

The Impact of Model Misspecification on Tree and Network Inference from Quartets

IMSI

Hector Baños

Department of Mathematics

Tuesday, August 12, 2025



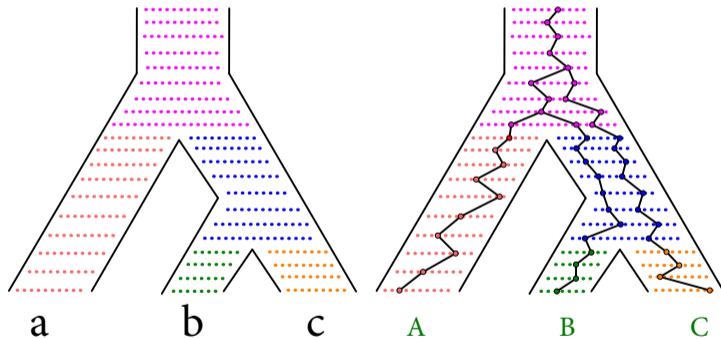


Vu Dinh



The Multispecies Coalescent Model

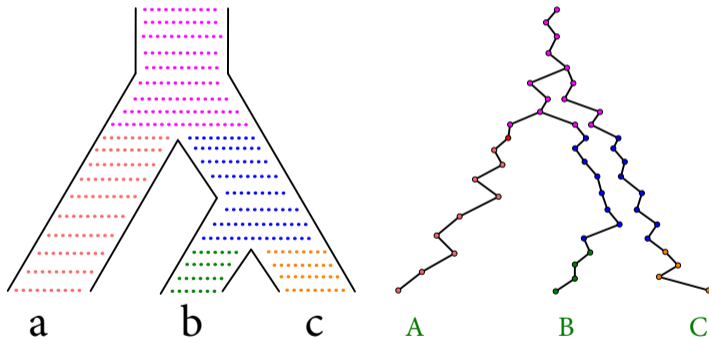
Rannala and Yang '03



The multi-species coalescent describes a stochastic model of gene tree generation accounting for gene tree incongruence.

The Multispecies Coalescent Model

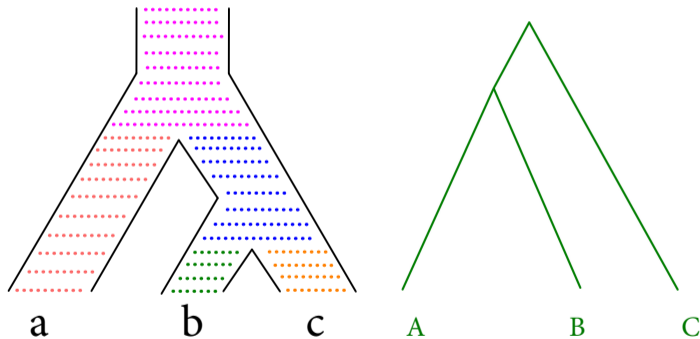
Rannala and Yang '03



The multi-species coalescent describes a stochastic model of gene tree generation accounting for gene tree incongruence.

The Multispecies Coalescent Model

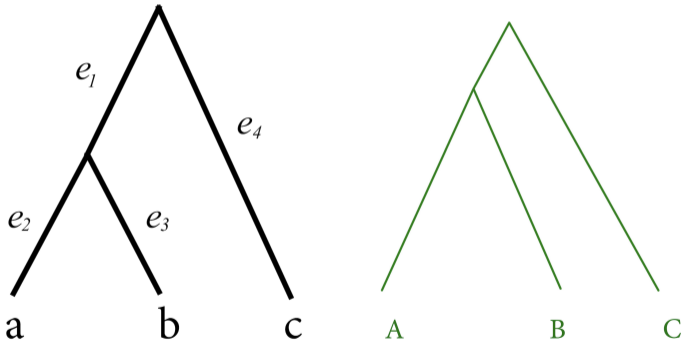
Rannala and Yang '03



The multi-species coalescent describes a stochastic model of gene tree generation accounting for gene tree incongruence.

