

# SPIDIR: Species informed gene phylogeny reconstruction across multiple complete genomes

Matt Rasmussen and Manolis Kellis  
{rasmus,manoli}@mit.edu

October 1, 2006

## Abstract

As comparative genomics is coming of age, dozens of fully sequenced genomes are becoming available in the mammalian, insect, worm, and fungal clades, offering an enormous opportunity to identify and analyze the functional elements of eukaryotic and animal genomes. Phylogeny is the most principled means of genome comparison, as a tree can represent nearly all possible evolutionary relationships, including orthologs, paralogs, gene duplication and loss, varying mutation rates, and highly constrained elements. However, current phylogeny methods have limited power to systematically and accurately reconstruct gene families at the genome-wide level.

In this paper, we present a new algorithm for phylogenetic reconstruction called SPIDIR, SPecies Informed DIstanced-based Reconstruction, that unlike previous methods incorporates the known species tree to reconstruct gene phylogenies with a high level of accuracy. SPIDIR uses a model of gene family evolution trained on a small subset of gene families with unambiguous phylogeny to reconstruct the phylogeny of genes with more complex cases of gene duplication and loss.

We evaluate our method on a variety of datasets, including gene phylogenies from 12 species of fly, 17 species of fungi, 7 species of yeast both preceding and following the whole-genome duplication, as well as simulated data. We find that inclusion of species information dramatically improves phylogenetic reconstruction accuracy enabling phylogenomics analysis of dozens of species and thousands of genes.

## 1 The challenge of phylogenomics

The increasing availability of fully-sequenced genomes from closely related species has led to a growing demand for principled methods for genome comparison. *Phylogenomic* has been proposed as a framework for applying the rigorous methods of phylogeny on the genome scale [6], resolving the evolutionary history of every gene family across multiple complete genomes. Given an accurate phylogeny, much of the desired comparative information can be inferred, including all orthologs, paralogs, gene duplications and losses, varying rates of mutation, and highly constrained elements. As opposed to methods based on pair-wise sequence comparisons (Inparanoid [22], COG [29], OrthoMCL [19]), gene phylogenies provide a consistent and scalable framework as well as the proper evolutionary context for studying gene evolution.

Progress has recently been made in realizing the full phylogenomic pipeline. The general approach is to coarsely identify clusters of genes with suspected homology [19] and then use phylogenetic methods to refine ancestral relations by constructing a tree for each cluster. Gene trees are then compared to the species tree in order to reconcile any differences they may have due to gene duplication or loss [9]. Simple and efficient algorithms exist to find the reconciliation that implies the fewest duplications [32], and gene trees can be efficiently rooted such that duplications and losses are minimized [31]. Projects that have implemented some or all of these steps include RIO [31], HOPS [27], and TreeFam [18].

Unfortunately, all of these phylogenomic strategies have encountered great difficulties in attaining accurate gene phylogenies. In the phylogenomics setting, these inaccuracies quickly propagate, resulting in extremely unlikely, large numbers of duplications and losses (usually 1-2 times the genome size)

and wide-spread cases of unrecognized orthology. Each of these methods handles gene tree inaccuracies differently, such as avoiding the gene tree and only focusing on orthology pairs (RIO), using a more relaxed reconciliation that infers fewer duplications (HOPS), or using human curators that use additional information to manually correct gene trees (TreeFam).

The performance of phylogenetic reconstruction has been well studied [25, 28]. It is also known that accuracy depends on sequence length [20, 23]. In our dataset of 12 flies, we have found that the typical gene (length  $\sim 950$ bp) is reconstructed correctly only 42% of the time, regardless of the phylogenetic method used. As sequence length increases to 2200bp, which is larger than 90% of all fly genes, reconstruction accuracy is still only 60%. To reconstruct the phylogeny of 8 species of yeast, Rokas *et. al.* concatenated 20 orthologs genes (15-30kb) in order to attain a suitable reconstruction accuracy [23].

These results indicate that the average gene simply does not contain enough informative characters for accurate phylogenetic reconstruction. This poses a significant problem for all phylogenomic efforts and rigorous genome comparisons. However, the phylogenomic problem is different from traditional phylogenetics, in that gene sequences are not anonymous; they are associated with species with known ancestry. Thus, leveraging knowledge of the species tree may hold the solution to accurate phylogenetic reconstruction. The connection of species and gene trees has been greatly investigated [9, 21], and several recent approaches for gene tree construction and reconciliation have been taken [1, 4], but a formulation that is both rigorous and practical remains to be proposed.

## 2 Methods

In this paper, we present a new phylogenetic reconstruction algorithm that is specifically designed for the phylogenomic setting. We frame phylogenetic reconstruction as a learning problem. We introduce a new model of gene family evolution that accurately captures the roles of gene-specific and species-specific mutation rates. We present a training procedure for estimating the parameters of our model from a dataset of genes with known phylogeny. Using this model, we formulate an algorithm for calculating the likelihood of any phylogeny for genes with unknown ancestry. We have implemented this algorithm in a new software package, SPIDIR (SPecies Informed DIstance-based Reconstruction), that combines our likelihood function with several tree search procedures to reconstruct high accuracy phylogenies. The software is available to the public online at <http://compbio.mit.edu/spidir>.

### 2.1 A new model of gene family evolution

In this paper, we introduce a new computational model for the evolution of gene families in the context of rich species phylogenies. We find that this model well approximates precise evolutionary properties measured across a large compendium of species, from both flies and fungi (see results).

In this section, we first describe our model considering only gene trees that contain no duplications or losses, and therefore have a one-to-one mapping between gene branches and species branches. In Section 2.3.2, we show how the model generalizes to any possible gene tree.

