

Supplementary Methods

1 Definitions

Let a *gene family* be defined as the set of genes that are descendants of a single gene in the last common ancestor (LCA) of all n species in consideration. A genetic sequence can undergo two types of evolution: micro-level (sequence mutation) and macro level (gene duplication, gene loss, horizontal transfer, etc.) [1].

Let T be a rooted phylogenetic tree of the a gene family. The leaf nodes of a gene-tree represent modern genes from modern species, where as the internal nodes represent ancestral genes from ancestral species. When there are no macro-level events, T will have the same topology as the species-tree. In such a case, the gene tree is said to be *congruent*. The edges of a rooted tree, often called *branches*, can be uniquely identified by the more recent of its two incident nodes. Using this convention, we can refer to a branch directly above node l as b_l .

From an alignment of the genes' sequences, we can estimate the number of substitutions (or length) along every branch of T using one of several methods, such as maximum likelihood [2] or least squared error [3, 4]. We will use b_l to refer to not only the branch, but also its length as inferred from one of these methods.

1.1 Reconciliation

The species of the modern genes is often a given. Although, the same is not true of ancestral genes, the species can be inferred by *reconciling* the gene-tree to a species-tree [1, 5]. This process produces a mapping R , called a *reconciliation*, that maps every gene-tree node to a species-tree node. Page was the first to give an algorithm for finding the reconciliation that infers the fewest duplications [5]. In this parsimonious formulation, the species $R(b_l)$ of a gene b_l is defined as follows:

$$R(b_l) = LCA(\cup_{c \in children(b_l)} R(c)) \quad (1)$$

In the absence of macro-level events (i.e. a congruent gene tree), each gene node b_i maps to a unique species node i

$$R(b_i) = i$$

However, when the gene tree is incongruent, multiple gene nodes will map to a single species node. In a case, each ancestral node in an incongruent gene-tree can be labeled as a *gene-duplication* or *speciation* [5]. In the parsimonious definition of reconciliation, a node b_l is a duplication if one of its children b_c reconciles to the same species (i.e. $R(b_c) = R(b_l)$).

1.2 Generalized reconciliation

In our model, we will need a more general definition of reconciliation than is commonly used [1, 5]. In particular, we need to map duplication nodes to the middle of branches (not to nodes). Thus, we redefine a reconciliation R to be a mapping from a gene node b_l to a species node i and *duplication point* k_l .

$$R(b_l) = (i, k_l)$$

If gene node b_l is a duplication, k_l defines the fraction along the species branch i at which the gene duplication occurred: $k_l = \epsilon$ if the duplication occurs immediately after speciation of species $parent(R(b_l))$ and $k_l = 1 - \epsilon$ if the duplication occurs immediately before species $R(b_l)$. If b_l represents a speciation we

define, k_l to be 1. We define k_l to be distributed uniformly over $(0,1)$, unless an ancestor $b_{l'}$ of b_l reconciles to the same species with duplication point $k_{l'}$, in which case $k_l \sim \text{Uniform}(k_{l'}, 1)$. For example, in Figure S20 k_4 must be greater (more recent) than k_2 .

For our model, we also define a reconciliation R_b that maps gene branches to species branches. One complication, is that a gene branch may map to a path of species branches and may use only a portion of the starting and ending species branch. Thus we define

$$R_b(b_l) = ((s_1, s_2, \dots, s_m), (p_1, p_2, \dots, p_m))$$

Where the vector s_1, \dots, s_m defines the path of branches in the species tree and p_1, \dots, p_m defines the portion of each species branch used by $R_b(b_l)$. Notice, that the internal branch portions p_2, \dots, p_{m-1} , if they exist, are always 1. Defining duplication points k_l immediately implies the values of p_j , and vice versa (See Figure S23 for details).

Figure S20 illustrates an example of a complex reconciliation. Notice how gene nodes x_2 and x_4 represent duplications (stars) and both reconcile to the species branch C . Also notice how when x_0 reconciles to E and x_2 reconciles to C , that there is an implied lost branch b (dashed red) that would have diverged off of branch x_2 . Also note how gene branch x_2 has a reconciliation path with species branches D, C and fractions p_0, p_1 , where $p_0 = 1$ and $p_1 < 1$. The reconciliation path of x_4 is simply C and p_2 .

