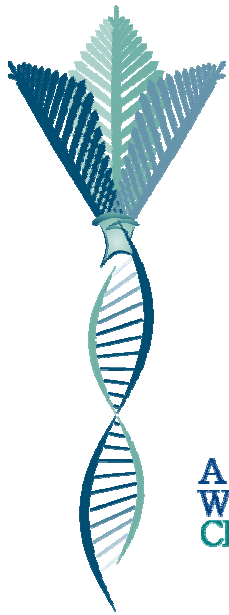# Phylogenetic closure operations and homoplasy-free evolution

IFCS, Chicago, June 2004

**Mike Steel**
Allan Wilson Centre for Molecular Biology and Evolution
Biomathematics Research Centre
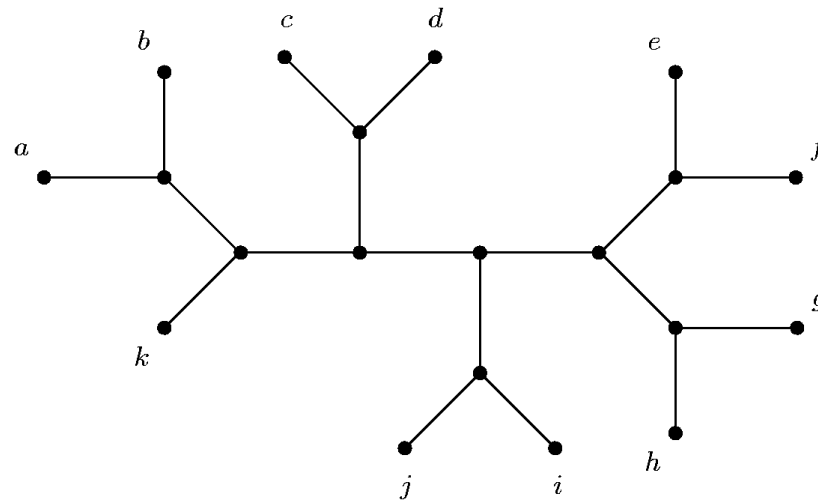University of Canterbury, Christchurch, New Zealand

ALLAN
WILSON
CENTRE

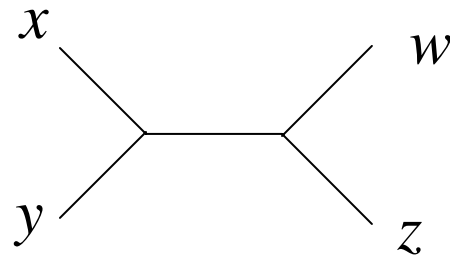Biomathematics
Research Centre

# Joint work with…

# Phylogenetic trees

■[Definition]     A **phylogenetic X-tree** is a tree $T=(V,E)$ with a set $X$ of labelled leaves, and all other vertices unlabelled and of degree $\geq 3$.

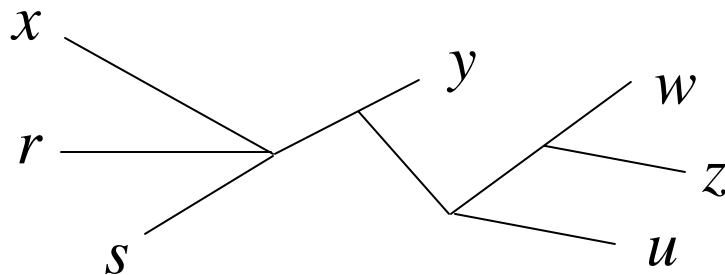■If all non-leaf vertices have degree 3 then $T$ is **binary**

# Quartet trees

- A **quartet tree** is a binary phylogenetic tree on 4 leaves (say, $x,y,w,z$) written $xy|wz$.
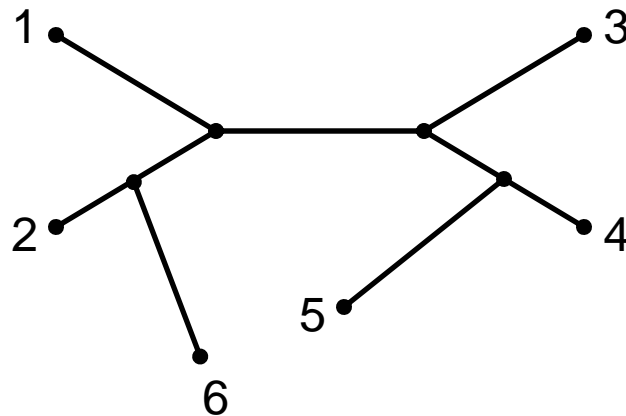
- A phylogenetic X-tree **displays** $xy|wz$ if there is an edge in $T$ whose deletion separates $\{x,y\}$ from $\{w,z\}$

# Compatibility

A set $Q$ of quartets is compatible if there is a phylogenetic $X$-tree T that **displays** each quartet of $Q$