The Multispecies Coalescent Model

Rannala and Yang '03

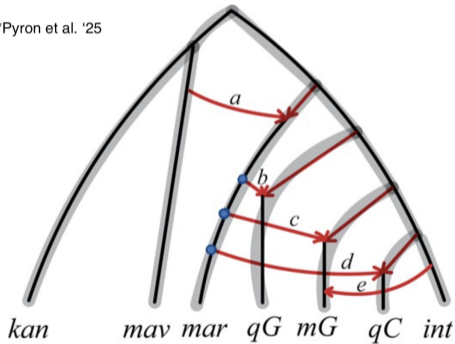


The multi-species coalescent describes a stochastic model of gene tree generation accounting for gene tree incongruence.

Hybridization

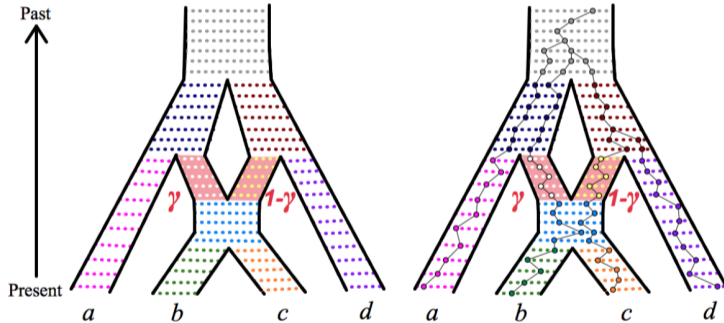
Hybridization occurs when two species merge genetically to create a new one.

*Pyron et al. '25



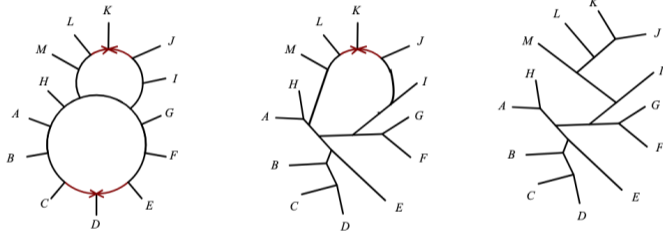
The Network Multispecies Coalescent Model

Meng & Kubatko '09 - Degnan, Yu, & Nakhleh '12



The network multi-species coalescent describes a stochastic model of gene tree generation in the presence of hybridization.

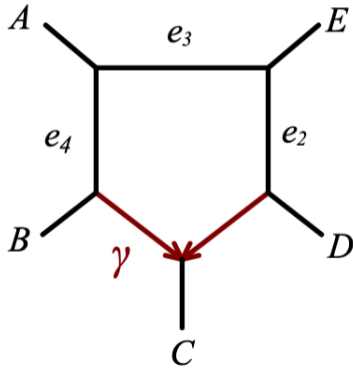
Question: How do inference methods behave when applied to data originating from a more complex topological structure than the one assumed by the method?



- Solís-Lemus, Yang & Ané. *Inconsistency of species tree methods under gene flow* '16.
- Long & Kubatko. *The effect of gene flow on coalescent-based species-tree inference* '18.
- Pang & Zhang. *Impact of ghost introgression on coalescent-based species tree inference and estimation of divergence time* '22.

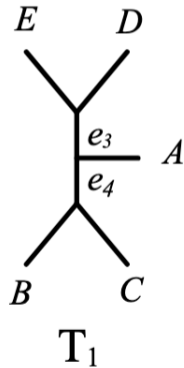
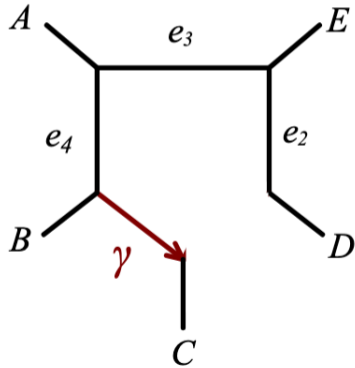
Network's Displayed Trees

A displayed tree is obtained by removing exactly one hybrid edge from each hybridization event.