2 Gene family evolution model: simple case

In our algorithm, we use a model of evolution that captures both gene-specific and species-specific substitution rates. In this section, we first describe a model for gene families with no gene duplications or losses. Later, we generalize the model to describe gene families that contain such events.

In our model, we assume that branch b_i can be decoupled into the product of two rates:

- A *gene-specific rate* g , which specifies an underlying rate of substitution specific to the gene family, distributed as a gamma distribution $G = \Gamma(\alpha, \beta)$.
- A *species-specific rate* s_i , which specifies a scaling of that substitution rate specific to each species i . Each s_i is distributed by a normal $S_i = N(\mu_i, \sigma_i^2)$. Each S_i is independent of the other and of the gene rate G . The mean μ_i indicates the species's average rate of substitution relative to the other species, and the variance σ_i^2 indicates to what degree a gene evolving in species i deviates from the underlying gene rate.

The choice of these particular distributions was motivated by the properties of real fungal and fly gene families analyzed in this paper. We refer to branch lengths b_i as *absolute* branch lengths and call branch lengths normalized by the gene rate $s_i = \frac{b_i}{g}$ *relative* branch lengths.

This model implies several properties of gene families that we find closely fit real data. First, by the definitions above, the sum of all of the branch lengths in any gene tree, which we call the *total tree length* L , should have the following expected value:

$$E[L] = E[\sum_i b_i] = E[\sum_i s_i g] = \left(\sum_i \mu_i \right) E[g] = 1 * E[g]$$

Thus, we can use the total tree length of a gene tree as a convenient estimate of its family rate.

Another implication of the model is that all branch lengths are correlated with each other. In practice, each species-specific rate has a standard deviation σ_i that is approximately one third of the mean μ_i (Figure S7, S10). Therefore, for each branch we have $b_i \approx \mu_i g$. This implies that branch lengths from a single

species should have a strong linear correlation with the gene rate g , and thus is also approximately gamma distributed. Also, when comparing branches to each other we have the following:

$$\frac{b_i}{b_k} \approx \frac{\mu_i g}{\mu_k g} = \frac{\mu_i}{\mu_k} \quad \forall i, k$$

Therefore, the lengths of any pair of branches b_i and b_k should have a consistent ratio across different families, and thus a strong linear correlation.

Although, a normal gives non-zero probability to negative branch lengths, such branches have only a small probability $< 3\%$ when $\sigma_i/\mu_i < 1/2$ and $< 0.2\%$ when $\sigma_i/\mu_i < 1/3$.

2.1 Gene tree likelihood in simple model

Using this model, we can calculate the likelihood of a observing gene tree T with branch lengths $b_1, \dots, b_{(2n-2)}$. If we are given the gene rate g , we can calculate the likelihood of the tree T as the product of the likelihood of its branches b_i , since s_i are sampled independently.

$$P(T|G = g) = \prod_i P(B_i = b_i|G = g) = \prod_i N(S_i = \frac{b_i}{g} | \mu_i, \sigma_i^2) = \prod_i N(S_i = s_i | \mu_i, \sigma_i^2) \quad (2)$$

The equation above gives the tree likelihood conditioned on a given family rate. The total likelihood of a tree can be found by multiplying the conditional likelihood by the gene rate prior $\Gamma(g|\alpha, \beta)$ and integrating over all possible gene rates. In practice, we calculate the conditional likelihood and use the gene rate that maximizes it (Section 3.5).

3 Gene family evolution model: General case

We now generalize our model to gene-trees with duplication and loss. With the ability to calculate the likelihood of any gene-tree, we can use a tree search to find the maximum likelihood tree for a set of given sequences.

