacement for Arthropoda.** *Mol Biol Evol* 2007, **24**(1):1-5.
17. Yang Z, Nielsen R, Masami H: **Models of amino acid substitution and applications to mitochondrial protein evolution.** *Mol Biol Evol* 1998, **15**(12):1600-1611.
18. Adachi J, Hasegawa M: **MOLPHY version 2.3: programs for molecular phylogenetics based in maximum likelihood.** *Comput Sci Monogr* 1996, **28**:1-150.
19. Dimmic MW, Rest JS, Mindell DP, Goldstein RA: **rtREV: an amino acid substitution matrix for inference of retrovirus and reverse transcriptase phylogeny.** *J Mol Evol* 2002, **55**(1):65-73.
20. Muller T, Vingron M: **Modeling amino acid replacement.** *J Comput Biol* 2000, **7**(6):761-776.
21. Whelan S, Goldman N: **A general empirical model of protein evolution derived from multiple protein families using a maximum-likelihood approach.** *Mol Biol Evol* 2001, **18**(5):691-699.
22. Hudson RR: **Properties of a neutral allele model with intragenic recombination.** *Theor Popul Biol* 1983, **23**:183-201.
23. Kingman JFC: **The coalescent.** *Stochastic Processes and their Applications* 1982, **13**:235-248.
24. Wiuf C, Posada D: **A coalescent model of recombination hotspots.** *Genetics* 2003, **164**(1):407-417.
25. Hudson RR: **Island models and the coalescent process.** *Mol Ecol* 1998, **7**:413-418.
26. Kimura M, Weiss GH: **The Stepping Stone Model of Population Structure and the Decrease of Genetic Correlation with Distance.** *Genetics* 1964, **49**(4):561-576.
27. Wright S: **Evolution in Mendelian populations.** *Genetics* 1931, **16**:97-159.
28. Navascues M, Depaulis F, Emerson BC: **Combining contemporary and ancient DNA in population genetic and phylogeographical studies.** *Mol Ecol Resour* 2010, **10**(5):760-772.
29. Arenas M, Posada D: **The effect of recombination on the reconstruction of ancestral sequences.** *Genetics* 2010, **184**(4):1133-1139.
30. Arenas M, Posada D: **Coalescent simulation of intracodon recombination.** *Genetics* 2010, **184**(2):429-437.
31. Griffiths RC, Marjoram P: **An ancestral recombination graph.** In: *Progress in population genetics and human evolution.* Edited by Donnelly P, Tavaré S, vol. 87. Berlin: Springer-Verlag; 1997: 257-270.
32. Beaumont MA: **Approximate Bayesian Computation in Evolution and Ecology.** *Annu Rev Ecol Evol Syst* 2010, **41**:379-405.
33. Beaumont MA, Zhang W, Balding DJ: **Approximate Bayesian computation in population genetics.** *Genetics* 2002, **162**(4):2025-2035.
34. Hudson RR: **Generating samples under a Wright-Fisher neutral model of genetic variation.** *Bioinformatics* 2002, **18**(2):337-338.
35. Yang Z: **Among-site rate variation and its impact on phylogenetic analysis.** *Trends Ecol Evol* 1996, **11**(9):367-372.
36. Arenas M, Patricio M, Posada D, Valiente G: **Characterization of phylogenetic networks with NetTest.** *BMC Bioinformatics* 2010, **11**(1):268.



37. Rambaut A, Grassly NC: **Seq-Gen: an application for the Monte Carlo simulation of DNA sequence evolution along phylogenetic trees.** *Comput Appl Biosciences* 1997, **13**(3):235-238.
38. Yang Z: **PAML: a program package for phylogenetic analysis by maximum likelihood.** *Comput Appl Biosciences* 1997, **13**(5):555-556.
39. Arenas M, Posada D: **Recodon: coalescent simulation of coding DNA sequences with recombination, migration and demography.** *BMC Bioinformatics* 2007, **8**:458.
40. Vendruscolo M, Kussell E, Domany E: **Recovery of protein structure from contact maps.** *Fold Des* 1997, **2**(5):295-306.
41. Bastolla U, Roman HE, Vendruscolo M: **Neutral evolution of model proteins: diffusion in sequence space and overdispersion.** *J Theor Biol* 1999, **200**(1):49-64.
42. Guerois R, Nielsen JE, Serrano L: **Predicting changes in the stability of proteins and protein complexes: a study of more than 1000 mutations.** *J Mol Biol* 2002, **320**(2):369-387.
43. Berezovsky IN, Zeldovich KB, Shakhnovich EI: **Positive and negative design in stability and thermal adaptation of natural proteins.** *PLoS Comput Biol* 2007, **3**(3):e52.
44. Noivirt-Brik O, Horovitz A, Unger R: **Trade-off between positive and negative design of protein stability: from lattice models to real proteins.** *PLoS Comput Biol* 2009, **5**(12):e1000592.
45. Minning J, Porto M, Bastolla U: **Detecting selection for negative design in proteins through an improved model of the misfolded state.** *Proteins* 2013.
46. Derrida B: **Random Energy Model: An exactly solvable model of disordered systems.** *Phys Rev B* 1981, **24**:2613-2626.
47. Goldstein RA: **The evolution and evolutionary consequences of marginal thermostability in proteins.** *Proteins* 2011, **79**(5):1396-1407.
48. Ewens WJ: **Mathematical Population Genetics**, vol. 9. Berlin: Springer-Verlag; 1979.
49. Arenas M, Valiente G, Posada D: **Characterization of reticulate networks based on the coalescent with recombination.** *Mol Biol Evol* 2008, **25**(12):2517-2520.
50. Arenas M: **Simulation of Molecular Data under Diverse Evolutionary Scenarios.** *PLoS Comput Biol* 2012, **8**(5):e1002495.
51. Fletcher W, Yang Z: **INDELible: a flexible simulator of biological sequence evolution.** *Mol Biol Evol* 2009, **26**(8):1879-1888.
52. Latter BD: **The island model of population differentiation: a general solution.** *Genetics* 1973, **73**(1):147-157.
53. Maruyama T: **A simple proof that certain quantities are independent of the geographical structure of population.** *Theor Popul Biol* 1974, **5**(2):148-154.
54. Matsen FA, Wakeley J: **Convergence to the island-model coalescent process in populations with restricted migration.** *Genetics* 2006, **172**(1):701-708.
55. Bastolla U, Farwer J, Knapp EW, Vendruscolo M: **How to guarantee optimal stability for most representative structures in the Protein Data Bank.** *Proteins* 2001, **44**(2):79-96.
56. Bastolla U, Demetrius L: **Stability constraints and protein evolution: the role of chain length, composition and disulfide bonds.** *Protein Eng Des Sel* 2005, **18**(9):405-415.

57. Kosakovsky Pond SL, Frost SD, Muse SV: **HYPHY: Hypothesis testing using phylogenies**. *Bioinformatics* 2005, **21**:676-679.
58. Abascal F, Zardoya R, Posada D: **ProtTest: selection of best-fit models of protein evolution**. *Bioinformatics* 2005, **21**(9):2104-2105.
59. Hudson RR: **Gene genealogies and the coalescent process**. *Oxf Surv Evol Biol* 1990, **7**:1-44.