

Mutation models used in PRESUME

Substitutions

In PRESUME, the substitution probabilities at different positions in each sequence are defined in a time-dependent manner using GTR-Gamma model or set to a certain rate as follows.

1. GTR-Gamma model (executed by `--gtrgamma`)

This model is commonly used in evolutionary biology to describe sequence diversification by modeling of heterogeneity in substitution rates across different sequence positions.

When the GTR-Gamma model is chosen, a position-specific relative substitution parameter γ_i for every different position i is first determined by a gamma distribution with user-defined parameters μ and α as follows:

$$\gamma_i \sim \Gamma(\alpha, \mu/\alpha).$$

Let $P_i(\Delta t)$ be a 4×4 matrix, where $P_i(\Delta t)_{x \rightarrow y}$ is the transition probability from a certain source nucleotide x to a destination nucleotide y ($x, y \in \{A, C, G, T\}$) within the time interval Δt . $P_i(\Delta t)$ is then defined using γ_i and the substitution rate matrix Q , as follows:

$$P_i(\Delta t) = e^{\gamma_i \Delta t Q},$$

where Q is given by

$$Q = \begin{pmatrix} - & a_{A \rightarrow C} & a_{A \rightarrow G} & a_{A \rightarrow T} \\ a_{C \rightarrow A} & - & a_{C \rightarrow G} & a_{C \rightarrow T} \\ a_{G \rightarrow A} & a_{G \rightarrow C} & - & a_{G \rightarrow T} \\ a_{T \rightarrow A} & a_{T \rightarrow C} & a_{T \rightarrow G} & - \end{pmatrix} \begin{pmatrix} \pi_A & 0 & 0 & 0 \\ 0 & \pi_C & 0 & 0 \\ 0 & 0 & \pi_G & 0 \\ 0 & 0 & 0 & \pi_T \end{pmatrix}.$$

The sum of the diagonal values of the right-side matrix of Q must be 1 ($\pi_A + \pi_C + \pi_G + \pi_T = 1$), and the left-side matrix needs to be symmetrical (i.e., the same values are assigned to the symmetrical nucleotide transition patterns), wherein the diagonal missing values are given to satisfy the condition that every row sum of Q becomes 0. Here, $a_{x \rightarrow y}$ and π_x are user-defined parameters.

In PRESUME, the GTR-Gamma model is executed by `--gtrgamma` with the following format to specify the parameters mentioned above:

$$\text{GTR}\{a_{A,C}/a_{A,G}/a_{A,T}/a_{C,G}/a_{C,T}/a_{G,T}\} + \text{FU}\{\pi_A/\pi_C/\pi_G/\pi_T\} + \text{G}\{\alpha\}$$

2. Time-independent model (executed by `--editprofile`)

Users can define time-independent substitutions per branch (or generation) using the substitution probability matrix $\Phi_{sub,i}$ for every sequence position which original position in the root sequence is i as follows:

$$\Phi_{sub,i} = \begin{pmatrix} - & \Phi_{i,A \rightarrow C} & \Phi_{i,A \rightarrow G} & \Phi_{i,A \rightarrow T} \\ \Phi_{i,C \rightarrow A} & - & \Phi_{i,C \rightarrow G} & \Phi_{i,C \rightarrow T} \\ \Phi_{i,G \rightarrow A} & \Phi_{i,G \rightarrow C} & - & \Phi_{i,G \rightarrow T} \\ \Phi_{i,T \rightarrow A} & \Phi_{i,T \rightarrow C} & \Phi_{i,T \rightarrow G} & - \end{pmatrix},$$

where all elements in the matrix are non-negative, and the diagonal missing values are given to satisfy the condition that every row sum becomes 1.

Indels

The time-independent indel simulation is as follows. Sequences are propagated along the template propagating unit (PU) tree and progressively accumulate indels at each generation. Let i' be the nucleotide position of a sequence at a given branch of the PU tree, where its corresponding position in the root sequence is i . Insertion events in each generation (branch) are modeled by insertion probability parameters $\Phi_{ins,i}$, representing the chance of having an insertion event at the 3' side of nucleotide position i' and a position-independent size probability distribution $\psi_{ins}(L)$. The sequence of the insert length L at every event is randomly generated with equal probability for each nucleotide at every position. Deletion events are modeled by deletion probability parameters $\Phi_{del,i}$ and a position-independent size probability distribution $\psi_{del}(L)$, where a deletion of size L centering at a nucleotide position i' spanning from position $(i' - \lfloor L/2 \rfloor)$ through $(i' - \lfloor L/2 \rfloor) + L - 1$ occurs with probability $\Phi_{del,i}$, where L is stochastically determined according to $\psi_{del}(L)$. Note that $\Phi_{ins,i}$ and $\Phi_{del,i}$ are defined only for nucleotide positions i of the root sequence, reflecting that genome editing occurs only for predefined targeting sequences in cell lineage tracing.