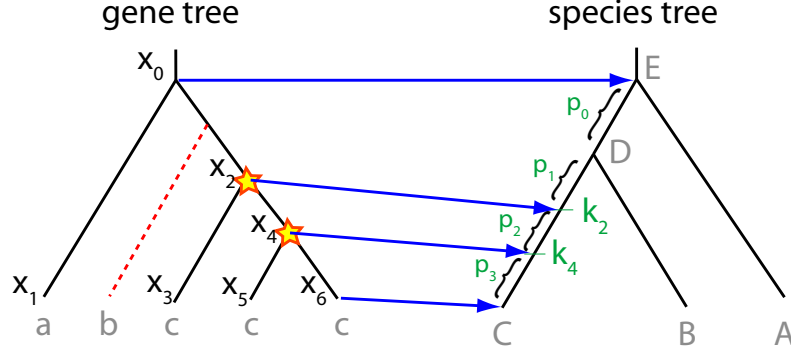
3.1 General branch length generation

Let b_l represent an absolute length and x_l its relative length (now distinct from a single species rate s_i). In the simple case, x_l was sampled from a single normal distribution S_i . In the general case, our reconciliation maps branch x_l across a path of species branches $R_b(b_l) = ((s_1, \dots, s_m), (p_1, \dots, p_m))$. Let y_j represent segments of the x_l gene branch, such that each y_j reconciles to exactly one species branch s_j with fraction p_j . Thus, x_l can viewed as the sum of m random variables y_1, \dots, y_m where each y_j is sampled from a different species branch.

For a branch $s_j, p_j = 1$, y_j is distributed by a normal distribution $y_j \sim N(\mu_{s_j}, \sigma_{s_j}^2)$. For the case where $p_j < 1$, we must now define a new distribution. We would like the new distribution to have the property that if two branches y_j and y_{j+1} reconcile to a whole species branch i (i.e. $p_j + p_{j+1} = 1$), then $y_j + y_{j+1} \sim N(\mu_i, \sigma_i^2)$. This can be achieved by defining the following:

$$y_j \sim N(p_j \mu_i, p_j \sigma_i^2)$$

Since the sum of independent normally distributed random variables is also normally distributed, we have



$$p_0 = 1; \quad p_1 = k_2; \quad p_2 = k_4 - k_2; \quad p_3 = 1 - k_4$$

$$x_1 \sim N(\mu_A, \sigma_A^2)$$

$$x_2 \sim N(p_1\mu_C, p_1\sigma_C^2) + N(p_0\mu_D, p_0\sigma_D^2) \\ = N(p_1\mu_C + p_0\mu_D, p_1\sigma_C^2 + p_0\sigma_D^2)$$

$$x_3 \sim N((1 - k_2)\mu_C, (1 - k_2)\sigma_C^2)$$

$$x_4 \sim N(p_2\mu_C, p_2\sigma_C^2)$$

$$x_5 \sim N(p_3\mu_C, p_3\sigma_C^2)$$

$$x_6 \sim N(p_3\mu_C, p_3\sigma_C^2)$$

Figure S20: Example of general reconciliation with duplications and loss. Stars represent duplications. Dashed red line represents gene loss. The k_i represent duplication points within the species tree.

$$y_j + y_{j+1} \sim N(p_j\mu_i, p_j\sigma_i^2) + N(p_{j+1}\mu_i, p_{j+1}\sigma_i^2) = N(\mu_i, \sigma_i^2)$$

Thus in general, the total length of the gene branch x_l with reconciliation path s_1, \dots, s_m and p_1, \dots, p_m is the sum of m independent normals.

$$x_l = \sum_j^m y_j \tag{3}$$

$$\sim \sum_j^m N(p_j\mu_{s_j}, p_j\sigma_{s_j}^2) = N\left(\sum_j^m p_j\mu_{s_j}, \sum_j^m p_j\sigma_{s_j}^2\right) \tag{4}$$

The benefit of this formulation is that x_l , although reconciling to a path of species branches, is still distributed by a normal distribution.

3.2 Gene tree likelihood in general model

From the model given above, we can see the length of any gene branch only depends on the species branches it reconciles to and the values of any nearby duplication points k_l . Thus, we can factor the likelihood of our model generating the entire gene tree into the likelihoods of *independent subtrees* as depicted in

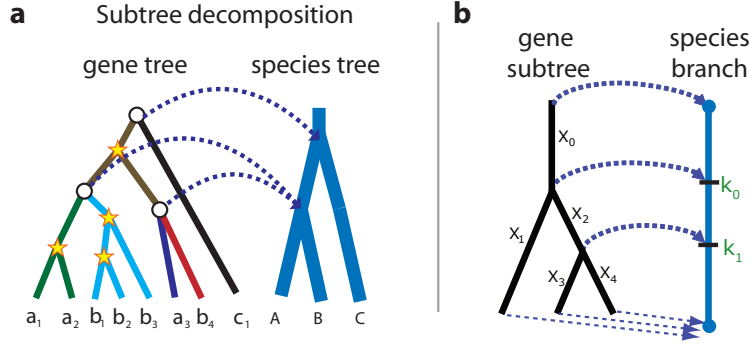


Figure S21: (a) The likelihood of any complex reconciliation can be evaluated by decomposition of the gene into independent subtrees (green, blue, brown, purple, red, black). In this example, three genes in species A ($a_1..a_3$), four genes in species B ($b_1..b_4$) and one gene in species C (c_1) arose from speciation events (circles) and duplication events (stars). Gene trees can only be decomposed at speciation nodes (circles). (b) An example of multiple duplication points needed on one species branch.