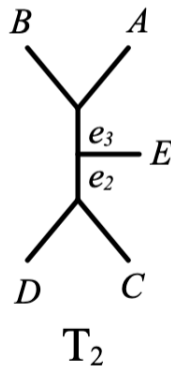
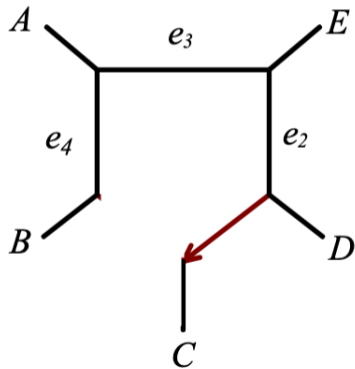
Network's Displayed Trees

A displayed tree is obtained by removing exactly one hybrid edge from each hybridization event.



Network's Displayed Trees

A displayed tree is obtained by removing exactly one hybrid edge from each hybridization event.

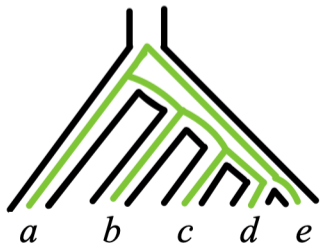


Anomalous tree

In the tree setting (i.e no hybridization),

Definition (Anomalous tree)

A species tree is said to be **anomalous** under the MSC if the most probable unrooted gene tree does not match the unrooted species tree topology.



Matches the unrooted species tree

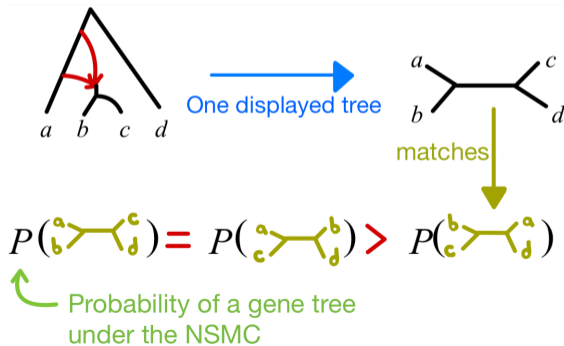
$$P\left(\begin{array}{c} b \\ \swarrow \quad \searrow \\ c \quad a \\ \swarrow \quad \searrow \\ d \quad e \end{array}\right) > P\left(\begin{array}{c} c \\ \swarrow \quad \searrow \\ a \quad e \\ \swarrow \quad \searrow \\ b \quad d \end{array}\right)$$

Probability of the unrooted gene tree under the NSMC

Definition (Anomalous network - Ané, et al.)

Let N^+ be a rooted metric network on X . N^+ is **anomalous**[†] under the NMSC if there is a gene tree whose topology is not displayed by N , that is more probable than a gene tree matching the topology of a displayed tree of N .

For example, Solís-Lemus and Ané '16



[†]Differs from Pang & Zhang '22

Anomalous networks are a problem for inference methods and identifiability results.

Effects of anomalies on tree inference methods:

- Solís-Lemus, Yang, Ané. *Inconsistency of species-tree methods under gene flow* '16.

Identifiability problems for some networks:

- B. *Identifying species network features from gene tree quartets under the coalescent model* '19.

Characterizations of 4-taxon anomalous networks:

- Ané, Fogg, Allman, B., Rhodes. *Anomalous networks under the multispecies coalescent: theory and prevalence* '24.