We model gene families as being shaped by two forces of evolution: a gene family specific mutation rate and species specific mutation rates. Let  $b_{ij}$  be the length of each branch in gene tree  $T_j$  that reconciles to species branch  $i$ . This branch is generated as the product of two rates:

- A *gene family rate*  $F_j$ , which specifies an underling rate of mutation specific to gene family  $j$ , distributed as a gamma distribution  $\Gamma(\alpha, \beta)$ . The parameters of the distribution allow for fast-evolving and slow-evolving gene families.
- A *species specific rate*  $S_i$ , which specifies a scaling of that mutation rate specific to each species  $i$ . Each  $S_i$  is independently distributed by a normal distributed  $N(\mu_i, \sigma_i^2)$ . The mean  $\mu_i$  indicates the species's average rate of mutation relative to the other species, and the variance indicates to what degree a gene evolving in species  $i$  deviates from the underling family rate  $F_j$ .

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  $j$ , which we call the *total*

tree length  $L_j$ , should have the following expected value:

$$E[L_j] = E\left[\sum_i b_{ij}\right] = E\left[\sum_i S_i F_j\right] = \left(\sum_i \mu_i\right) E[F_j] \propto E[F_j]$$

Thus, we can use the total tree length of a gene tree to estimate 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  $\sigma_i$  that is approximately one fourth of the mean  $\mu_i$ . Therefore, for each branch we have  $b_{ij} \approx \mu_i F_j$ . This implies that each branch has a strong linear correlation with the family  $F_j$ , and thus is also gamma distributed. Also, when comparing branches to each other we have the following:

$$\frac{b_{ij}}{b_{kj}} \approx \frac{\mu_i}{\mu_k} \quad \forall i, k$$

Therefore, the lengths of any pair of branches  $b_{ij}$  and  $b_{kj}$ , in the same gene family  $F_j$ , should have a consistent ratio across different species  $i$  and  $k$ , and thus a strong linear correlation.

Although, a normal gives non-zero probability to negative branch lengths, we find that in practice this is only true for a small fraction of the total probability (about 3%). All of these properties appear in real gene families from fly and fungal species (see Section 3.1). These model assumptions also lend themselves to efficient computation.

## 2.2 Model training

Given a set of  $N$  known gene phylogenies, we can learn the parameters of our gene family model. The model consists of the parameters  $\alpha$  and  $\beta$  for family rate distribution and the  $\mu_i$  and  $\sigma_i^2$  for each species. Here, we show how given either family or species parameters, one can estimate the other. Therefore, by using an Expectation Maximization (EM) algorithm, we can estimate both sets of parameters simultaneously.

Let each tree  $T_j$  in the training set be a gene tree with no duplications or losses. For each  $T_j$ , we make an initial estimate of its family rate  $F_j$  to be the total tree length  $L_j$ . We then fit the family rate parameters  $\alpha$  and  $\beta$  by determining their maximum likelihood estimates. Given the family rate, we can also estimate each species specific rate to be  $s_{ij} = b_{ij}/F_j$ . This produces an  $s_{ij}$  for each tree. To estimate the normal distribution  $S_i$  we find the maximum likelihood parameters  $\mu_i = \frac{1}{N} \sum_j s_{ij}$  and  $\sigma_i^2 = \frac{1}{N-1} \sum_j (s_{ij} - \mu_i)^2$ .

At each EM iteration, we use our estimated species specific rate distributions to improve our family rate estimates. The maximum likelihood estimate (MLE) of the family rate  $\hat{F}_j$  for any tree  $T_j$  can be found by finding the largest root of the following cubic equation (see Section A.1 for derivation).

$$\hat{F}_j^3 - \left(\frac{\alpha - 1}{\beta}\right) \hat{F}_j^2 + \left(\frac{1}{\beta} \sum_i \frac{\mu_i b_{ij}}{\sigma_i^2}\right) \hat{F}_j - \left(\frac{1}{\beta} \sum_i \frac{b_{ij}^2}{\sigma_i^2}\right) = 0 \quad (1)$$

Using these equations, we can iterate our EM and estimate every parameter in our model. During training on our fly dataset (Section 3.1), we find that the family rates  $F_j$  converge to a value that remains proportional to total tree length and is roughly twice as large.

## 2.3 Calculating the likelihood of any gene tree

Using our gene family evolution model, we can estimate the likelihood of any gene tree given the sequence data. This likelihood function combined with a tree search that proposes topologies can find the maximum likelihood phylogeny for a set of genes with unknown ancestry. In this section, we first show how to calculate the likelihood of a proposed topology that has a simple reconciliation to the species tree, such that there are no duplications or losses. Next, we show how the likelihood of topologies with more complex reconciliations can be calculated with only a few extensions.

For each proposed topology  $T$ , we fit branch lengths using an input distance matrix and the Least Squared Error algorithm [2]. The resulting tree is then rooted and reconciled to the species tree such

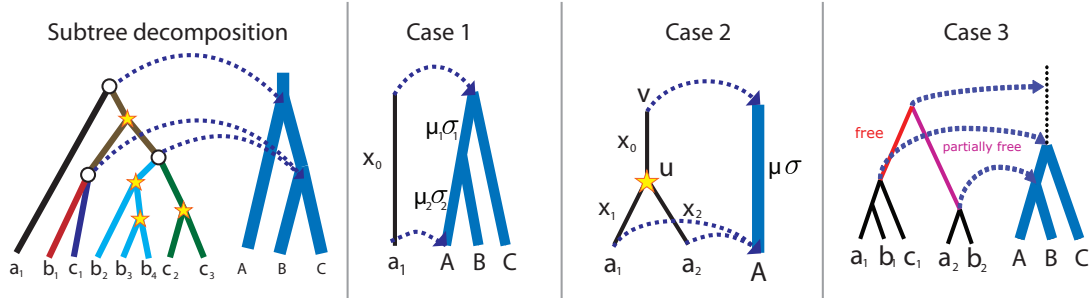


Figure 1: The likelihood of any complex reconciliation can be evaluated by decomposition into one of three cases shown here, and described in section 2.3. Gene-tree nodes (left) are mapped to species-tree nodes (right) by reconciliation. In this example, one gene in species A ( $a_1$ ), four genes in species B ( $b_1..b_4$ ) and three genes in species C ( $c_1..c_3$ ) arose from speciation events (circles) and duplication events (stars). Gene trees can only be decomposed at speciation nodes (circles).

that number of duplications and losses are minimized [32]. Using the branch lengths  $b_i$ , we determine the MLE family rate  $\hat{f}$ . Next, we *normalize* the gene tree by dividing each branch by the family rate. Thus, for each normalized branch length  $x_i$ , we have  $x_i = b_i/\hat{f}$ .