Figure S21a. Each subtree has a root that is either the root of the whole tree or is a speciation node, and each subtree has leaves that are either genes or speciations. All the internal nodes of a subtree must be duplications. Since subtrees can be generated without knowledge of the the branch lengths or duplication points in other subtrees, they are independent, and thus, the likelihood of the gene tree is simply the product of the likelihoods of its independent subtrees.

The likelihood of a subtree T' with branch lengths x_l can be calculated as follows.

$$P(T'|\mu, \sigma) = \int \cdots \int P(T'|\mu, \sigma, k_0, \dots, k_m) P(k_0, \dots, k_m) dk_0 \dots dk_m \quad (5)$$

$$= \int \cdots \int \left[\prod_l P(x_l|\mu, \sigma, k_0, \dots, k_m) \right] P(k_0, \dots, k_m) dk_0 \dots dk_m \quad (6)$$

Example tree likelihood. For the example given in Figure S21b, the likelihood is as follows.

$$\begin{aligned} P(T'|\mu, \sigma^2) &= P(X_0, X_1, X_2, X_3, X_4|\mu, \sigma^2) \\ &= \int_0^1 \int_{k_0}^1 P(X_0, X_1, X_2, X_3, X_4|k_0, k_1, \mu, \sigma^2) P(k_0, k_1, |\mu, \sigma^2) dk_1 dk_0 \\ &= \int_0^1 \int_{k_0}^1 P(X_0|k_0, \mu, \sigma^2) P(X_1|k_0, \mu, \sigma^2) P(X_2|k_0, k_1, \mu, \sigma^2) P(X_3|k_1, \mu, \sigma^2) P(X_4|k_1, \mu, \sigma^2) \\ &\quad P(k_1|k_0) P(k_0) dk_1 dk_0 \end{aligned}$$

where

$$\begin{aligned} P(k_0) &= 1 \text{ if } k_0 \in (0, 1), \text{ otherwise } 0 \\ P(k_1) &= \frac{1}{1 - k_0} \text{ if } k_1 \in (k_0, 1), \text{ otherwise } 0 \end{aligned}$$

and each $P(X_l|k_j, \dots, \mu_j, \sigma_j^2)$ is a normal of the form in (4).

3.3 Likelihood calculation implementation

In the following pseudo-code, we provide an algorithm for computing likelihood of any subtree. In our implementation, we perform the integration in (6) using sampling Figure S25. The samples are generated by the same process given in Section 3.1.

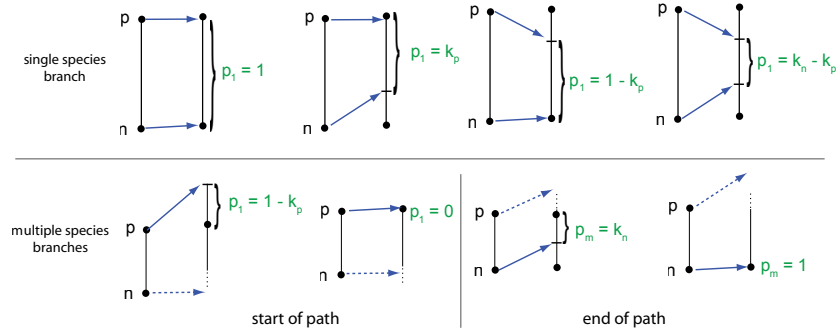


Figure S22: Cases for a gene branch reconciliation path. In each case, the gene branch n with parent p (left) reconciles to a species branch (right), sometimes with duplication points k_n or k_p . From the duplication points, the species branch fractions p_1 (starting fraction) and p_m (ending fraction) can be calculated. All interior fractions p_2, \dots, p_m , if they exist, are 1.

