

6.047/6.878 Lecture 21: Phylogenomics II

Guest Lecture by
Matt Rasmussen
Orit Giguzinsky and Ethan Sherbondy

November 15, 2012

Contents

1	Introduction	3
2	Inferring Orthologs/Paralogs, Gene Duplication and Loss	4
2.1	Species Tree	4
2.2	Gene Tree	4
2.3	Gene Family Evolution	4
2.4	Reconciliation	5
2.4.1	Definitions	5
2.4.2	Maximum Parsimony Reconciliation (MPR) algorithm	5
2.4.3	Reconciliation Examples	6
3	Reconstruction	8
3.1	Species Tree Reconstruction	8
3.2	Improving Gene Tree Reconstruction and Learning Across Gene Trees	8
4	Modeling Population and Allele Frequencies	9
4.1	The Wright-Fisher Model	9
4.2	The Coalescent Model	10
5	SPIDIR:Background	12
6	SPIDIR: Method and Model	13
7	Conclusion	14
8	Current Research Directions	14
9	Further Reading	14
10	Tools and Techniques	14
11	What Have We Learned?	14

List of Figures

1	Species Tree	4
2	Gene Tree	4
3	Gene Tree Inside a Species Tree	4
4	Gene Family Evolution: Gene Trees and Species Trees	5
5	Mapping Diagram	5
6	Nesting Diagram	5
7	Maximum Parsimony Reconciliation (MPR)	6
8	Maximum Parsimony Reconciliation Recursive Algorithm	6
9	Reconciliation Example 1, simple mapping case	6
10	Reconciliation Example 2, parsimonious reconciliation for complex case	7
11	Reconciliation Example 3, non parsimonious reconciliation for complex case	7
12	Reconciliation Example 4, invalid Reconciliation	7
13	Species Tree Reconstruction	8
14	Using species trees to improve gene tree reconstruction.	8
15	We can develop a model for what kind of branch lengths we can expect. We can use conserved gene order to tell orthologs and build trees.	9
16	Branch length can be modeled as two different rate components: gene specific and species specific.	9
17	The Wright-Fisher model	10
18	The Wright-Fisher model continued over many generations and ignoring the ordering of chromosomes.	11
19	The coalescent model.	11
20	Geometric probability distribution for coalescent events in k lineages.	11
21	Multispecies Coalescent Model. Leaf branches track one lineage. There is a lag time from when population separated and when two actual gene lineages find a common ancestor. The rate of coalescent slows down as N gets bigger and for short branches. Deep coalescent is depicted in light blue for three lineages. The species and gene tree are incongruent since C and D are sisters in gene tree but not the species tree. There is a $\frac{2}{3}$ chance that incongruence will occur because once we get to the light blue section, the Wright-fisher is memory less and there is only $\frac{1}{3}$ chance that it will be congruent. Effect of incongruence is called incomplete lineage sorting.	12
22	MPR reconciliation of genes and species tree.	13
23	Inaccuracies in gene tree.	13

1 Introduction

In the previous chapter, we covered techniques for reasoning about evolution in terms of trees of descent. The algorithms we covered for tree-building, UPGMA and neighbor-joining, assumed that we were comparing fully aligned sections of sequences.

In this section, we present additional models for using phylogenetic trees in different contexts. Here we clarify the differences between species and gene trees. We then cover a framework called reconciliation which lets us effectively combine the two by mapping gene trees onto species trees. This mapping gives us a means of inferring gene duplication and loss events.

We will also present a phylogenetic perspective for reasoning about population genetics. Since population genetics deals with relatively recent mutation events, we offer the Wright-Fisher model as a tool for representing changes in whole populations. Unfortunately, when dealing with real-world data, we usually are only able to sequence genes from the current living descendants of a group. As a remedy to this shortcoming, we cover the Coalescent model, which you can think of as a time-reversed Wright-Fisher analog.

By using coalescence, we gain a new means for estimating divergence times and population sizes across multiple species. At the end of the chapter, we touch briefly on the challenges of using trees to model recombination events and summarize recent work in the field along with frontiers open for exploration.

2 Inferring Orthologs/Paralogs, Gene Duplication and Loss

There are two commonly used trees, Species tree and Gene tree. This section explains how these trees can be used and how to fit a gene tree inside a species tree (reconciliation).

