

RANGER-DTL v1.0 Manual

Description of the Software

RANGER-DTL (short for Rapid ANalysis of Gene family Evolution using Reconciliation-DTL) is a software package for inferring gene family evolution by speciation, gene duplication, horizontal gene transfer, and gene loss. The software takes as input a gene tree (rooted or unrooted) and a rooted species tree and reconciles the two by postulating speciation, duplication, transfer, and loss events. RANGER-DTL implements efficient algorithms that can efficiently analyze trees with even thousands of taxa, and can handle both undated and dated species trees (i.e., both cladograms and chronograms). Further technical details on the reconciliation model and algorithms used by RANGER-DTL are available in the following paper:

“Efficient Algorithms for the Reconciliation Problem with Gene Duplication, Horizontal Transfer, and Loss”, Mukul S. Bansal, Eric J. Alm, and Manolis Kellis. ISMB 2012; Bioinformatics 28: i283-i291, 2012.

Using the software

The RANGER-DTL software package consists of two separate programs, `ranger-dtl-U` and `ranger-dtl-D`. The first, `ranger-dtl-U`, works with undated species trees (so any branch lengths on the species tree will be ignored), while the other, `ranger-dtl-D` works with fully-dated species trees (i.e., chronograms) and restricts any postulated transfer events to only occur between coexisting species.

Both programs take as input a single file (specified using the `-i` command line option) containing first the species tree, followed by one or more gene trees. All input trees must be expressed using the Newick format terminated by a semicolon, and they must be fully binary (fully resolved). Species names in the species tree must be unique, and the species tree must be rooted.

e.g., `((speciesA,speciesB),speciesC);`

Each leaf in the gene trees must be labeled with the name of the species from which that gene was sampled. Both programs allow for the analysis of either rooted or unrooted gene trees. An unrooted gene tree must still be written as a rooted tree in newick format, but must have the tag [&U] included at the beginning of that gene tree.

e.g., [&U](d,(e,((a,b),c)));

Any unrooted gene trees in the input will be optimally rerooted to perform the reconciliation.

Besides the input trees, both programs must also be given the costs for duplication, transfer, and loss events (specified using `-D`, `-T`, and `-L` options respectively) to be used for the reconciliation. If a cost is not specified, its default value is used instead. Note that each of these three costs must be positive integers. Thus, if one desires to use costs 1, 2.5 and 3.7 for loss, duplication, and transfer, they must use 10, 25, and 37, respectively, instead.

Test files

The four included test files illustrate the proper input format. Test files 1, 2 and 3 are meant to be used with `ranger-dtl-U`, while test file 4 is meant to be used with `ranger-dtl-D` (its species tree is a chronogram). Test files 2 and 3 illustrate the use of unrooted gene trees.

Using variable transfer costs in `ranger-dtl-D`

`ranger-dtl-D` allows the use of variable transfer costs. Specifically, it allows for two types of variable transfer costs, invoked using the `--type 1` and `--type 2` options. (`--type 0` is the default option, and means that the regular, fixed transfer cost is used.) When using Type 1 transfer costs the user specifies a threshold and an “additional transfer cost” parameter. If the distance (in the number of edges) between the donor and recipient species for the transfer is less than the threshold, the regular transfer cost is used, and for any distance larger than or equal to the threshold, the (regular transfer cost) + (additional transfer cost) is used. Type 2

variable transfer costs are a generalized version of type 1 costs, in that the transfer cost for a transfer is given by $(\text{regular transfer cost}) + \lfloor (\text{distance between donor and recipient}) / \text{threshold} \rfloor * (\text{additional transfer cost})$.

Interpreting the output

Executing the four included test files should give the user a good idea of what the output looks like and how it is to be interpreted. In general, the output consists of reconciliation blocks, one for each input gene tree. Each block begins with the species tree, same as given in the input, but for which the internal nodes have been labeled. This is followed by the gene tree, again with internal nodes labeled. (If this gene tree was originally unrooted, then the optimally rooted version of that gene tree is output here.) This is followed by the reconciliation itself, which consists of the nodes of the gene tree, output one at a time in post order. Each of these node is labeled either as a leaf node, a speciation node, a duplication node, or a transfer node, and is assigned a mapping to some node from the species tree (this mapping defines the embedding of the gene tree into the species tree). For transfer events, in addition to the mapping (which, in this case, represents the donor species), the recipient species is also specified. Finally, at the end of each reconciliation block we output the total reconciliation cost as well as the number of duplication, transfer, and loss events.

Important Notes

The program `ranger-dtl-U` implements the fast $O(mn)$ -time algorithm (where 'm' and 'n' are the number of leaves in the gene tree and species tree respectively) given in the cited paper. However, we have not yet implemented the fast $O(mn \log n)$ -time algorithm for fully-dated trees presented in the paper, and the time complexity of `ranger-dtl-D` is $O(mn^2)$. Note also that, while `ranger-dtl-D` allows variable transfer costs, `ranger-dtl-U` does not. Please let us know if you are interested in using variable transfer costs with undated species trees; we should be able to send you an executable for doing that. Finally, the current version of `ranger-dtl-D` does not error-check to make sure that the input species tree is a chronogram (or ultrametric), and will perform incorrectly if it isn't. So please ensure that any species tree you use

with ranger-dtl-D is such that all its leaf nodes are equally distant from the root in terms of the given branch lengths.

Available command line options for ranger-dtl-U

-i, --input	input file
-o, --output	output file
-D,	Duplication cost (whole number only, default value 2)
-T,	Transfer cost (whole number only, default value 3)
-L,	Loss cost (whole number only, default value 1)
--seed <integer number>	Set a user defined random number generator seed.
-q, --quiet	no process output
-s, --summary	Only output summary statistics
-v, --version	version number
-h, --help	brief help message

Available command line options for ranger-dtl-D

-i, --input	input file
-o, --output	output file
-D,	Duplication cost (whole number only, default value 2)
-T,	Transfer cost (whole number only, default value 3)
-L,	Loss cost (whole number only, default value 1)
--type 0 1 2	type of transfer cost to use
	0 - Fixed transfer cost [default]
	1 - Simple threshold based variable transfer cost
	2 - General threshold based variable transfer cost
--thr <whole number>	threshold value for use with types 1 and 2 (default 10)
--add <whole number>	additional transfer cost for use with types 1 and 2 above (default value is transfer cost divided by 2, rounded down)
--seed <whole number>	set a user defined random number generator seed.
-q, --quiet	no process output

-s, --summary	only output summary statistics
-v, --version	version number
-h, --help	brief help message