## Parsimonious Reconciliation of Non-binary Trees

Louxin Zhang National University of Singapore matzlx@nus.edu.sg

### Gene Tree vs the (Containing) Species Tree.

A species tree S represents the evolutionary history of species

- S can be binary or non-binary.
  - -- Hard polytomy
  - -- Soft polytomy.
  - Each leaf represents a modern species.



A gene tree G is over the members of a gene family

G can be non-binary.

--- Non-binary nodes are soft polytomies.

Each leaf represents a family gene and is labeled by the species where it resides.

S1g S2g S1g S3g S4g Gene tree G

A **reconciliation** between gene tree G and species tree S is a map from V(G) to V(S) with the following properties:

Leaf-preserving: f(x) is a leaf with the same label if x ∈ Leaf(G);
Order-preserving: u ≤ v ⇒ f(u) ≤ f(v), ∀u, v ∈ V(G).



# Gene tree and species tree reconciliation is an important method for

- Inferring duplications, losses, and horizontal transfers
- Inferring orthology and paralogy gene relationship





**Lowest Common Ancestor Reconciliation**  $\lambda$  (Goodman et al, 1979, Page, 1994)

- -- It maps leaves to respective leaves with the same label;
- -- it maps internal node g is the lowest common ancestor of the images of its children.



Let  $u \in V(G)$  with children  $u_1$  and  $u_2$ .

► A duplication is inferred with *u* iff

$$\lambda(u) = \lambda(u_1), \text{ or } \lambda(u) = \lambda(u_2).$$

The inferred duplication occurred in the branch entering  $\lambda(u)$ 

• Gene losses are computed using the # of branches on the path from  $\lambda(u)$  to  $\lambda(u_i)$ , i = 1, 2.

(The duplication cost of  $\lambda$ ) = (the # of inferred duplications). The gene loss cost is defined similarly.  $\Re(G, S)$ : All the reconciliations of G and S.  $\lambda$ : The lca reconciliation of G and S.

**Theorem** Let G and S be binary. Then,

- i).  $\lambda$  has the smallest duplication cost in  $\mathcal{R}(G, S)$  (Gorecki & Tiuryn, 2006);
- ii).  $\lambda$  is the unique one with the smallest loss cost in  $\mathcal{R}(G, S)$  (Chauve et al, 2009);
- iii).  $\lambda$  is the unique one with the smallest deep coalescence cost in  $\mathcal{R}(G, S)$  (Wu & Zhang, 2011);
- iv).  $\lambda$  is computable in O(|G|+|S|) (Zhang, 1997).

### $\lambda$ is the parsimonious solution for binary trees.

### Part II: Reconciliation with Non-binary Trees

#### **General Reconciliation Problem:**

**Instance**: A gene tree G and a species tree S;

**Solution**: Binary refinement  $\hat{G}$  of G and  $\hat{S}$  of S such that the lca reconciliation of  $\hat{G}$  and  $\hat{S}$  minimizes a reconciliation cost.



#### **Our Heuristic Reconciliation Procedure**



**Step 1** Refine S based on structural inform. of gene tree G **Step 2** Refine G based on the refinement Ŝ of species tree S

**Step 3** Reconcile  $\hat{G}$  and  $\hat{S}$  to infer the evolution of the gene family

### Refine S based on the Structural Information of G

**Instance**: A (binary or non-binary) *G* and a non-binary *S*. **Solution**: A binary refinement  $\hat{S}$  of *S* that minimizes a reconciliation cost of reconciling *G* and  $\hat{S}$ .

The above refinement problem is NP-hard even for binary gene trees, which is proved for the duplication cost via a reduction from

### **Species tree problem**

**Instance:** A set of gene trees  $G_i$  (0<*i*<*n*). **Solution:** A species tree S' that minimizes  $\sum_{0 < i < n} c(S', G_i)$ 

> Ma, Li, & Zhang, SIAM J. Computing, 2000 Zhang, IEEE TCBB, 2011 Bansal & Shamir, IEEE TCBB, 2011 Than & Nakhleh, PLoS Comput. Biol. 2009





 $A_i \cap P = \varphi \text{ or } A_i \cap Q = \varphi \text{ for every } i = 1,2$ 

partial partitions (in red)



### Performance of First-Partition Algorithm

We repeat the following test 1000 times for each of 8 combinations of k (# of splits over a ground set) and t (the size of the ground set)

# of	# of	# of errors	# of errors
$ ext{elements}$ ( $t$ )	$\operatorname{splits}(k)$	by FP <sup>+</sup>	by HC <sup>‡</sup>
5	5	7	15
	10	0	18
10	5	0	4
	10	1	2
	20	0	0
15	7	0	3
	15	0	1
	30	0	1

- + An algorithm made an error if it output a non-optimal partition on the input.
- HC is an algorithm designed through a reduction to the Min Hypergraph Cut problem (Ouangraoua, Swenson, Chauve, 2009).

### Resolving non-binary nodes in G based on Ŝ



The following duplication inference rule does not work for non-binary nodes:

A duplication is associated with *u* having children  $u_1$  and  $u_2$  iff  $\lambda(u) = \lambda(u_1)$ , or  $\lambda(u) = \lambda(u_2)$ .

We present an extension of above rule to non-binary nodes. The whole process takes O(|G|+|Ŝ|) time.

### $\lambda$ : The lca reconciliation of G and $\hat{S}$



The node v and its children are mapped to a subtree (in blue) under λ, which is expanded into a binary subtree (by adding dark blue edges).