2.1 Species Tree

Species trees that show how different species evolved from one another. These trees are created using morphological characters, fossil evidence, etc. The leaves of each tree are labeled as species and the rest of the tree shows how these species are related. An example of a species tree is shown in Figure 1.

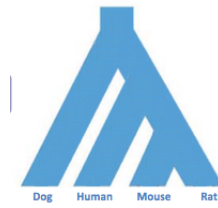


Figure 1: Species Tree

2.2 Gene Tree

Gene trees are trees that look at specific genes in different species (leaves are genes). The leaves of gene trees are labeled with gene sequences or gene ids associated with specific sequences. Figure 2 shows an example of a gene tree that has 4 genes (leaves). The sequences associated with each gene are presented on the right side of Figure 2.

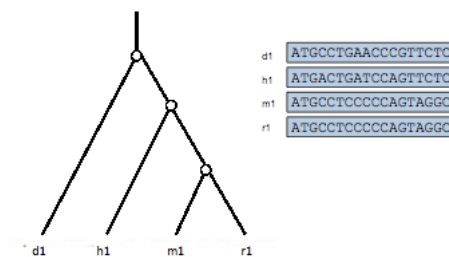


Figure 2: Gene Tree

2.3 Gene Family Evolution

Gene trees evolve inside a species tree. An example of a gene tree contained in a species tree is shown in Figure 3 below.

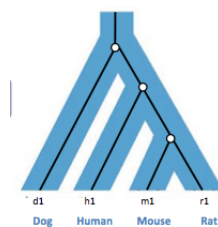


Figure 3: Gene Tree Inside a Species Tree

The next sub section explains how we can fit gene trees inside a species trees using Reconciliation.

2.4 Reconciliation

Reconciliation is an algorithm that helps compare gene trees to genome trees by fitting a gene tree fits inside a species tree. This is done by by mapping the vertices in the gene tree to vertices in the species tree. This sub section will focus on Reconciliation, related definitions, algorithm (Maximum Parsimony Reconciliation algorithm) and examples.

2.4.1 Definitions

Two genes are **orthologs** if their recent common ancestor (MRCA) is a speciation (splitting into different species).

Paralogs are genes whose MRCA is a duplication.

Figure 4 below illustrates how these types of genes can be represented in a gene tree. The tree below has 4 speciation nodes, one duplication and one loss.

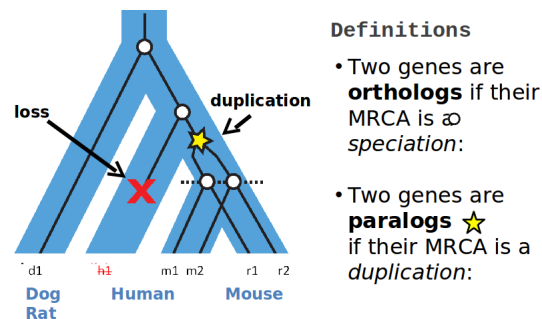


Figure 4: Gene Family Evolution: Gene Trees and Species Trees

A mapping diagram is a diagram that shows the node mapping from the gene tree to the species tree. Figure 5 shows an example of a mapping diagram.

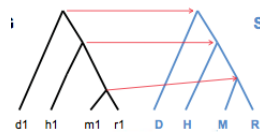


Figure 5: Mapping Diagram

A nesting diagram shows how the gene tree can be nested inside the species tree. For every mapping diagram there is a nesting diagram. Figure 6 shows an example of a possible nesting diagram for the mapping diagram in Figure 5.

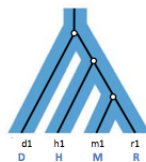


Figure 6: Nesting Diagram

2.4.2 Maximum Parsimony Reconciliation (MPR) algorithm

MPR is an algorithm that fits a gene tree in a species tree while minimizing the number of duplications and deletions.

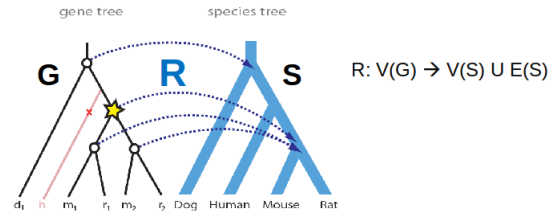


Figure 7: Maximum Parsimony Reconciliation (MPR)

Given a gene tree and a species tree, the algorithm finds the reconciliation that minimizes the number of duplications and deletions. Figure 7 above shows an example of a possible mapping from a gene tree to a species tree. Figure 8 presents the pseudocode for the MPR algorithm.