The details of calculating the branch likelihood term $P(X_l | k_j, \dots, \mu_j, \sigma_j^2)$ are given in the following figures and pseudo-code examples. In Figure S22, we illustrate all the cases that can occur for a reconciliation path. Only the fractions p_1 and p_m are indicated since all interior fractions p_2, \dots, p_{m-1} , if they exist, are 1. In Figure S23, is example pseudo-code that determines the fractions and necessary normal parameters (μ , σ^2) for a gene branch. Next, in Figure S24, these fractions and parameters are used to construct a single normal from which the likelihood of the branch length x_l is inferred.

3.4 Likelihood of subtrees that contain the root

It is possible that the genes in our gene tree arose from two or more ancestral genes in the LCA of the species, by duplications that precede the root of the species tree. These duplication nodes will reconcile “above” the root of the species tree, joining together two or more paralogous groups (Figure S26). Because the learned evolutionary model does not have any distributions for branch lengths above the species tree root, we can assume that all mutation rates are equally likely, and exclude these branches when evaluating the likelihoods of proposed gene trees. We call such branches *free branches* (no influence on likelihood).

A gene branch can occur such that only one end reconciles above the species tree root, thus containing both a free and a non-free portion (Figure S26). We segment such *partially free* branches into a free and non-free segment such that the likelihood of the non-free segment is maximized. If the branch length x_l is larger than the mean μ of its corresponding distribution, then the non-free segment is length μ , otherwise the non-free segment is set to capture the entire branch length x_l .

Partially free branches can usually be handled separately from all other branches, making their likelihood calculation easy to implement. There is only one case where this is not true. When rooting a tree, our algorithm chooses the midpoint of the branch that minimizes duplication and losses. One effect is that the distributions of the top two gene branches are always the same, which is valid as long as both branches are either free or non-free. However, if only one branch is free, then positioning the root at the midpoint may distort the non-free portion of the branch. We handle this case by *unfolding* the branch containing the root. Since only one child of the root needs the branch for calculating likelihood (the other child is free), we can extend the non-free child branch by the length of the free sibling and classify it as partially free.

```

def reconBranch(node, recon, events, params):

    # types of fractional branches
    enum { FRAC_NONE, FRAC_DIFF, FRAC_PARENT, FRAC_NODE }

    # types of events
    enum { EVENT_GENE, EVENT_SPEC, EVENT_DUP }

    # set fractional branches
    if recon[node] == recon[node.parent]:
        # start reconciles to a subportion of species branch
        if events[node] == EVENT_DUP:
            startfrac = FRAC_DIFF #  $k[node] - k[node.parent]$ 
        else:
            startfrac = FRAC_PARENT #  $1.0 - k[node.parent]$ 
            startparams = params[recon[node]]

        # there is only one frac
        endfrac = FRAC_NONE; endparams = NULL
    else:
        if events[node.parent] == EVENT_DUP:
            # start reconciles to last part of species branch
            startfrac = FRAC_PARENT #  $1.0 - k[node.parent]$ 
            startparams = params[recon[node.parent]]
        else:
            startfrac = FRAC_NONE; startparams = NULL

        if events[node] == EVENT_DUP:
            # end reconciles to first part of species branch
            endfrac = FRAC_NODE #  $k[node]$ 
            endparams = params[recon[node]]
        else:
            # end reconcile to at least one whole species branch
            endfrac = FRAC_NONE; endparams = NULL

    # set midparams
    if recon[node] == recon[node.parent]:
        # we begin and end on same branch, there are no midparams
        midparams = NULL
    else:
        # we begin and end on different branches
        totmean = 0; totvar = 0

        # determine most recent species branch which we fully recon to
        if events[node] == EVENT_DUP:
            snode = recon[node].parent
        else:
            snode = recon[node]

        # walk up specie path until starting species branch
        # starting species branch is either fractional or NULL
        parent_snode = recon[node.parent]
        while snode != parent_snode:
            totmean += params[snode].mu
            totvar += params[snode].sigma ** 2
            snode = snode.parent

        midparams = [totmean, sqrt(totvar)]

    return startparams, startfrac, midparams, endparams, endfrac