### 2.3.1 Likelihood of simple reconciliation

A gene tree that has a simple reconciliation reconciles every gene branch to exactly one species branch. The distribution for every species branch length  $x_i$  is a normal with parameters  $\mu_i$  and  $\sigma_i^2$ . Therefore, the likelihood of seeing each branch can be computed as  $N(x_i|\mu_i, \sigma_i^2)$ . Since the normalized branches are independent, the likelihood of seeing a tree  $T$  is the product of the likelihood of its branches.

$$P(T|\mu, \sigma) = \prod_i N(x_i|\mu_i, \sigma_i^2)$$

### 2.3.2 Likelihood of complex reconciliation

To calculate the likelihood of a tree with complex reconciliation, we only need a few reasonable assumptions. First, we assume that gene loss does not affect the mutation rate of any other branch within the gene family. Second, we assume that when a gene duplication node occurs, the lengths of paths traveling through it from an ancestral speciation node to a more recent one are still distributed by a normal, but we allow sibling branches extending below a duplication node to be independent given their parent branch. This would enable for example genes to take on new function and selective regimes after a duplication event.

The likelihood of any gene tree can be factored into the likelihoods of its subtrees. We will show that these subtrees are all independent and provide a procedure for each likelihood calculation. We decompose a gene tree into subtrees, such that the root and leaves of a subtree are either gene or a speciation event (Figure 1). Each of these subtrees, reconcile to an integer number of species branches. Each subtree falls into atleast one of three cases: either (1) they consist of a single branch, (2) they contain multiple branches with duplication nodes, or (3) it contains the root of the tree, which represents a gene duplication.

**Case 1: Likelihood of a single branch subtree** Subtrees with a single gene branch  $x$  either reconcile to a single species branch  $y$  or to multiple species  $y_1, y_2, \dots, y_k$ . If  $x$  reconciles to a single species branch  $y$ , we can calculate its likelihood with the normal distribution  $N(x|\mu_i, \sigma_i^2)$  like we did in the simple gene tree. A single gene branch  $x$  that reconciles to multiple species branches  $y_1, \dots, y_k$ , implies at least  $k - 1$  gene loss events, that gene branch  $x$  has evolved through each of the  $k$  species. Therefore, the gene branch length should actually be the sum of all  $k$  gene branches each of which

is independent and normally distributed. The distribution of our gene branch  $x$  is thus the sum of independent normals:

$$x \sim N(\mu_{i_1}, \sigma_{i_1}^2) + \dots + N(\mu_{i_k}, \sigma_{i_k}^2) = N(\mu_{i_1} + \dots + \mu_{i_k}, \sigma_{i_1}^2 + \dots + \sigma_{i_k}^2)$$

Using this distribution, we can again directly calculate the likelihood of seeing a branch of length  $x$ . In both cases, the likelihood of  $x$  only depends on the species to which it reconciles, and therefore is independent of other subtrees.

**Case 2: Likelihood of a multi-branch subtree** To calculate the likelihood of subtrees with duplication nodes, we can use some of the techniques developed in case 1. By definition, the root of the gene subtree corresponds to a speciation node. The subtree leaves are either genes or internal speciation nodes. Therefore, every path from the subtree root to a subtree leaf reconciles to an integer number of species branches, similarly to case 1. However, each path shares some branches with other paths, and thus their likelihoods are not independent, and a simple product of path likelihoods cannot be used.

Before developing a general equation for the likelihood of a subtree with an arbitrary number of duplications and losses, first consider the following tree (Figure 1 case 2). Let  $a_1$  and  $a_2$  be two paralogs. Their immediate ancestor  $u$  represents a duplication event, and the root of subtree  $v$  represents a speciation event. Both paths  $v \rightarrow a_1$  and  $v \rightarrow a_2$  reconcile to the same species branch which has parameters  $\mu$  and  $\sigma^2$ . Let the lengths of the paths ending at  $a_1$  and  $a_2$  be random variables  $X_1$  and  $X_2$ , respectively. Therefore, we need to compute the following likelihood:

$$P(X_1 = x_0 + x_1, X_2 = x_0 + x_2) = P(X_2 = x_0 + x_2 | X_1 = x_0 + x_1)P(X_1 = x_0 + x_1)$$

Conditioning  $X_2$  on  $X_1$  imposes that  $X_2$  be greater than the shared branch length,  $x_0$ . (assuming branch lengths are non-negative). Thus, the likelihood of  $X_2$  is simply scaled by the probability that  $X_2 > x_0$  (see Section A.2). We can thus write:

$$P(X_2 = x_0 + x_2 | X_1 = x_0 + x_1)P(X_1 = x_0 + x_1) = \frac{P(X_2 = x_0 + x_2)}{P(X_2 > x_0)}P(X_1 = x_0 + x_1) \quad (2)$$

$$= \frac{N(x_0 + x_2 | \mu, \sigma^2)N(x_0 + x_1 | \mu, \sigma^2)}{C(x_0 | \mu, \sigma^2)} \quad (3)$$

Where  $C(x_0 | \mu, \sigma^2) = \int_{x_0}^{\infty} N(a | \mu, \sigma^2) da$  is the cumulative normal distribution which also has a closed form solution. Notice that the final likelihood does not depend on which path we condition upon. Either way, the likelihood is a product of each path's likelihood divided by the probability of seeing a branch longer than  $x_0$ .