Solve recursively:

- $R[v] = \text{species of } v$ if $v \in L(G)$
- $R[v] = \text{LCA}(R[\text{right}(v)], R[\text{left}(v)])$ if $v \in I(G)$
- LCA = “least common ancestor”
(also called “most recent common ancestor”)

Labeling events:

- v is a dup if $R[v] = R[\text{right}(v)]$ or $R[\text{left}(v)]$
- Branch above v has at least one loss if $R[\text{parent}(v)] \neq R[v]$ or $\text{parent}[R[v]]$

Figure 8: Maximum Parsimony Reconciliation Recursive Algorithm

We map the arrows low as possible, since lower mapping usually results in fewer events. However, we cannot map too low. We map as low as we can without violating the descendent-ancestor relationships. The algorithm goes recursively from bottom up, starting from the leaves. We already know the mapping for the leaves, so we can easily map them. To map the ancestors, for each node (going recursively up the tree) we look at the right child and left child and take the least common ancestor (LCA) of the species that they map to. If a node maps to its right or left child, we know there is a duplication. An expected branch that does not exist indicates a loss.

2.4.3 Reconciliation Examples

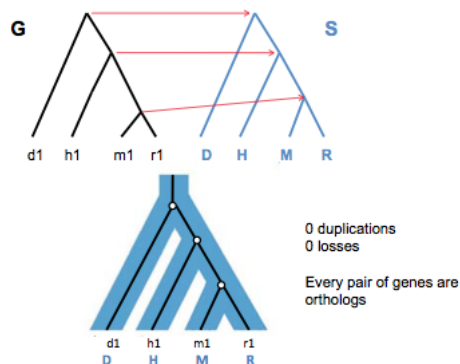


Figure 9: Reconciliation Example 1, simple mapping case

In Figure 9, the nodes can be mapped straight across, since there are no duplications or losses.

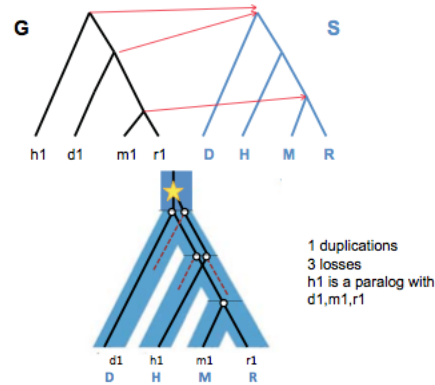


Figure 10: Reconciliation Example 2, parsimonious reconciliation for complex case

In Figure 10, we see a parsimonious (minimum number of losses and duplications) reconciliation for a case in which nodes from the gene tree cannot be mapped straight across.

does a non-parsimonious reconciliation look like?

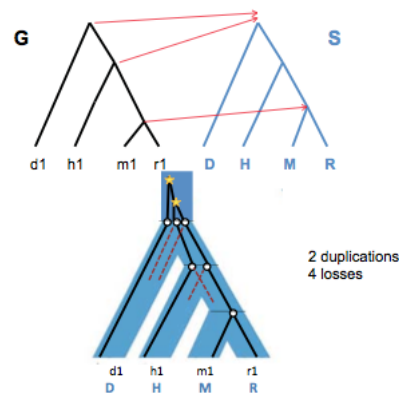


Figure 11: Reconciliation Example 3, non parsimonious reconciliation for complex case

Figure 11 shows a non-parsimonious reconciliation. The parsimonious mapping for the same trees is shown in Figure 9.

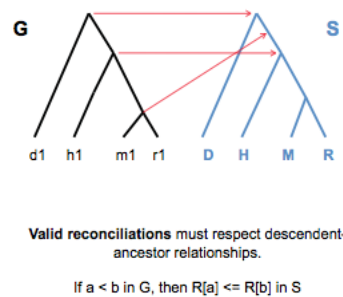


Figure 12: Reconciliation Example 4, invalid Reconciliation

Figure 12 shows an invalid reconciliation. This reconciliation is invalid since it does not respect descent-ancestor relationships. In order for this reconciliation to be possible, the descendent would have to travel

back in time and be created before its ancestor. Clearly, such a scenario would be impossible. A valid reconciliation must satisfy the following: **If $a < b$ in G , then $R[a] \leq R[b]$ in S .**

3 Reconstruction

In the previous section we learned how to compare gene trees and species trees. In this section, we will use this information to reconstruct gene trees and species trees.