Identifiability results which specifically require non anomalous scenarios:

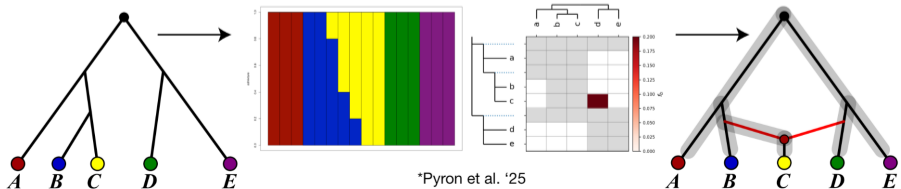
- Rhodes, B., Xu, Ané. *Identifying circular orders for blobs in phylogenetic networks* '25.

For the rest of the talk we assume NO quartet anomalous scenarios

Many network inference algorithms either infer “simple” networks or are not easily scalable.

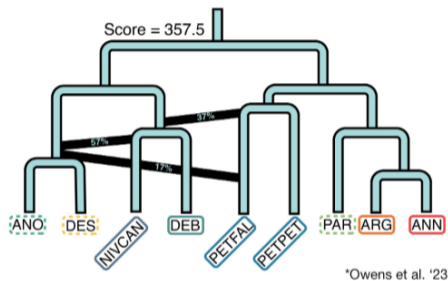
A reasonable approach avoiding either oversimplification of the network or scalability issues is to:

- 1) Infer a “displayed” tree, representing underlying tree-like relationships among species
- 2) then inferring hybridization events on top of it using different techniques (for e.g. Dsuite)



The use of ASTRAL for Network Inference

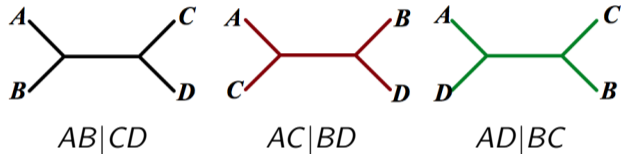
Many works have used ASTRAL for Inferring a “displayed” tree,



Owens et al., 2023; Zhou et al., 2022; Singh et al., 2022; Jensen et al., 2023; Sanderson et al., 2023; Ciezarek et al., 2024; Scherz et al., 2022; Yang et al., 2023; Lopes et al., 2023; Feng et al., 2022; Bernhardt et al., 2020; DeRaad et al., 2022; Herrig et al., 2024; Zhang et al., 2023; Zhou et al., 2023.

Quartets

A *quartet* is a binary tree on 4 taxa. There are 3 different quartet trees on 4 taxa:



Given a sample of gene trees, for any subset of four taxa, each gene tree displays exactly one quartet on those taxa.



ASTRAL is a powerful and widely used tool for species tree inference, known for its computational speed and consistency under the MSC (Mirarab, Warnow, et al. '14.).

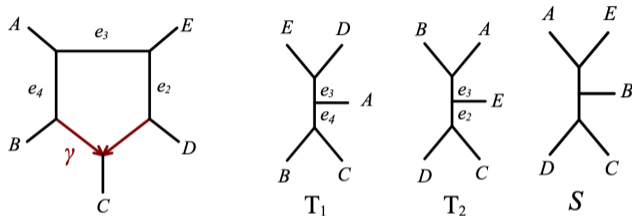
Given a collection of gene trees $\mathcal{T}_m = \{T_1, T_2, \dots, T_m\}$, it searches for a species tree \mathbb{T} such that

$$A(\mathbb{T}) = \frac{1}{m} \sum_{q \in Q(\mathbb{T})} w_m(q, \mathcal{T}_m) \text{ is **maximized**,}$$

where $Q(\mathbb{T})$ is the set of quartet trees induced by \mathbb{T} and $w_m(q, \mathcal{T}_m)$ is the number of the trees in \mathcal{T}_m that induce quartet topology q .

Theorem (Allman, Degnan, Rhodes '11)

Under the MSC, the most probable quartet gene tree has the same topology as the quartet species tree.

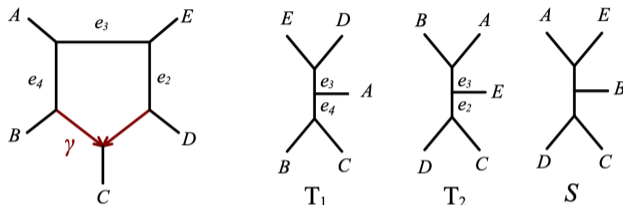


Theorem (Dinh, B.)