This calculation can be extended to subtrees with any number of duplication nodes and gene losses. For each path length  $p_k$  in the subtree, we first determine the species branches  $R_k$  to which it reconciles. The likelihood of  $p_k$  is  $N(p_k | \sum_{i \in R_k} \mu_i, \sum_{i \in R_k} \sigma_i^2)$ . Next, we determine how to condition on path overlaps, using a depth first search (DFS), starting at the root  $v$  of the subtree. For the first path found to a leaf  $p_1$ , we use an unconditioned likelihood. We then continue the DFS until a new unvisited node is found, and record the shared (conditioned) path as  $c_2$  and the set of species  $Q_2$  to which  $v \rightarrow c_2$  reconciles. We then continue until the DFS reaches a new leaf  $p_2$ . Let  $X_1 = \sum_{i \in R_1} S_i$ ,  $X_2 = \sum_{i \in R_2} S_i$ , and  $Y_2 = \sum_{i \in Q_2} S_i$ . Accordingly, the conditional probability of  $p_2$  given  $p_1$  is

$$P(X_2 = p_2 | X_1 = p_1) = P(X_2 = p_2 | Y_1 > c_2) \quad (4)$$

$$= \frac{N(p_2 | \sum_{i \in R_2} \mu_i, \sum_{i \in R_2} \sigma_i^2)}{C(c_2 | \sum_{i \in Q_2} \mu_i, \sum_{i \in Q_2} \sigma_i^2)} \quad (5)$$

Notice that the previous subtree example had no gene loss, therefore  $Q_k = R_k$ . This simplifies equation 5 to equation to the conditioned term in equation 3. Therefore, equation 5 is just a more generalized equation for conditional branch likelihood. After DFS, all leaf terminating paths  $p_k$  and conditioned paths  $c_k$  will be found. Using this information, we can calculate the likelihood of the entire subtree  $T'$ :

$$\begin{aligned}
P(T') &= P(X_1 = p_1) * P(X_2 = p_2 | X_1 = p_1) * P(X_3 = p_3 | X_1 = p_1, X_2 = p_2) * \dots * \\
&\quad P(X_n = p_n | X_1 = p_1, \dots, X_{n-1} = p_{n-1}) \\
&= P(X_1 = p_1) * P(X_2 = p_2 | Y_2 > c_2) * P(X_3 = p_3 | Y_3 > c_3) * \dots * P(X_n = p_n | Y_n > c_n) \\
&= \frac{\prod_{k \in \text{leaves}(T')} N(p_k | \sum_{i \in R_k} \mu_i, \sum_{i \in R_k} \sigma_i^2)}{\prod_{k \in \text{internals}(T')} C(c_k | \sum_{i \in R_k} \mu_i, \sum_{i \in R_k} \sigma_i^2)}
\end{aligned}$$

Notice, that no matter how we choose to visit the nodes in the DFS, that each leaf is visited once, and each internal node is recorded as the end of a conditioned path once. Also the final product disassociates all numerators from denominators. Therefore, it does not matter how we choose to guide our DFS. This equation, encompasses all cases we have previously seen.

**Case 3: Likelihood of subtrees that contain the root** Our algorithm does not assume that the given set of genes are all orthologs to each other. Instead, it allows the possibility that the gene set represents a collection of paralogous families with duplication nodes that precede all speciation events. These duplication nodes will reconcile “above” the root of the species tree, joining together two or more paralogous families. 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 simply ignore such branches when comparing the likelihoods of proposed gene trees. We call such branches *free branches*, because they get to grow for free (at no likelihood cost). All other branches are called *non-free*.

A gene branch can also occur such that only one end reconciles above the species tree root, thus containing both a free portion and a non-free portion (Figure 1 case 3). SPIDIR chooses the midpoint of these *partially free* branches such that the likelihood of the non-free branch is maximized, which can be done by simply choosing the minimum of the gene branch length and the mean of the branch distribution.

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 works out fine as long as both branches are either free or non-free. However, if only one branch is non-free, then positioning the root at the center of the rooting branch may be distorting the non-free portion of the branch. We handle this case by *unfolding* the rooting branch. 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. Now, the branch can be split such that its likelihood is maximized.

**Bringing it all together** Using these cases, we can compute the likelihood of every possible subtree. The likelihood of each subtree only depends on its branches and the species to which it reconciles. Therefore, subtrees are independent, and the product of all subtree likelihoods equals the likelihood for the entire gene tree. All need equations can be represented in closed form, and can be computed in a small number of constant operations for each branch. The runtime for the entire likelihood function is linear in the number of genes.

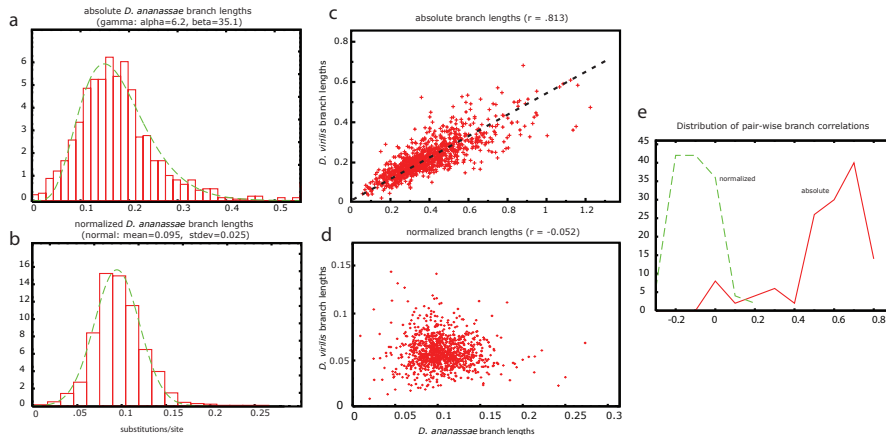


Figure 2: (a) Distribution of 5154 absolute branch lengths ( $b_{ij}$ ) follow a gamma distribution. (b) Distribution of normalized branch lengths ( $b_{ij}/F_j$ ) follows a normal distribution. (c) Absolute branch lengths show strong pairwise linear correlations across different species (in this example, *D. ananassae* and *D. virilis* branch lengths show a Pearson’s correlation of  $R = .81$ ). (d) Normalized branch lengths are largely independent ( $R = .08$ ) (e) Normalization by the family rate  $F_j$  largely eliminates dependencies of all pair-wise branch correlations in the fly tree.