```

Figure S23: Pseudo-code for determining the reconciliation path of a gene branch. In the above pseudo-code, p_1 is represented as startfrac, p_m is endfrac, and all means and variances for the middle of the path are immediately combined since their fractions do not depend on any duplication points k_l . node is an object with fields node.parent and node.dist. recon is a mapping from gene nodes to species nodes. events is a mapping from gene nodes to one of the three events EVENT_*. Lastly, params is a mapping from species nodes to a tuple of parameters mu and sigma.

```

def branchLikelihood(dist, node, k, startparams, startfrac,
                    midparams, endparams, endfrac):
    totmean = 0.0
    totvar = 0.0

    if startfrac == FRAC_DIFF:
        totmean += (k[node] - k[node.parent]) * startparams.mu
        totvar += (k[node] - k[node.parent]) * startparams.sigma ** 2
    elif startfrac == FRAC_PARENT:
        totmean += (1.0 - k[node.parent]) * startparams.mu
        totvar += (1.0 - k[node.parent]) * startparams.sigma ** 2
    else startfrac == NULL:
        pass # do nothing

    if midparams != NULL:
        totmean += midparams.mu
        totvar += midparams.sigma ** 2

    if endfrac == FRAC_PARENT:
        totmean += (1.0 - k[node.parent]) * endparams.mu
        totvar += (1.0 - k[node.parent]) * endparams.sigma ** 2
    elif endfrac == FRAC_NODE:
        totmean += k[node] * endparams.mu
        totvar += k[node] * endparams.sigma ** 2
    else endfrac == NULL:
        pass # do nothing

    return normal(dist, totmean, totvar)

```

Figure S24: Pseudo-code for calculating the likelihood of a gene branch. Using the result from the previous function `reconBranch`, we can now quickly change the value of the duplication points k to get the probability of the branch length `dist`.

```

def subtreeLikelihood(subtree, recon, events, params):
    startparams = {}
    startfrac = {}
    midparams = {}
    endparams = {}
    endfrac = {}

    # determine all reconciliation paths
    for node in subtree:
        startparams[node], startfrac[node], midparams[node],
        endparams[node], endfrac[node] =
            reconBranch(node, recon, events, params)

    # integrate likelihood
    prob = 0.0
    for i in [1 to NSAMPLES]:
        k = randomSampleOfK(subtree, recon, events)

        for node in subtree:
            prob += branchLikelihood(dist, node, k, startparams[node],
                                    startfrac[node], midparams[node],
                                    endparams[node], endfrac[node]) / NSAMPLES

    return log(prob / NSAMPLES)

```

Figure S25: Pseudo-code for calculating the log likelihood of a subtree using integration by sampling. The function `randomSampleOfK` randomly samples values for the duplication points k_l as described in Section 1.2.

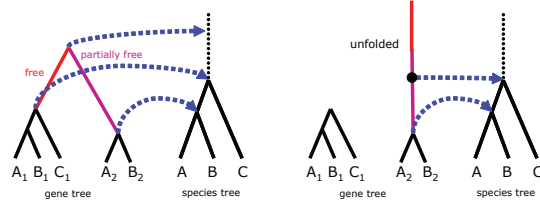


Figure S26: Left: example of free and partially free branches. Right: example of branch unfolding.

```

def treeLikelihood(tree, recon, events, params):
    # This function iterates over duplication-only subtrees of tree

    logl = 0 # log likelihood

    for subtree in getSubtrees(tree):
        logl += subtreeLikelihood(subtree, recon, events, params)

    return logl

```

Figure S27: Pseudo-code for calculating the likelihood of a tree.

3.5 Gene rate estimation

SPIDIR uses the maximum likelihood estimate (MLE) of the gene-species rate \hat{g} for each proposed tree. This is currently done as an approximation. Ideally, the tree likelihood should be integrated over all possible gene-specific rates. The MLE can be estimated using the known prior distribution of gene rates $gamma(g|\alpha, \beta)$ and the distribution of species-specific rates $normal(s_i|\mu_i, \sigma_i^2)$ for each branch i of the tree. Using this information, \hat{g} can be found by solving the following cubic equation (derivation in Appendix A).

$$\hat{g}^3 - \left(\frac{\alpha - 1}{\beta}\right) \hat{g}^2 + \left(\frac{1}{\beta} \sum_i \frac{\mu_i b_i}{\sigma_i^2}\right) \hat{g} - \left(\frac{1}{\beta} \sum_i \frac{b_i^2}{\sigma_i^2}\right) = 0 \quad (7)$$

A convenient benefit of this equation is that one can find the MLE gene rate even for gene trees with gene loss and duplication events. That is, the summations in the third and fourth terms normalize the gene rate for branches that do not appear at all (due to loss) or appear multiple times (due to duplication). With a gene rate estimate, normalized branch lengths x_i can be calculated by b_i/\hat{g} .

