

Parameters for simulations using **retrocombinator**

Anindya Sharma

2021-03-23

Parameters and their defaults

- **SequenceParams** represents what set of retrotransposons the simulation starts off with. It is constructed from one of the following:
 - **num_seq** : **numeric** Number of initial retrotransposons (**default = 100**)
 - **seq_length** : **numeric** Sequence length (number of nucleotides) of the retrotransposons (**default = 5000**)
- **MutationParams** represents how what nucleotide substitution model will be used to modify the sequences during the simulation. It is one of the following **character** literals:
 - "JC69" Jules and Cantor - 1969
 - "K80" Kimura - 1980 (**default**)
 - "F81" Felsenstein - 1981
 - "HKY85" Hasegawa, Kishino and Yano - 1985
 - "TN93" Timura and Nei - 1993
 - "GTR" General Time Reversible Model, Tavaré - 1986
- **FlagParams** represents the parameters used during a simulation in which we keep track of a retrotransposon's 'active' status. A retrotransposon that is active is capable of transposition or bursting, but potentially loses its active status if a point mutation affects a critical part of the sequence.
 - **length_critical_region**: **numeric** Sequence length (number of nucleotides) of the critical region of a retrotransposon (**default = 10**)
 - **prob_inactivity_when_mutated** : **numeric** The probability that a point mutation to the critical region causes a sequence to become inactive (**default = 0.001**)
 - **max_inactive_copies** : **numeric** What is the largest population size of active sequences to keep track of? (**default = 500**)
- **BurstParams** represents how an active transposon will burst during transposition. It comprises of the following:
 - **burst_probability** : **numeric** The probability that an active retrotransposon will increase in copy number during a time jump of one timestep (**default = 0.1**)
 - **burst_mean** : **numeric** The Poisson mean for the distribution that specifies how many new sequences an active sequence will create during bursting (**default = 1**)
 - **recomb_mean** : **numeric** The expected number of template switches during recombination between two sequences (chosen from a Poisson distribution with this as its mean) (**default = 0.5**)
 - **recomb_similarity** : **numeric** How similar does an active sequence have to be with another sequence for them to be allowed to be recombine during transposition? (**default = 0.85**)
- **SimulationParams** represents how long the simulation will run for, and at what timescale. It comprises of the following:
 - **num_jumps** : **numeric** The number of steps (jumps) in our simulation (**default = 20**)
 - **timestep** : **numeric** How much real time does one step (jump) in our simulation measure, in millions of years (**default = 1**)
 - **max_active_copies** : **numeric** What is the largest population size of active sequences to keep track of? (**default = 500**)
- **SpeciationParams** represents how we keep track of families during the simulation. It comprises of the

following:

- **selection_threshold** : **numeric** What sequence similarity to the original sequence do we wish to maintain? Sequences that diverge beyond this similarity level are dropped over the course of simulation (**default = 0.5**)
- **family_similarity** : **numeric** What sequence similarity do two sequences have to be to each other for them to be considered to be of the same family? (**default = 0.7**)
- **family_coherence** : **numeric** What proportion of the overall sequence similarity matrix of a family needs to score be above **family_similarity** before we decide the family has split into further sub-families? (**default = 0.5**)
- **OutputParams** represents how and where the output of the simulation will be saved. It comprises of the following:
 - **file_out** : **character** Where should the simulation be saved? (**default = ‘file.out’**)
 - **num_out_seqs** : **numeric** How many times during the event of the simulation should we output the raw sequences (as nucleotide strings) themselves? (**default = 2**)
 - **num_out_init** : **numeric** How many times during the event of the simulation should we output the distance of sequences to the initial sequence? (**default = 10**)
 - **num_out_pair** : **numeric** How many times during the event of the simulation should we output pairwise distances between all pairs of sequences? (**default = 10**)
 - **num_out_fam** : **numeric** How many times during the event of the simulation should we output the family tags of each of the sequences? (**default = 10**)
- **SeedParams** represents how to select the seed for randomisation for the simulation. It comprises of the following:
 - **to_randomise** : **logical** Should this simulation be run with a random seed to begin with? (The seed is based on system time) (**default = false**)
 - **to_seed** : **logical** Should this simulation be run with a specified seed to begin with? (This overwrites **to_randomise**) (**default = true**)
 - **seed** : **numeric** If **to_seed** is TRUE, what should the initial seed for the random number generator be? (**default = 0**)

How to use

To run a simulation, just call the main function `simulate()` and it will run with the default parameters.

```
simulate()
```

To overwrite any parameters, first create objects to represent the parameters you wish to overwrite. Say for example you want to change the recombination mean, and the flagging parameters, run the following code. Whenever a parameter is not specified explicitly, the default will be used.

```
burst_params <- BurstParams(recomb_mean = 1.0)
flag_params <- FlagParameters(length_critical_region = 0,
                              prob_inactivity_when_mutated = 0)
simulate(BurstParams = burst_params, FlagParams = flag_params)
```