## 2.4 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.

## 3 Results

For our training and testing sets, we used syntentically aligned orthologs from the largest groups of fully sequenced eukaryotic species: 12 flies (*D. melanogaster*, *D. sechellia*, *D. simulans*, *D. yakuba*, *D. erecta*, *D. ananassae*, *D. pseudoobscura*, *D. persimilis*, *D. willistoni*, *D. mojavensis*, *D. virilis*, *D. grimshawi*) [3] spanning 40 million years of evolution and 17 fungi (*S. cerevisiae*, *S. paradoxus*, *S. mikatae*, *S. bayanus*, *S. castellii*, *C. glabrata*, *K. lactis*, *A. gossypii*, *K. waltii*, *C. albicans*, *C. dublinensis*, *C. tropicalis*, *C. parapsilosis*, *L. elongisporus*, *C. lusitaniae*, *C. guillermondii*, *D. hansenii*) [14, 12, 13, 15, 11, 16] spanning 150 million years. We use synteny information to identify sets of unambiguous orthologs such that there is exactly one gene present from each species in consideration. For each set of orthologs, we can conclude that all divergence is due to speciation and assume that no gene duplication or loss has occurred. Therefore, the true gene phylogeny is exactly the one matching the species phylogeny. Using this technique, we have identified 5154 sets of fly genes and 739 sets of fungal genes with unambiguous gene phylogenies. For the fungi, syntenic segments were found separately in the Saccharomycota and Candidiae clades and merged using Best Bidirectional BLAST hits. We aligned these genes using the MUSCLE sequence alignment program on protein sequences [5]. To attain nucleotide alignments, we map the nucleotide sequence on the peptide alignments, substituting every amino acid by the corresponding codon and every gap by a triplet of gaps. Maximum likelihood branch lengths were estimated for each tree using the PHYML software package [10].

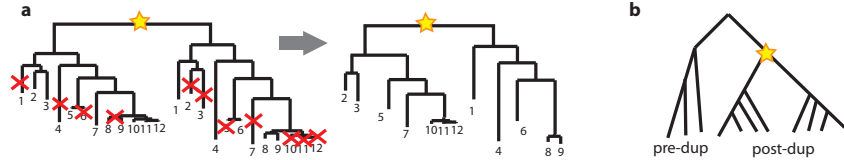


Figure 3: (a) Example of constructing large fly families with simulated loss. (b) Yeast Whole-Genome Duplication gene tree

### 3.1 Evolutionary model captures properties of fungal and fly trees

We found that key properties of actual gene families match very closely the properties of our model of gene evolution. Indeed, absolute branch lengths are distributed as a gamma distribution (Figure 2a), as expected for a mutation rate. The normalized branch lengths  $b_{ij}/F_j$  are indeed distributed according to a normal (Figure 2b). Furthermore, we find a significant correlation between the rates of orthologous genes across any pair of species (Figure 2c), namely that fast genes in one species are also fast in any other species. However, after normalizing by the family rate  $F_j$ , we find that pairs of branches become largely independent (Figure 2d). The average Pearson’s correlation drops from 0.6-0.8, to less than 0.1 (Figure 2e).

### 3.2 Preparation of test set from recently sequenced species

We trained our model of gene family evolution on both the fly and fungal clades using only a subset of our trusted trees (1000 for flies, 300 for fungi) allowing us to use the remainder for cross-validation. In addition to the 1-to-1 ortholog gene trees found by synteny, we also identified two kinds of gene trees with duplication and loss. First, we chose 7 species surrounding the whole-genome duplication (WGD) event [17] in the yeast clade (Figure 3b). We identified syntenic regions with exactly 1 chromosomal segment from each pre-duplication species and exactly 2 segments from each post-duplication species. From these regions, we obtained 220 sets of genes syntentically aligned across the WGD. Post-duplicated genes were separated into two groups by Neighbor Joining. If paralogs failed to separate, possibly due to gene conversion, we filtered them from our test set. Since we know the phylogenetic position of the WGD event, we can conclude that the correct gene phylogeny must be the one shown in Figure 3b.

We also identified fly trees with duplications predating the divergence of the species, by determining pairs of 1-to-1 orthologous sets that had best bidirectional BLAST hits between each other. From such pairs, which we call *neighbors*, we can conclude that the full phylogeny must be the one shown in Figure 3a. To simulate gene loss events, we randomly removed genes from one orthologous group or the other. The result is a gene tree with actual divergence distances, an ancestral duplication, and simulated losses.

### 3.3 Simulating gene gain and loss in the fly phylogeny

We also constructed a simulated dataset in order to evaluate SPIDIR’s ability to reconstruct trees with many duplication and loss events. We developed our own simulator in order to generate trees with branch lengths that followed the properties found in real gene phylogenies. Since our model of gene family evolution is generative, we could simply use it along with parameters trained from real data to generate new simulated trees, which closely match actual gene evolution while also allowing us to generate a configurable number of gene duplications and losses. We ensured that the percent identity of our alignments equalled that of an average real fly alignment (36%) and that they had the same gapless alignment length ( $\sim 900$  bp). For a more detailed explanation of our simulations, see the appendix A.3.



(a)

Program	distance metric	12 Flies 1-to-1	12 Flies neighbors	17 Fungi 1-to-1	7 Yeasts WGD
SPIDIR	HKY	82%	71%	81%	82%
	JTT	88%	82%	93%	93%
PHYML+assist	HKY	36%	8%	0.8%	12%
PHYML	HKY	38%	16%	0.9%	18%
	JTT	21%	13%	0.6%	25%
MrBayes	HKY	34%	10%	0.2%	12%
	JTT	19%	5%	0.8%	24%
PHYLIP ML	HKY	35%	8%	0.2%	17%
	JTT	19%	5%	0.5%	24%
PHYLIP Parsimony	HKY	23%	14%	0%	10%
	JTT	10%	6%	0.1%	11%
BIONJ	HKY	42%	15%	0.4%	12%
	JTT	23%	7%	0.4%	23%
	dS	2%	3%	0%	1%

(b)

Program	1-to-1	1 dup	2 dup	3 dup	2 dup 1 loss	2 dup 2 loss	3 dup 2 loss
SPIDIR(nt)	100%	92%	74%	66%	78%	69%	65%
PHYML(nt)	60%	52%	36%	35%	49%	51%	47%
MrBayes(nt)	60%	45%	33%	28%	47%	49%	43%
DNAML(nt)	56%	45%	33%	29%	46%	52%	42%
DNAPARS(nt)	43%	36%	23%	21%	34%	38%	36%
HKY+BIONJ(nt)	60%	47%	35%	30%	46%	41%	42%

Figure 4: (a) Accuracy of phylogenetic reconstruction on several real datasets: 5154 alignments of 12 fly 1-to-1 orthologs, 200 alignments of fly ortholog neighbors, 739 alignments of 17 fungal 1-to-1 orthologs, 220 alignments of 7 yeast Whole-Genome-Duplication orthologs and paralogs (b) Accuracy of phylogenetic reconstruction on several simulated fly datasets with various numbers of duplication and losses events.