3.1 Species Tree Reconstruction

In the past, it was really hard to identify a marker gene for a specific species. As sequencing improved we started having lots of sequencing data, people started building trees for different loci. The tree you got highly dependent on the tree you used. Possible reasons why trees differ include noise (from statistical estimate errors and noise), hidden duplications and losses and allele sorting in a population.

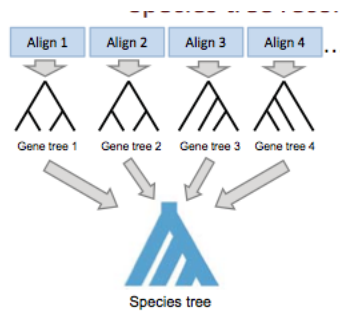


Figure 13: Species Tree Reconstruction

Given lots of different gene trees that disagree, our goal is to make them into once species tree (as shown in Figure 13). There are lots of different algorithms that reconstruct species trees. These algorithms include Supermatrix methods (Rokas 2003, Ciccareli 2006), Supertree methods (Creevey & McInerney 2005), Minimizing Deep Coalescence (Maddison & Knowles 2006) and Modeling coalescence (Liu & Pearl 2007).

3.2 Improving Gene Tree Reconstruction and Learning Across Gene Trees

We can use methods similar to those described above to build better gene trees. This can be done by using information from a species tree to study a gene tree of interest. For example, species trees can be used to determine when losses and duplications occurred. The idea is that we can use the fact that species trees are often built from the entire genome, to obtain more information about related gene trees.

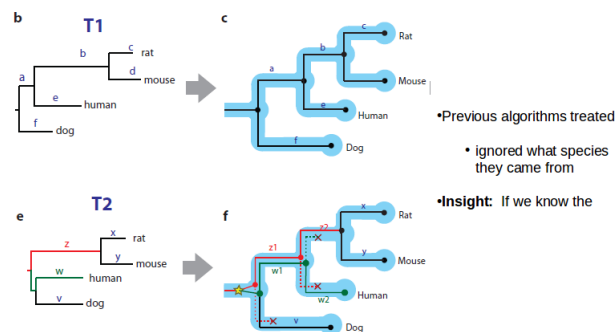


Figure 14: Using species trees to improve gene tree reconstruction.

If we know the species tree, we can develop a model for what kind of branch lengths we can expect. We can use conserved gene order to tell orthologs and build trees.

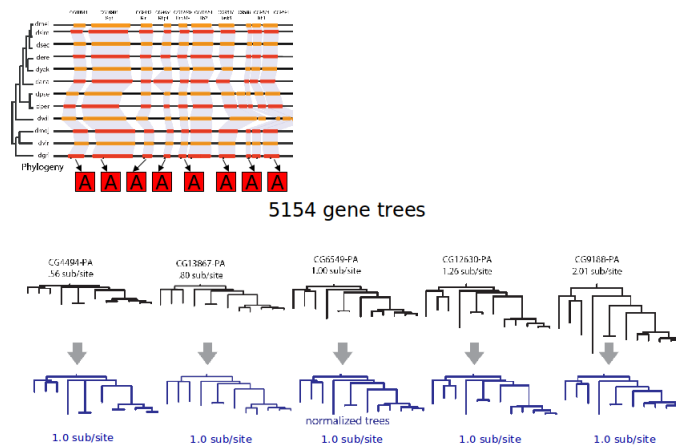


Figure 15: We can develop a model for what kind of branch lengths we can expect. We can use conserved gene order to tell orthologs and build trees.

When gene is fast evolving in one species, it is fast evolving in all species. We can model a branch length as two different rate components. One is gene specific (present across all species) and a species specific which is customized to a specific species.

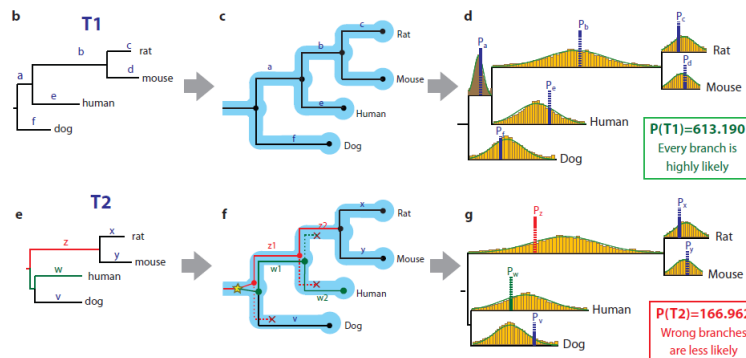


Figure 16: Branch length can be modeled as two different rate components: gene specific and species specific.