Let N be the semi-directed network on 5 taxa above and let $y_i = \exp(-e_i)$. If

$$(1 - y_2) > \frac{\gamma}{1 - \gamma} (1 - y_3 y_4) \quad \text{and} \quad (1 - y_3) < \min \left(\frac{\gamma}{2 - \gamma} (1 - y_4), \frac{1 - \gamma}{1 + \gamma} (1 - y_2) \right),$$

under the NMSC, the tree S above, which is not displayed by N , has a higher expected ASTRAL score than the displayed trees of N .



Quartet in T_1	Quartet in T_2
AB DE	AB DE
AD BC	AB CD
BC DE	BE CD
AE BC	AB CE
AC DE	AE CD

$$A(T_1) = P(\mathbf{AB|DE}) + P(AD|BC) + P(BC|DE) + P(\mathbf{AE|BC}) + P(AC|DE)$$

$$A(T_2) = P(\mathbf{AB|DE}) + P(\mathbf{AB|CD}) + P(\mathbf{BE|CD}) + P(AB|CE) + P(\mathbf{AE|CD})$$

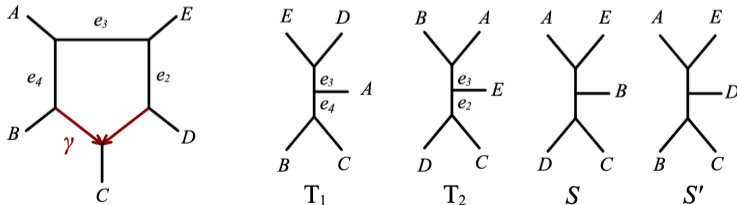
$$A(S) = P(AE|BD) + P(\mathbf{AB|CD}) + P(\mathbf{BE|CD}) + P(\mathbf{AE|BC}) + P(\mathbf{AE|CD})$$

Corollary (Dinh, B.)

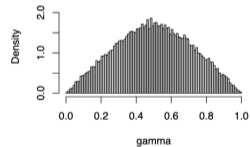
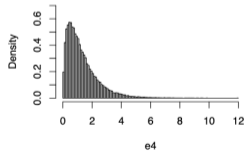
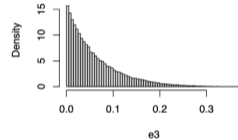
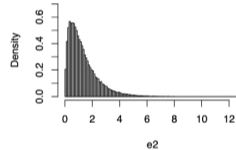
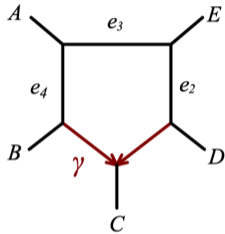
Let N be the semi-directed network on 5 taxa below and let $y_i = \exp(-e_i)$. If

$$(1 - y_4) > \frac{\gamma}{1 - \gamma} (1 - y_3 y_2) \quad \text{and} \quad (1 - y_3) < \min \left(\frac{\gamma}{2 - \gamma} (1 - y_4), \frac{1 - \gamma}{1 + \gamma} (1 - y_2) \right),$$

under the NMSC, the tree S' below, which is not displayed by N , has a higher expected ASTRAL score than the displayed trees of N .



We sampled 10^6 sets of parameters from the network below.