### 3.4 Evaluation

We compared the performance of SPIDIR to the most popular and successful phylogeny algorithms. For maximum likelihood methods, we used both PHYML v2.4.4 [10] and the PHYLIP programs DNAML/PROML [7]. Nucleotides substitutions were modeled with the HKY model and peptide substitutions were modeled with JTT. For parsimony methods, we used PHYLIP’s DNAPARS and PROTPARS programs. MrBayes v3.1.1, a bayesian-based method, was used to find the maximum a posterior phylogenetic tree [24]. We ensured MrBayes reported the most likely binary tree. Lastly, we used the Neighbor-Joining program BIONJ [8] with a variety of distance matrices, including HKY built with PUZZLE-TREE [26], JTT built with PROTDIST, and dS built with the YN00 program from PAML v3.15 [30]. For most programs, default options were used.

We found that SPIDIR significantly out-performs all existing phylogenetic methods on each of our tests for both real and simulated datasets (Figure 4).

The accuracy of other methods appears to drop significantly as the number of genes increases, whereas SPIDIR’s accuracy remained high. One factor is certainly that the number of possible trees increases exponentially with the number genes. However, inaccuracies are not due solely to the tree search: we found that assisting the search by providing the correct tree did not increase accuracy (*PHYML+assist* row). Even when PHYML was initialized on the correct tree topology, it only chose it 0.8% of the time in the 17-fungi dataset, suggesting that erroneous trees in fact obtained higher likelihood scores. For each method we also ensured that the correct tree was still the most frequently inferred tree, despite its very low frequency.

SPIDIR’s high accuracy in datasets with many duplications and losses indicates that our model is able accurate model duplicated topologies which significantly differ from the learned species tree, in both the fly neighbor and WGD datasets (Figure 4a), as well as 7 simulated datasets with increasing numbers of duplication and loss events (Figure 4b). Note that only a small number of duplications early in the tree can lead to large numbers of duplicated genes, up to 30 in our datasets.

Finally, we found that these results were consistent using both nucleotide and protein alignments. We also used silent mutations  $dS$  as a distance metric, often praised for representing a neutral rate of mutation. However, in our datasets performance dropped, as they may be too divergent for the metric to perform accurately.

## 4 Contributions

We have introduced a new method of phylogenetic reconstruction that successfully addresses the phylogenomic problem. Our approach differs significantly from previous methods, in that we view phylogenetic reconstruction as a learning problem. By observing the evolution of a small portion of gene families with unambiguous phylogeny, we learn a model for gene and species evolution, which guides the phylogenetic reconstruction of all gene families across multiple genomes.

To achieve this, we have developed a statistical model which captures two evolutionary forces of gene family evolution: gene-specific rates, likely to come from the selective pressures on individual gene functions, and species-specific mutations rates, likely to come from the effective population size and dynamics of a species. This model has the advantage to both closely fit the evolution of real gene phylogenies, and lend itself to efficient computation for rapid and accurate reconstruction of thousands of phylogenetic trees.

We have also presented a training procedure to estimate our model parameters and showed a practical means of constructing training and test datasets by using synteny information, which is strongly informative of gene ancestry, and also orthogonal to gene mutation rates in phylogenetic reconstruction.

We formulated an efficient likelihood function for our model, which can evaluate large phylogenies in time linear in the number of genes. Coupled with an efficient tree search, we implemented our likelihood function as maximum likelihood algorithm, which we have made available in the software package SPIDIR.

Lastly, we extensively evaluated our method in comparison to the most successful phylogenetic algorithms on a wide variety of real and simulated datasets, and found that it largely outperforms existing methods which cannot incorporate or learn a species tree in their modeling.

Our model of gene family evolution also provides a new framework for addressing the long standing question of how to relate species and gene phylogeny. By normalizing a gene tree by its family rate, the relative evolutionary rates of different species became apparent, and we were able to learn species-specific rates for each branch.

We expect that this method will be applicable to many diverse species, and in particular the upcoming 32 mammalian genomes, and the multitude of fungal genomes important in health and biotechnology. We tested the model extensively in the two largest eukaryotic datasets currently available, 12 fly genomes and 17 fungal genomes, and showed that our modeling of gene mutation rate correlation closely fits gene phylogenies that span over 150 million years of evolution.

As technological advances continue to drive sequencing costs down, each biologically important clade of life is likely to be rich with multiple closely related species for comparative sequence analysis. The demand for rigorous and practical comparative methods will only increase with such datasets, and the phylogenomics approach developed here is likely to play a significant role in the comparison and understanding of these species.

## 5 Acknowledgements

The authors would like to thank Mike Lin, Alex Stark, Joshua Grochow, Radek Sklarczyk for help, advice and comments. We are indebted to the fly community, and in particular Doug Smith, Bill Gelbart, Mike Eisen, and Thom Kaufman, for early release of data and annotations, and permission to use in our benchmarks. We also thank the NIAID, the Sanger Center, and our collaborators at the Broad Institute, and especially Bruce Birren and Christina Cuomo, for discussions and permission to use fungal genomes in our benchmarks. One of us (MDR) was partially funded by the Ruth L. Kirschstein National Research Service Award (NRSA).

## A Derivations

### A.1 Estimating the gene family rate

Let  $F$  be the family rate. For each branch  $i$ , we have branch length  $x_i$ , species specific rate  $S_i \sim N(\mu_i, \sigma_i^2)$ .