This method greatly improves reconstruction accuracy.

4 Modeling Population and Allele Frequencies

With the advent of next-gen sequencing, it is becoming economical to sequence the genomes of many individuals within a population. In order to make sense of how alleles spread through a population, it's helpful to have a model to compare data against. The **Wright-Fisher** reproduction model has filled this role for the past 70 years.

4.1 The Wright-Fisher Model

Like HMMs, Wright-Fisher is a Markov process: at each step, the system randomly progresses, and the current state of the system depends only on the previous state. In this case, state transitions represent reproduction. By modeling the transmission of chromosomes to offspring, we can study genetic drift.

The model makes a number of simplifying assumptions:

1. Population size, N , is constant at each generation.
2. Only members of the same generation reproduce (no overlap).
3. Reproduction occurs at random.
4. The gene being modeled only has 2 alleles.
5. Genes undergo neutral selection.

Note that Wright-Fisher is not an appropriate choice if you're trying to model the change in frequency of a gene that is positively or negatively selected for. If we use Wright-Fisher to model the chromosomes of diploid individuals, the population size of the model becomes $2N$.

In English, here's how Wright-Fisher works:

At every generation, for each child, we randomly select from the parents (with replacement). The allele of the child becomes that of the randomly selected parent.

We repeat this process for many generation, with the children serving as the new parents, ignoring the ordering of chromosomes.

It really is that simple. To determine the probability of k copies of an allele existing in the child generation when it had a frequency of p in the parent generation, we can use this formula:

$$\binom{2N}{k} p^k q^{2N-k} \quad (1)$$

Here, $q = (1 - p)$. It is the frequency of non- p alleles in the parent generation.

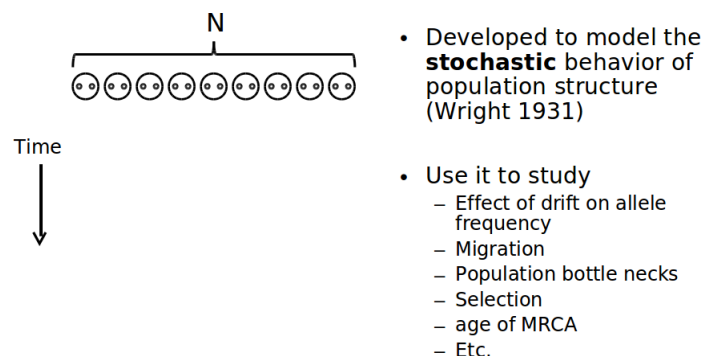


Figure 17: The Wright-Fisher model

Now we can begin to explore such questions as: how probable is it and how many generations is it expected to take for a given allele to become **fixed**, meaning the allele is present in *every* member of the population?

The expected time (in generations) for fixation, given the assumptions made by Wright-Fisher, is proportionate to $4N_E$, where N_E is the effective population size.

Again, it's important to keep in mind the limitations of this model and ask if it actually makes sense for the system you're trying to represent. Consider how you could tweak the proposed model to account for a selection coefficient ranging between -1 (lethal negative selection) and 1 (strong positive selection).

4.2 The Coalescent Model

The coalescent model only focuses on the genealogy. It only is concerned about the lineages we have sequences for; do not have to worry about others. It is a probabilistic model that works backwards in time to find when they have common ancestors.