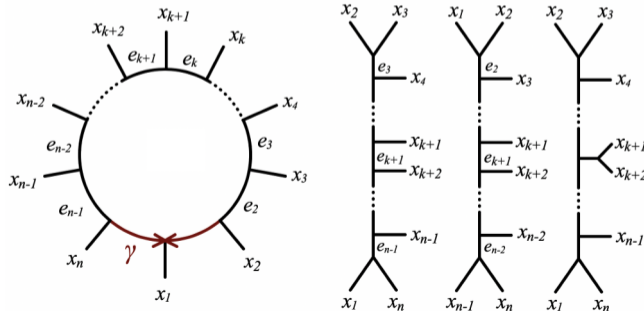
Range for γ	(0,1)	(0.2,0.8)	(0.4,0.6)	$(0,0.1) \cup (0.9,1)$
Proportion of parameters in Θ	0.06	0.08	0.10	0.01

This is Behavior can be Generalized

For sets of parameters satisfying the inequalities, we simulated gene trees using PhyloCoalSimulations (Ané, Fogg, Allman '24).

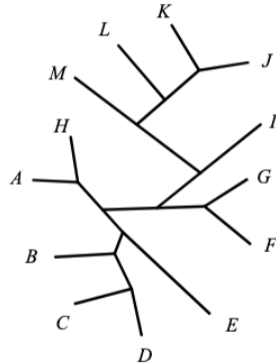
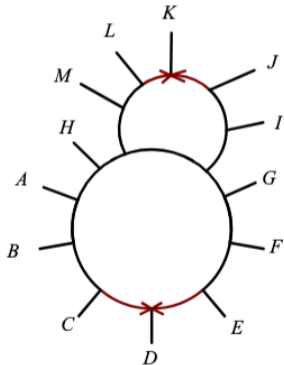
In all these, ASTRAL did not recover a displayed tree.

We generalized the results for bigger networks:



This is Behavior can be Generalized

We showed this behavior occurs in more complex networks



Not displayed!

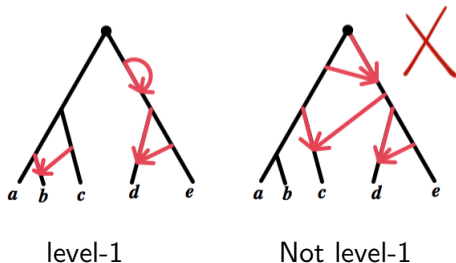
Why things behave like this?

This is a problem of misspecification on quartet-based methodologies not ASTRAL itself. For example:

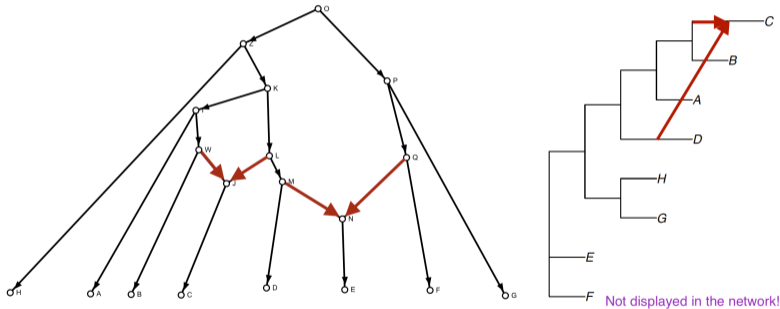
SNaQ is a quartet-based method for network inference under the NMSC. It uses a pseudo-likelihood framework. One key assumption is that the network is **level-1**

Definition

A network \mathcal{N} is **level-1** if no pair of cycles in \mathcal{N} share an edge.



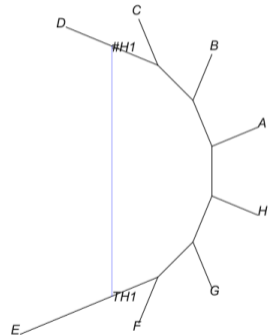
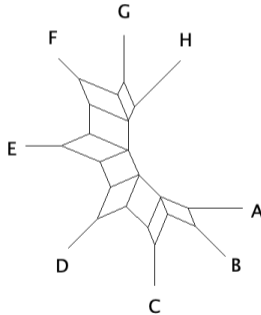
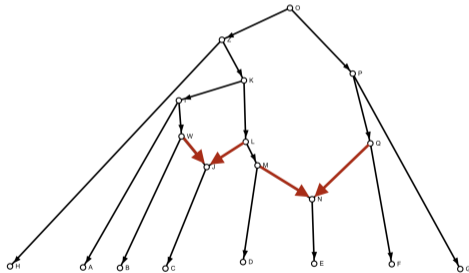
We have preliminary results showing that same issues can occur for SNaQ when inferring networks.