Given several  $x_i$  we want to find maximum likelihood of  $F$ .

$$\begin{aligned}
 X_i &= S_i * F \\
 P(X_i = x_i | \mu_i, \sigma_i) &= \text{normal}(x_i | \mu_i, \sigma_i^2) \\
 F &= X_i / S_i \quad \forall i \\
 P(F = f | \alpha, \beta) &= \text{gamma}(f | \alpha, \beta) \\
 P(F = f | x_i, \mu_i, \sigma_i \quad \forall i, \alpha, \beta) &= P(x_1 / S_1 = f, \dots, x_n / S_n = f) * P(F = f | \alpha, \beta)
 \end{aligned}$$

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

$$\begin{aligned}
 P(F = f | x_i, \mu_i, \sigma_i \quad \forall i) &= P(F = f) \prod_i P(x_i / S_i = f) \\
 &= P(F = f) \prod_i P(S_i = x_i / f) \\
 &= \text{gamma}(f | \alpha, \beta) \prod_i \text{normal}(x_i / f | \mu_i, \sigma_i^2)
 \end{aligned}$$

The maximum likelihood estimator of  $F$  given  $\alpha$ ,  $\beta$ ,  $x_i$ ,  $\mu_i$ , and  $\sigma_i$  for all  $i$  is then

$$\begin{aligned}
 \arg \max_b P(F = f | x_i, \mu_i, \sigma_i \quad \forall i, \alpha, \beta) &= \arg \max_f \text{gamma}(f | \alpha, \beta) \prod_i \text{normal}(x_i / f | \mu_i, \sigma_i^2) \\
 &= \arg \max_f f^{\alpha-1} e^{-\beta f} \frac{\beta^\alpha}{\Gamma(\alpha)} \prod_i \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left(-\frac{(\frac{x_i}{f} - \mu_i)^2}{2\sigma_i^2}\right)
 \end{aligned}$$

Applying a natural log gives

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

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

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

Which gives the following cubic polynomial

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

## A.2 Conditional branch likelihood

First consider the following tree (Figure 1 case 2). Let  $a_1$  and  $a_2$  be two paralogs. Their immediate ancestor  $u$  represents a duplication event, and the root of subtree  $v$  represents a speciation event. Both paths  $v \rightarrow a_1$  and  $v \rightarrow a_2$  reconcile to the same species branch which has parameters  $\mu$  and  $\sigma^2$ . Let the lengths of the paths ending at  $a_1$  and  $a_2$  be random variables  $X_1$  and  $X_2$ , respectively. Therefore, we need to compute the following likelihood:

$$P(X_1 = x_0 + x_1, X_2 = x_0 + x_2) = P(X_2 = x_0 + x_2 | X_1 = x_0 + x_1) P(X_1 = x_0 + x_1)$$

Conditioning  $X_2$  on  $X_1$  imposes that  $X_2$  be greater than the shared branch length,  $x_0$ . (assuming branch lengths are non-negative). Thus, the likelihood of  $X_2$  is simply scaled by the probability that  $X_2 > x_0$ . We can thus write:

$$P(X_2 = x_0 + x_2 | X_1 = x_0 + x_1) P(X_1 = x_0 + x_1) = P(X_2 = x_0 + x_2 | X_2 > x_0) P(X_1 = x_0 + x_1) \quad (6)$$

$$= \frac{P(X_2 = x_0 + x_2, X_2 > x_0)}{P(X_2 > x_0)} P(X_1 = x_0 + x_1) \quad (7)$$

$$= \frac{P(X_2 = x_0 + x_2)}{P(X_2 > x_0)} P(X_1 = x_0 + x_1) \quad (8)$$

$$= \frac{N(x_0 + x_2 | \mu, \sigma^2) N(x_0 + x_1 | \mu, \sigma^2)}{C(x_0 | \mu, \sigma^2)} \quad (9)$$

Where  $C(x_0 | \mu, \sigma^2) = \int_{x_0}^{\infty} N(a | \mu, \sigma^2) da$  is the cumulative normal distribution which also has a closed form solution. Notice that the final likelihood does not depend on which path we condition upon. Either way, the likelihood is a product of each path's likelihood divided by the probability of seeing a branch longer than  $x_0$ .

## A.3 Simulation

The simulation takes as input a species tree and our species specific mutation model, which consists of a normal distribution for each species branch. The simulation begins with a random sequence at

the root of the species tree. The root sequence is chosen by concatenating randomly selected codons according to the codon usage found in *D. melanogaster*. The sequence then evolves down each branch of the species tree.

For each species branch, the simulation chooses a gene branch length from the normal distribution for that species. Given the branch length, the sequence is mutated according to a codon substitution model that was specifically trained for the fly genomes. This process repeats until the evolving sequences reach the leaves. For simplicity, inserts and deletes are not simulated.

To simulate duplication and loss, we use two parameters for how many duplications  $D$  or losses  $L$  to create. Duplications are randomly assigned to  $D$  species tree branches. When a sequence reaches a branch assigned with a duplication, it chooses a duplication midpoint  $k$  between 0 and 1 uniformly, which represents how far down the species branch the sequence should evolve before it duplicates. Mathematically, this done by choosing a length from a normal with mean and variance scaled by  $k$ . When the sequence reaches the duplication point, two new branches are created that recursively follow the same process. If more duplication midpoints are needed, they are chosen between  $k$  and 1. After a tree is fully simulated,  $L$  branches are randomly picked to be lost. These branches along with their descendants are removed. Also any leaves that are not reconciled to modern species are also pruned. The resulting tree is a simulated gene tree with at most  $D$  duplications and  $L$  losses.