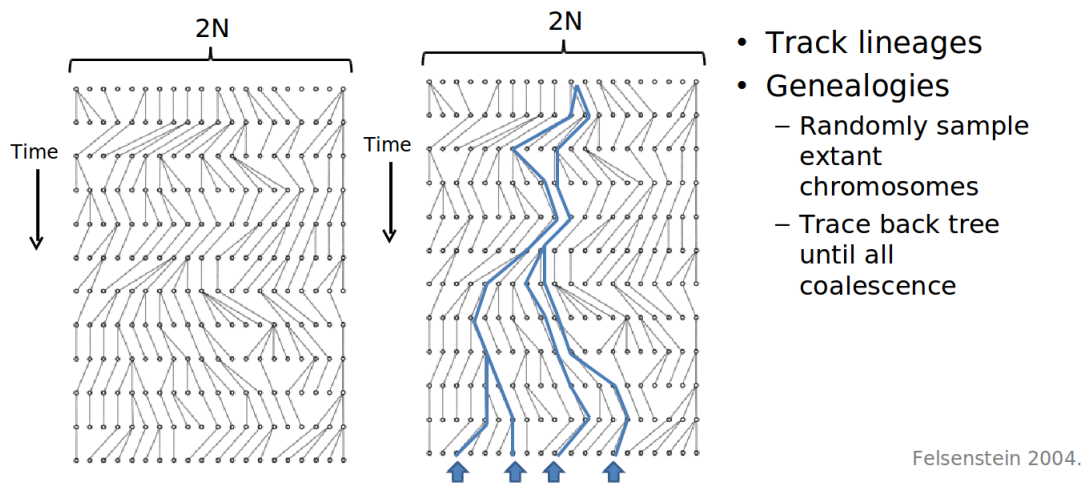


Figure 18: The Wright-Fisher model continued over many generations and ignoring the ordering of chromosomes.

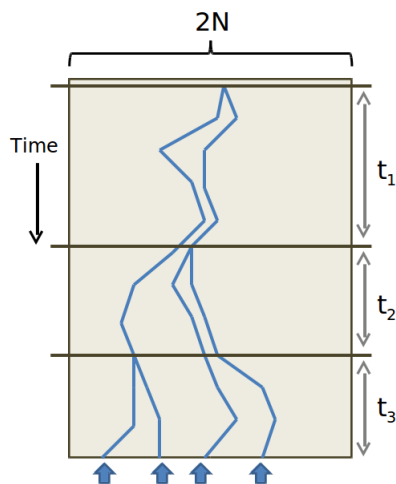


Figure 19: The coalescent model.

$$\begin{aligned}
 & \frac{(2N-1)}{2N} \frac{(2N-2)}{2N} \dots \frac{(2N-k+1)}{2N} = \prod_{i=1}^{k-1} \left(1 - \frac{i}{2N}\right) \\
 & = 1 - \sum_{i=1}^{k-1} \frac{i}{2N} + o\left(\frac{1}{N^2}\right) = 1 - \binom{k}{2} \frac{1}{2N} + o\left(\frac{1}{N^2}\right),
 \end{aligned}$$

For $k \ll N$, $O(1/N^2)$ is very small

Figure 20: Geometric probability distribution for coalescent events in k lineages.

Say we have $2N$ individuals, what is the probability that k lineages do not have any coalescent events in parental generation? What is the probability that the first coalescent of k lineages is at t generations? This process can be seen as a geometric distribution. Can repeat to find when all individuals coalesce. Each branch of species tree can be seen as having its own Wright-Fisher inside of it.

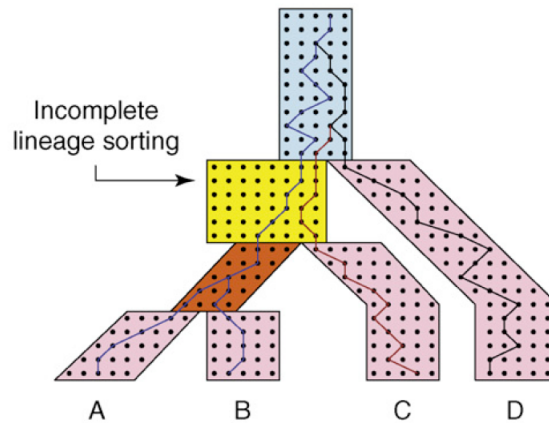


Figure 21: Multispecies Coalescent Model. Leaf branches track one lineage. There is a lag time from when population separated and when two actual gene lineages find a common ancestor. The rate of coalescent slows down as N gets bigger and for short branches. Deep coalescent is depicted in light blue for three lineages. The species and gene tree are incongruent since C and D are sisters in gene tree but not the species tree. There is a $\frac{2}{3}$ chance that incongruence will occur because once we get to the light blue section, the Wright-fisher is memory less and there is only $\frac{1}{3}$ chance that it will be congruent. Effect of incongruence is called incomplete lineage sorting.