3.6 Tree search

To propose topologies for our likelihood function, we have implemented several search strategies. Since the tree space is exponential in the number of genes, we must rely on heuristics to search the tree space. We have implemented many tree search heuristics popular in other phylogenetic reconstruction algorithms, including nearest neighbor interchange (NNI) and subtree pruning regrafting (SPR) in both a greedy and MCMC framework. In our evaluations, we used a modified version of the MCMC tree search.

In our program, the most parsimonious reconciliation is always reconsidered, but this is only used for reducing search space.

3.7 Distance errors

Each proposed tree topology is fitted with branch lengths using the distance matrix and Least Square Error (LSE). However, if the a topology disagrees drastically with the distance matrix, there may not exist any reasonable way to fit branch lengths. To detect such cases, we use the distance error $error(T)$ of a tree T . We penalize trees with high error since their branch lengths do not reflect the original sequence data. Here, we outline our use of distance error.

Let D be a square $n \times n$ matrix of pair-wise distances between sequences i and j in a multiple sequence alignment. Let D' be a square $n \times n$ matrix such that $D'_{i,j}$ equals the total path length between the pair of leaves i and j in T . Unless the distances are perfectly additive and the tree topology is correct, these two matrices will differ. In real data, distances are rarely perfectly additive. LSE assigns branch lengths to T that minimize the difference (error) between these two matrices.

$$\sum_{i,j} (D_{i,j} - D'_{i,j})^2$$

This error is dependent on the number and divergence of the genes. Therefore, we define a relative distance error

$$error(T) = \sqrt{error(T) / \sum_{i,j} D_{i,j}}$$

We find that this quantity is more consistent across trees with different family rates. The probability of a correct tree having a relative distance error of $error(T)$ quickly diminishes as $error(T)$ increases. We were able to improve the accuracy of our algorithm by incorporating the distance error into our likelihood function. For a given tree T , let $P(T)$ be the likelihood as determined by our gene family model. The adjusted likelihood of the tree T is then

$$P'(T) = P(T) * e^{-K_e * error(T)}$$

where K_e is a parameter for how much to penalize the distance error. In our experience, reconstruction accuracy is most improved for K_e values in the range of 100 to 600.

3.7.1 Duplication and loss

We have implemented a simple penalty for duplication and loss. Using a parsimony approach the number of gene duplications N_d and losses N_l are estimated from a reconciled gene-tree. Let D and L be the probabilities that a duplication or loss appears in a gene-tree. Using a simplifying assumption that duplications and losses occur independently, the adjusted likelihood of a gene-tree T with likelihood $P(T)$ is then

$$P'(T) = L^{N_l} * D^{N_d} * P(T)$$

In this paper, we have only explored adjusting the probability of D . Gene loss probabilities were not used.

References

- [1] Goodman M, Czelusniak J, Moore G, Romero-Herrera A, Matsuda G (1979) Fitting the gene lineage into its species lineage, a parsimony strategy illustrated by cladograms constructed from globin sequences. *Systematic Zoology* 28:132–163.

- [2] Felsenstein J (1981) Evolutionary trees from dna sequences: A maximum likelihood approach. *Journal of Molecular Evolution* 17:368–376.
- [3] Cavalli-Sforza L, Edwards A (1963) Analysis of human evolution. *Proc 11th Internat Congr Genetics, The Hague* :923?933.
- [4] Desper R, Gascuel O (2002) Fast and accurate phylogeny reconstruction algorithms based on the minimum-evolution principle. *J Comput Biol* 9:687–705. doi:10.1089/106652702761034136. URL <http://dx.doi.org/10.1089/106652702761034136>.
- [5] Page R (1994) Maps between trees and cladistic analysis of historical associations among genes, organisms, and areas. *Systematic Biology* 43:58–77.
- [6] Stoye J, Evers D, Meyer F (1998) Rose: generating sequence families. *Bioinformatics* 14:157–163.
- [7] Hasegawa M, Kishino H, Yano T (1985) Dating of the human-ape splitting by a molecular clock of mitochondrial dna. *J Mol Evol* 22:160–174.

A Estimating gene tree base rates

Let g be the gene rate. For each branch i , we have branch length b_i , species rate s_i , relative rate mean μ_i and standard deviation σ_i .