## References

- [1] Lars Arvestad, Ann-Charlotte Berglund, Jens Lagergren, and Bengt Sennblad. Gene tree reconstruction and orthology analysis based on an integrated model for duplications and sequence evolution. *RECOMB*, 2004.
- [2] David Bryant and Peter Waddell. Rapid evaluation of least-squares and minimum-evolution criteria on phylogenetic trees. *Mol. Biol. Evol.*, 15(10):1346–1359, 1998.
- [3] Andrew Clark, Greg Gibson, Thom Kaufman, Eugene Myers, and Patrick OGrady. Proposal for drosophila as a model system for comparative genomics. <http://flybase.net/.data/docs/communitywhitepapers/genomeswp2003.html>.
- [4] J. F. Dufayard, L. Duret, S. Penel, M. Gouy, F. Rechenmann, and G. Perriere. Tree pattern matching in phylogenetic trees: automatic search for orthologs or paralogs in homologous gene sequence databases. *Bioinformatics*, 21:2596–603, June 2005.
- [5] Robert C. Edgar. Muscle: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Research*, 32(5):1792–1797, 2004.
- [6] Jonathan A. Eisen. Phylogenomics: improving functional predictions for uncharacterized. *Genome Research*, pages 8(3):163–7, 1998.
- [7] J Felsenstein. Phylip (phylogeny inference package) version 3.6. *Distributed by the author. Department of Genome Sciences, University of*, 2005.
- [8] Olivier Gascuel. Bionj: an improved version of the nj algorithm based on a simple model of sequence data. *Molecular Biology and Evolution*, 14:685–695, 1997.
- [9] M. Goodman, J. Czelusniak, G. W. Moore, and A. E. Romero-Herrera G. Matsuda. Fitting the gene lineage into its species lineage, a parsimony strategy illustrated by cladograms constructed from globin sequences. *Syst. Zool.*, 28:132–168, 1979.
- [10] Stephane Guindon and Olivier Gascuel. A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. *Systematic Biology*, 52, 2003.
- [11] Ashbya Genome Database. <http://agd.unibas.ch/>.
- [12] Genolevures. <http://cbi.labri.fr/Genolevures/>.
- [13] Washington University Genome Sequencing Center. <http://genome.wustl.edu/projects/yeast/>.
- [14] Broad Institute. Fungal Genome Initiative. <http://www.broad.mit.edu/annotation/cgi/>.
- [15] Sanger Candida Project. <http://www.sanger.ac.uk/sequencing/Candida/>.
- [16] Saccharomyces Genome Database. <http://www.yeastgenome.org/>.
- [17] Manolis Kellis, Bruce Birren, and Eric Lander. Proof and evolutionary analysis of ancient genome duplication in the yeast *saccharomyces cerevisiae*. *Nature*, April 2004.
- [18] Heng Li, Avril Coghlan, Jue Ruan, Lachlan James Coin, Jean-Karim Heriche, Lara Osmotherly, Ruiqiang Li, Tao Liu, Zhang Zhang, Lars Bolund, Gane Ka-Shu Wong, Weimou Zheng, Paramvir Dehal, Jun Wang, and Richard Durbin. Treefam: a curated database of phylogenetic trees of animal gene families. *Nucleic Acids Research*, 34, 2006.
- [19] Li Li, Christian J Jr Stoeckert, and David S Roos. Orthomcl: identification of ortholog groups for eukaryotic genomes. *Genome Res*, 13(9):2178–2189, September 2003.
- [20] Masatoshi Neidagger, Sudhir Kumar, and Kei Takahashi. The optimization principle in phylogenetic analysis tends to give incorrect topologies when the number of nucleotides or amino acids used is small. *Proc Natl Acad Sci*, 95:12390–12397, October 1998.
- [21] R.D.M. Page and M. A. Charleston. Reconciled trees and incongruent gene and species trees. *Mathematical Hierarchies in Biology, DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, 37:57–70, 1997.
- [22] M Remm, CEV Storm, and Sonnhammer ELL. Automatic clustering of orthologs and in-paralogs from pairwise species comparisons. *JMB*, 314:1041–1052, 2001.

- [23] A. Rokas, B. L. Williams, N. King, and S. B. Carroll. Genome-scale approaches to resolving incongruence in molecular phylogenies. *Nature*, 425:798–804, 2003.
- [24] Huelsenbeck JP. Ronquist F. Mrbayes 3: Bayesian phylogenetic inference under mixed models. *Bioinformatics*, pages 19(12):1572–4, 2003 Aug 12.
- [25] N. Saitou and T. Imanishi. Relative efficiencies of the fitch-margoliash, maximum-parsimony, maximum-likelihood, minimum-evolution, and neighbor-joining methods of phylogenetic tree construction in obtaining the correct tree. *Mol. Biol. Evol.*, 6:514–525, 1989.
- [26] H. A. Schmidt, K. Strimmer, M. Vingron, and A von Haeseler. Tree-puzzle: maximum likelihood phylogenetic analysis using quartets and parallel computing. *Bioinformatics*, 18, 2002.
- [27] Christian E.V. Storm and Erik L.L. Sonnhammer. Comprehensive analysis of orthologous protein domains using the hops database. 2003.
- [28] Y. Tateno, N. Takezaki, and M. Nei. Relative efficiencies of the maximum-likelihood, neighbor-joining, and maximum-parsimony methods when substitution rate varies with site. *Mol. Biol. Evol.*, 11:261–277, 1994.
- [29] Roman L Tatusov, Natalie D Fedorova, John D Jackson, Aviva R Jacobs, Boris Kiryutin, Eugene V Koonin, Dmitri M Krylov, Raja Mazumder, Sergei L Mekhedov, Anastasia N Nikolskaya, B Sridhar Rao, Sergei Smirnov, Alexander V Sverdlov, Sona Vasudevan, Yuri I Wolf, Jodie J Yin, and Darren A Natale. The cog database: an updated version includes eukaryotes. *BMC Bioinformatics*, 4:41, September 2003.
- [30] Z. Yang. Paml: a program package for phylogenetic analysis by maximum likelihood. *Comput Appl Biosci*, 13, October 1997.
- [31] Christian M Zmasek and Sean R Eddy. Rio: analyzing proteomes by automated phylogenomics using resampled inference of orthologs. *BMC Bioinformatics*, 3:14, May 2002.
- [32] Eddy SR. Zmasek CM. A simple algorithm to infer gene duplication and speciation events. *Bioinformatics*, 2001.