5 SPIDIR:Background

As presented in the supplementary information for SPIDIR, a gene family is the set of genes that are descendents of a single gene in the most recent common ancestor (MRCA) of all species under consideration. Furthermore, genetic sequences undergo evolution at multiple scales, namely at the level of base pairs, and at the level of genes. In the context of this lecture, two genes are orthologs if their MRCA is a speciation event; two genes are paralogs if their MRCA is a duplication event.

In the genomic era, the species of a modern genes is often known; ancestral genes can be inferred by reconciling gene- and species-trees. A reconciliation maps every gene-tree node to a species-tree node. A common technique is to perform Maximum Parsimony Reconciliation (MPR), which finds the reconciliation R implying the fewest number of duplications or losses using the recursion over inner nodes v of a gene tree G . MPR first maps each leaf of the gene tree to the corresponding species leaf of the species tree. Then the internal nodes of G are mapped recursively:

$$R(v) = MRCA(R(right(v)), R(left(v)))$$

If a speciation event and its ancestral node are mapped to the same node on the species tree. Then the ancestral node must be an duplication event.

Using MPR, the accuracy of the gene tree is crucial. Suboptimal gene trees may lead to an excess of loss and duplication events. For example, if just one branch is misplaced (as in ??) then reconciliation infers 3 losses and 1 duplication event. In [6], the authors show that the contemporaneous current gene tree methods perform poorly (60% accuracy) on single genes. But if we have longer concatenated genes, then accuracy may go up towards 100%. Furthermore, very quickly or slowly evolving genes carry less information as compared with moderately diverging sequences (40-50% sequence identity), and perform correspondingly worse. As corroborated by simulations, single genes lack sufficient information to reproduce the correct species tree. Average genes are too short and contains too few phylogenetically informative characters.

While many early gene tree construction algorithms ignored species information, algorithms like SPIDIR capitalize on the insight that the species tree can provide additional information which can be leveraged for gene tree construction. Synteny can be used to independently test the relative accuracy of different gene tree reconstructions. This is because syntenic blocks are regions of the genome where recently diverged organisms have the same gene order, and contain much more information than single genes.

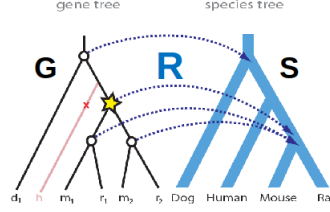


Figure 22: MPR reconciliation of genes and species tree.

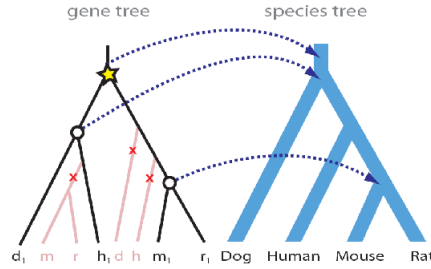


Figure 23: Inaccuracies in gene tree.

There have been a number of recent phylogenomic algorithms including: RIO [2], which uses neighbor joining (NJ) and bootstrapping to deal with incongruencies, Orthotrappier [7], which uses NJ and reconciles to a vague species tree, TreeFAM [3], which uses human curation of gene trees as well as many others. A number of algorithms take a more similar track to SPIDIR [6], including [4], a probabilistic reconciliation algorithm [8], a Bayesian method with a clock [9], and parsimony method using species tree, as well as more recent developments: [1] a Bayesian method with relaxed clock and [5], a Bayesian method with gene and species specific relaxed rates (an extension to SPIDIR).

6 SPIDIR: Method and Model

SPIDIR exemplifies an iterative algorithm for gene tree construction using the species tree. In SPIDIR, the authors define a generative model for gene-tree evolution. This consists of a prior for gene-tree topology and branch lengths. SPIDIR uses a birth and death process to model duplications and losses (which informs the prior on topology) and then then learns gene-specific and species-specific substitution rates (which inform the prior on branch lengths). SPIDIR is a *Maximum a posteriori* (MAP) method, and, as such, enjoys several nice optimality criteria.