Given several b_i we want to find maximum likely g .

$$\begin{aligned}
 b_i &= s_i g \\
 P(S_i = s_i | \mu_i, \sigma_i) &= \text{normal}(s_i | \mu_i, \sigma_i^2) \\
 g &= b_i / s_i \quad \forall i \\
 P(G = g | \alpha, \beta) &= \text{gamma}(g | \alpha, \beta) \\
 P(G = g | b, \mu, \sigma, \alpha, \beta) &= P(b | g, \mu, \sigma, \alpha, \beta) P(G = g | \mu, \sigma, \alpha, \beta) / P(b | \mu, \sigma, \alpha, \beta) \\
 &= P(b | g, \mu, \sigma) P(G = g | \alpha, \beta) / P(b | \mu, \sigma, \alpha, \beta) \\
 &= P(b_1 / s_1 = g, \dots, b_n / s_n = g | \mu, \sigma) P(G = g | \alpha, \beta) / P(b | \mu, \sigma, \alpha, \beta)
 \end{aligned}$$

Since s_i are independent from each other, we can factor this likelihood for each branch.

$$\begin{aligned}
 P(G = g | b, \mu, \sigma, \alpha, \beta) &= P(G = g | \alpha, \beta) / P(b | \mu, \sigma, \alpha, \beta) \prod_i P(b_i / s_i = g | \mu, \sigma) \\
 &= P(G = g | \alpha, \beta) / P(b | \mu, \sigma, \alpha, \beta) \prod_i P(s_i = b_i / g | \mu, \sigma) \\
 &= \text{gamma}(g | \alpha, \beta) / P(b | \mu, \sigma, \alpha, \beta) \prod_i \text{normal}(b_i / g | \mu_i, \sigma_i^2)
 \end{aligned}$$

The maximum likelihood estimator of g given α , β , x_i , μ_i , and σ_i for all i is then

$$\begin{aligned}
 \arg \max_g P(G = g | b, \mu, \sigma, \alpha, \beta) &= \arg \max_g \text{gamma}(g | \alpha, \beta) \prod_i \text{normal}(b_i / g | \mu_i, \sigma_i^2) \\
 &= \arg \max_g g^{\alpha-1} e^{-\beta g} \frac{\beta^\alpha}{\Gamma(\alpha)} \prod_i \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left(-\frac{(b_i/g - \mu_i)^2}{2\sigma_i^2}\right)
 \end{aligned}$$

Applying a natural log, which does not change the value of argmax, gives

$$\begin{aligned}
 &= \arg \max_g (\alpha - 1) \log g - \beta g + \log \frac{\beta^\alpha}{\Gamma(\alpha)} + \sum_i \log \frac{1}{\sigma_i \sqrt{2\pi}} + \sum_i \left(-\frac{(b_i/g - \mu_i)^2}{2\sigma_i^2} \right) \\
 &= \arg \max_g (\alpha - 1) \log g - \beta g + K + \sum_i \left(-\frac{(b_i/g - \mu_i)^2}{2\sigma_i^2} \right)
 \end{aligned}$$

Where K is the constant $\log \frac{\beta^\alpha}{\Gamma(\alpha)} + \sum_i \log \frac{1}{\sigma_i \sqrt{2\pi}}$. Setting the derivate to zero gives

$$\begin{aligned}
0 &= \frac{d}{dg}(\alpha - 1) \log g - \beta g + K + \sum_i \frac{-1}{2\sigma_i^2} \left(\frac{b_i}{g} - \mu_i \right)^2 \\
&= \frac{\alpha - 1}{g} - \beta + \sum_i \frac{-1}{\sigma_i^2} \left(\frac{b_i}{g} - \mu_i \right) (-b_i g^{-2}) \\
&= \frac{\alpha - 1}{g} - \beta + \sum_i \left(\frac{b_i^2}{g^3 \sigma_i^2} - \frac{\mu_i b_i}{g^2 \sigma_i^2} \right) \\
&= (\alpha - 1)g^2 - \beta g^3 + \sum_i \frac{b_i^2}{\sigma_i^2} - \sum_i \frac{\mu_i b_i g}{\sigma_i^2}
\end{aligned}$$

Which gives the following cubic polynomial

$$g^3 - \frac{\alpha - 1}{\beta} g^2 + \left(\frac{1}{\beta} \sum_i \frac{\mu_i b_i}{\sigma_i^2} \right) g - \frac{1}{\beta} \sum_i \frac{b_i^2}{\sigma_i^2} = 0$$