The image subtree in  $\hat{S}$ 



Step1: Compute m(u), the maximum number of child images on a path from an internal node u to a leaf descendant.

#### Algorithm

 $\omega(u)$  is the # of children mapped to *u*.  $m(u) = \max\{ m(u_1), m(u_2)\} + \omega(u)$ 



#### Algorithm

$$\begin{split} &\alpha(r) = 1\\ &\alpha(u) = \beta(p(u)) - \omega(p(u))\\ &\beta(u) = \begin{cases} m(u) \text{ if } \alpha(u) \ge m(u) \text{ or } u \text{ is a leaf;}\\ &\max\{\alpha(u), \min\{m(u_1), m(u_2)\} + \omega(u), 1 + \omega(u)\}. \end{cases} \end{split}$$

 $\alpha(u)$ : the # of genes entering a branch.

 $\beta(u)$ : the # of genes leaving a branch.



**Step1:** Compute m(u),

:

**Step2:** Compute  $\alpha(u) / \beta(u)$  using m(u).

**Step 3:** Infer duplications and losses.

If  $\alpha(u) < \beta(u)$ , duplications ( $\square$ ) are postulated.

If  $\alpha(u) > \beta(u)$ , losses (  $\blacksquare$  ) are postulated.

# **Theorem** The above algorithm resolves a non-binary node v with m(r)-1 duplications.

**Sketch of Proof.** Assume the following path from the root to a leaf contains the largest number, m(r), of child images,

$$P: r = u_0, u_1, \cdots u_k$$

4=m(r)

(1) There are no gene losses on P, i.e.,  $\alpha(u_j) \le \beta(u_j)$ . (2) All the duplications are postulated on P.

# of duplications

$$= \sum_{i=0}^{k} [\beta(u_{i}) - \alpha(u_{i})]$$
  
=  $\beta(u_{0}) - 1 + \sum_{i=1}^{k} [\beta(u_{i}) - \beta(u_{i-1}) + \omega(u_{i-1})]$   
=  $\sum_{i=0}^{k} \omega(u_{i}) - 1$   
=  $m(\lambda(r)) - 1$ .

Thm (i) The obtained reconciliation of a non-binary node v has the optimal dup. cost.
(ii) It also has the smallest loss cost over all the reconciliations with the same duplication cost.

**Idea of Proof.** Let *v* have children  $v_1, v_2, \ldots v_k$ . We consider partially ordered set:

 $\mathcal{P} = (\{L(\lambda(v_i)) : 1 \le i \le k\}, \subseteq).$ 

(1) L: The size of the longest chain in *P*P: The min. # of antichains into which *P* may be partitioned.

#### **Dual of Dilworth Theorem (Mirsky, 1971): L=P.**

(2) At least p-1 duplications are needed to produce all the children of v (Berglund-Sonnhammer et al, '06, Chang & Eulenstein'06)

### Part III: Software and Experiments

Our algorithms have been implemented in Python.

- Our program reconciles one or more gene trees and a species tree in the duplication cost, or the duplication and then loss cost.
- It is executed from command line to allow for automated analysis of large data sets.
- ► No limitation on the number of species.
- Automatically rerooting gene trees.

We validated our program on simulated and real data; we also compared it with NOTUNG (Durand et al'08), which requests either gene tree or species is binary. Repeated the following experiment 1000 times for n=20, 40, 60, 80, 100:

- -- Generate a binary species tree S over n species, and a non-binary species tree S' from S by randomly contracting edges.
- -- Generate 16 binary gene trees over totally about 1.2n genes in n species using proper duplication and gene loss rates in S; and divided them into 4 groups containing 1, 3, 5, 7 gene trees respectively.
- -- Find the reconciliation of each group of gene trees and S'.
- -- Check whether S and the refinement of S' are identical or not.

### **Performance Analysis**





# of gene trees in a data set

We tested our program on

- 1. A gene tree of Tor in 13 fungal species (Shertz et al. 2011) and a non-binary species tree from NCBI taxonomy database.
  - --- Inferred duplication events are identical to those reported in the paper;
  - --- Output refinement of the species tree is consistent with a large binary fungal species tree appearing in literature (www.broadinstitute.org)
- 2. A non-binary STAT gene tree and a binary species tree.
  - --- Co-evolution of STAT and other proteins in its signaling pathway.

#### The JAK/STAT Signaling Pathway

Jason S. Rawlings, Kristin M. Rosler and Douglas A. Harrison



STAT: Signal Transducer and Activator of Transcription



#### 





**Co-evolution of Egfr, Jak and STAT genes**.

Nakatani, Takeda, Kohara, & Morishita, 2007

### Conclusion

- We developed a software for reconciliation with non-binary trees
  - -- For binary gene tree and non-binary species tree, our program output a reconciliation with much less duplications than NOTUNG.
  - -- Durand et al (2005, 2008), Chang & Eulenstein (2006), Berglund-Sonnhammer et al.(2006)
- Parsimony approach vs Bayesian approach
  - -- Akerborg et al (2009); Arvestad et al, (2009)
- Study how to reconcile non-binary gene tree and HGT (horizontal gene transfer) networks in future.

### Acknowledgment

Reconciliation work:

-- Taoyang Wu -- Yu Zheng

STAT evolution work:

- -- Xin-Yuan Fu (Dept of Biochemistry)
- -- C. Pawan Kumar Patro
- -- Choong Yong Ung

### Thank You!

\_\_\_\_