In terms of the estimation problem, the full SPIDIR model appears as follows:

$$\operatorname{argmax}_{L, T, R} P(L, T, R | D, S, \Theta) = \operatorname{argmax}_{L, T, R} P(D | T, L) P(L | T, R, S, \Theta) P(T, R | S, \Theta)$$

The parameters in the above equation are: D = alignment data, L = branch length, T = gene tree topology, R = reconciliation, S = species tree (expressed in times), Θ = (gene and species specific parameters [estimated using EM training], λ , μ dup/loss parameters). This model can be understood through the three terms in the right hand expression, namely:

1. the sequence model– $P(D | T, L)$. The authors used the common HKY model for sequence substitutions,

which unifies Kimura's two parameter model for transitions and transversions with Felsenstein's model where substitution rate depends upon nucleotide equilibrium frequency.

2. the first prior term, for the rates model– $P(L|T, R, S, \Theta)$, which the authors compute numerically after learning species and gene specific rates.
3. the second prior term, for the duplication/loss model– $P(T, R|S, \Theta)$, which the authors describe using a birth and death process.

Having a rates model is very useful, since mutation rates are quite variable across genes. In the lecture, we saw how rates were well described by a decomposition into gene and species specific rates. In lecture we saw that an inverse gamma distribution appears to parametrize the gene specific substitution rates, and we were told that a gamma distribution apparently captures species specific substitution rates. Accounting for gene and species specific rates allows SPIDIR to build gene trees more accurately than previous methods. A training set for learning rate parameters can be chosen from gene trees which are congruent to the species tree. An important algorithmic concern for gene tree reconstructions is devising a fast tree search method. In lecture, we saw how the tree search could be sped up by only computing the full $\arg\max_{L, T} RP(L, T, R|D, S, \Theta)$ for trees with high prior probabilities. This is accomplished through a computational pipeline where in each iteration 100s of trees are proposed by some heuristic. The topology prior $P(T, R|D, S, \Theta)$ can be computed quickly. This is used as a filter where only the topologies with high prior probabilities are selected as candidates for the full likelihood computation. The performance of SPIDIR was tested on a real dataset of 21 fungi. SPIDER recovered over 96% of the synteny orthologs while other algorithms found less than 65%. As a result, SPIDER invoked much fewer number of duplications and losses.

7 Conclusion

Incorporating species tree information into the gene tree building process via introducing separate gene and species substitution rates allows for accurate parsimonious gene tree reconstructions. Previous gene tree reconstructions probably vastly overestimated the number of duplication and loss events. Reconstructing gene trees for large families remains a challenging problem.

8 Current Research Directions

9 Further Reading

10 Tools and Techniques

11 What Have We Learned?

References

- [1] O. Akerborg, B. Sennblad, L. Arvestad, and J. Lagergren. Bayesian gene tree reconstruction and reconciliation analysis. *Proc Natl Acad Sci*, 106(14):5714–5719, Apr 2009.
- [2] Zmasek C.M. and Eddy S.R. Analyzing proteomes by automated phylogenomics using resampled inference of orthologs. *BMC Bioinformatics*, 3(14), 2002.
- [3] Li H, Coghlan A, Ruan J, Coin LJ, Heriche JK, Osmotherly L, Li R, Liu T, Zhang Z, Bolund L, Wong GK, Zheng W, DEhal P, Wang J, and Durbin R. Treefam: a curated database of phylogenetic trees of animal gene families. *Nucleic Acids Res*, 34, 2006.
- [4] Arvestad L., Berglund A., Lagergren J., and Sennblad B. Bayesian gene/species tree reconciliation and orthology analysis using mcmc. *Bioinformatics*, 19 Suppl 1, 2003.

- [5] M. D. Rasmussen and M. Kellis. A bayesian approach for fast and accurate gene tree reconstruction. *Mol Biol Evol*, 28(1):273290, Jan 2011.
- [6] Matthew D. Rasmussen and Manolis Kellis. Accurate gene-tree reconstruction by learning gene and species-specific substitution rates across multiple complete genomes. *Genome Res*, 17(12):1932–1942, Dec 2007.
- [7] C.E.V. Storm and E.L.L. Sonnhammer. Automated ortholog inference from phylogenetic trees and calculation of orthology reliability. *Bioinformatics*, 18(1):92–99, Jan 2002.
- [8] Hollich V., Milchert L., Arvestad L., and Sonnhammer E. Assessment of protein distance measures and tree-building methods for phylogenetic tree reconstruction. *Mol Biol Evol*, 22:2257–2264, 2005.
- [9] Wapinski, I. A. Pfeffer, N. Friedman, and A. Regev. Automatic genome-wide reconstruction of phylogenetic gene trees. *Bioinformatics*, 23(13):i549–i558, 2007.