Along the Solís-Lemus Lab, we are exploring the behavior of SNaQ via a simulation study (In progress).

NANUQ and NANUQ+

NANUQ and NANUQ+ are quartet-based level-1 network inference methods (Allman, B., Rhodes, Wicke).



We suspect quartet-based methods such as SVDquartets and PhyNEST have similar behavior (joint work with Dinh, Allman, and Rhodes).

Can we avoid this issues?

We believe that this behavior can be extended.

That is, for $m > k$, there are parameters such that inferring a level- k network from data that came from a level- m network is problematic.

How could we overcome this?

- Displayed tree inference method

Along Pyron and some colleagues we have a 'proof of concept' method

Systematic Biology, Volume 74, Issue 1, January 2025, Pages 124–140

(Not easily scalable, theoretical limitations)

Quartet Concordance Factors

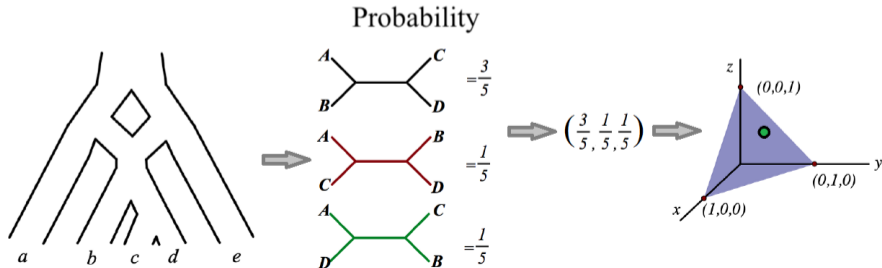
Following the approach of C. Solís-Lemus and C. Ané 2016

Definition

Let \mathcal{N} be a species network. Then

$$CF_{abcd} = (CF_{AB|CD}, CF_{AC|BD}, CF_{AD|BC})$$

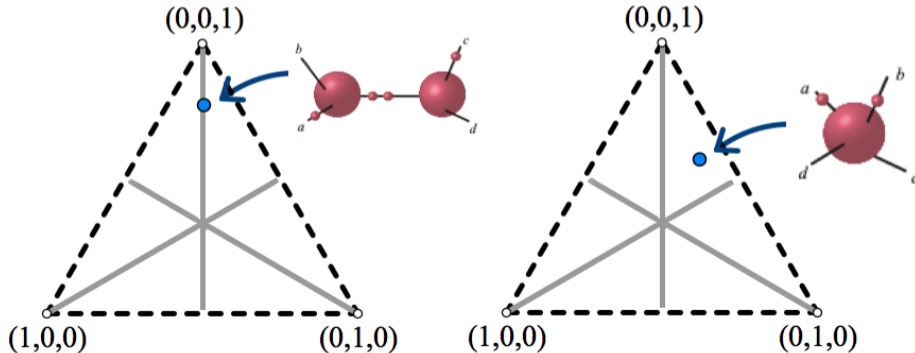
is the triplet of probabilities of gene trees quartets under the NMSC.



Cut and Non-cut Concordance Factors

The concordance factor CF_{abcd} is:

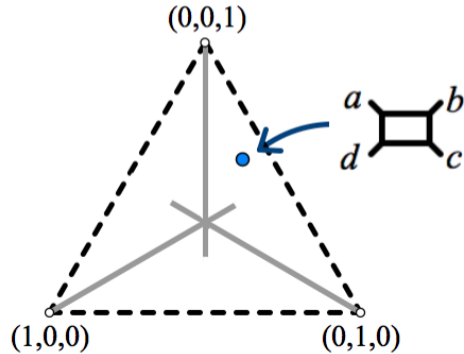
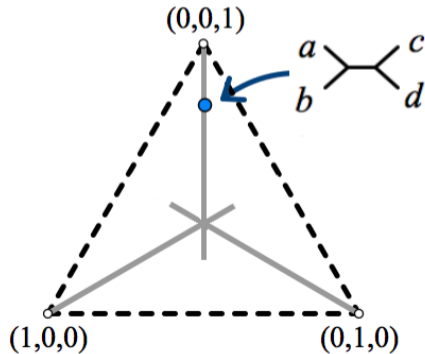
- a **cut** CF if two of its entries are equal, in addition the third is distinct, or
- a **non-cut** CF if it is not cut.



Cut and Non-cut Concordance Factors

The concordance factor CF_{abcd} is:

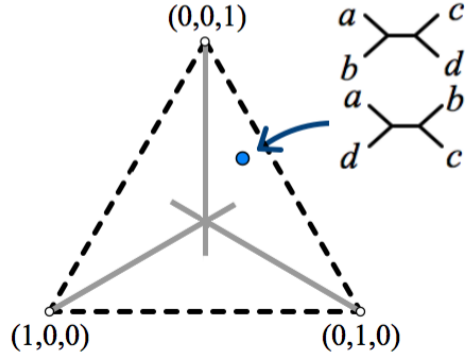
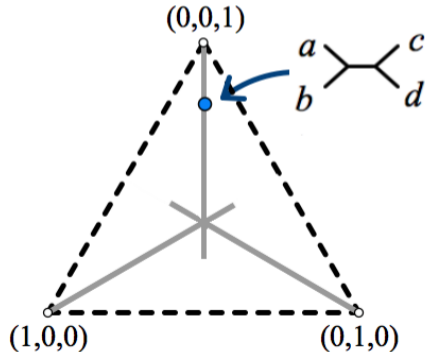
- a **cut** CF if two of its entries are equal, in addition the third is distinct, or
- a **non-cut** CF if it is not cut.



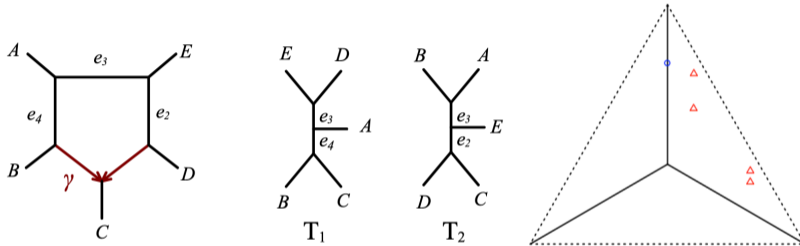
Cut and Non-cut Concordance Factors

The concordance factor CF_{abcd} is:

- a **cut** CF if two of its entries are equal, in addition the third is distinct, or
- a **non-cut** CF if it is not cut.



Cut and Non-cut Concordance Factors



$$A(T_1) = P(\mathbf{AB}|\mathbf{DE}) + P(AD|BC) + P(BC|DE) + P(\mathbf{AE}|\mathbf{BC}) + P(AC|DE)$$

$$A(T_2) = P(\mathbf{AB}|\mathbf{DE}) + P(\mathbf{AB}|\mathbf{CD}) + P(\mathbf{BE}|\mathbf{CD}) + P(AB|CE) + P(\mathbf{AE}|\mathbf{CD})$$

$$A(S) = P(AE|BD) + P(\mathbf{AB}|\mathbf{CD}) + P(\mathbf{BE}|\mathbf{CD}) + P(\mathbf{AE}|\mathbf{BC}) + P(\mathbf{AE}|\mathbf{CD})$$

Vu an I believe that by weighting the gene tree quartets that arise from a **cut** CFs will lead to better displayed tree estimation.

Things to figure out:

- Optimal weight
- Identify the family of networks for which the algorithm is consistent
- Implementation

This is an open invitation

Thank you!



Questions?

Hector Banos

hector.banos@csusb.edu