- Example: $Q=\{12|34, 13|45, 14|26\}$

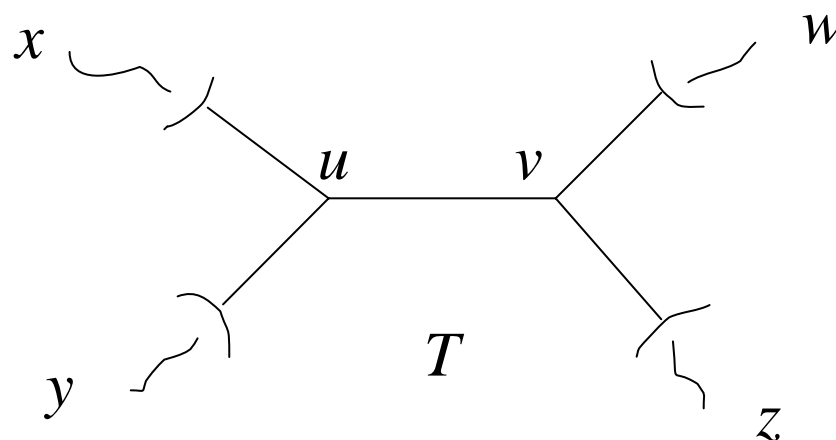# Defining sets

If *T* is the only phylogenetic *X*-tree that displays *Q* (and *X*= L(*Q*)) then we say *Q* **defines** *T* .

■Let *Q*(*T*) be the set of **all** quartets displayed by (any) *T*.

If *T* is binary, then *Q*(*T*) defines *T*.

# A necessary condition for $Q$ to define $T$

- **Definition:** For a binary phylogenetic tree $T$, a collection $Q$ of induced quartet trees *distinguishes* an interior edge $\{u,v\}$ of $T$ if there exists a quartet $xy|wz$ in $Q$ that looks like this:



**Observation:** If $Q$ defines $T$ then $T$ is binary and $Q$ distinguishes every interior edge of $T$ (so $|Q| \geq n$-3).

# Warning:

$Q =\{12|45, 56|23, 34|16\}$ distinguishes each interior edge of the tree:



and also                    !

# Sufficient condition for $Q$ to define T:

- Suppose $Q$ is compatible and distinguishes every interior edge of a binary phylogenetic $X$-tree $T$.

**Proposition:** If there is an element of $X$ that is a leaf of every tree in $Q$ then $Q$ defines $T$.

**Corollary:**

There are subsets of $Q(T)$ that define $T$ of size $|X|$-3.

# Character data

- **Type** **States** **Transitions**

- Morphology $\quad W(\text{ings}), \neg W, \ -W$ $\qquad \neg W \rightarrow W \rightarrow -W$
- Sequences $\qquad$ A,C,G,T $\qquad\qquad\qquad x \leftrightarrow y$

- Gene order $\quad g_1 g_2 g_3 g_4 g_5 g_6 g_7 \cdots$ $\quad g_1 g_2 \boxed{g_5 g_4 g_3} g_6 g_7 \cdots$
- Gene content $\quad G = \{g_1, \ldots, g_k\}$ $\qquad\quad +g_i / -g_i$
- SINEs $\qquad\qquad \ldots g \ldots$ $\qquad \ldots \rightarrow \ldots g \ldots \rightarrow \ldots ? \ldots$
- Oligonucleotides $\quad \ldots g_1 g_2 g_3 \ldots g_k \ldots = 1$ $\quad$ 0 to 1 (once), 1 to 0

# Definitions:

- [Character] A character is any function

$$f : X \to S$$

- [Convexity] Given a character $f : X \to S$
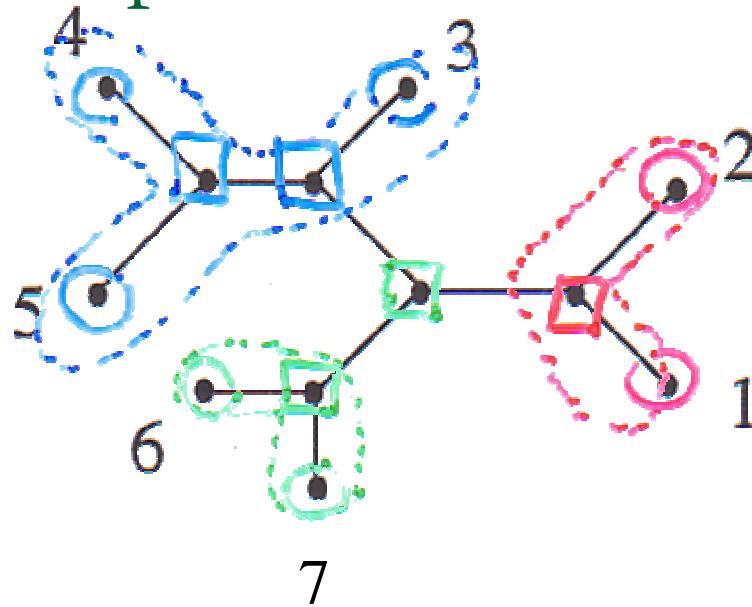and a phylogenetic $X$-tree $T=(V,E)$, we say $f$ is **convex on $T$**
if $f$ extends to $f':V \to S$
so that $f'\,|\,X = f$
and $\{v \in V : f'(v) = s\}$ is connected for all $s$ in $S$.
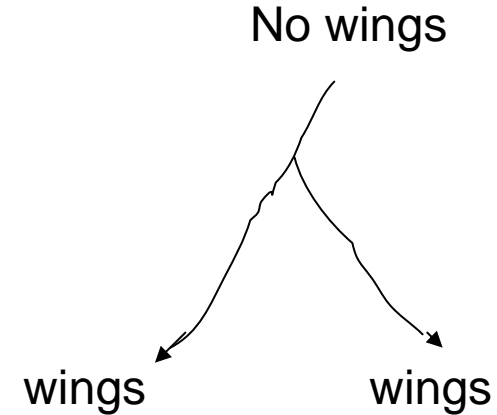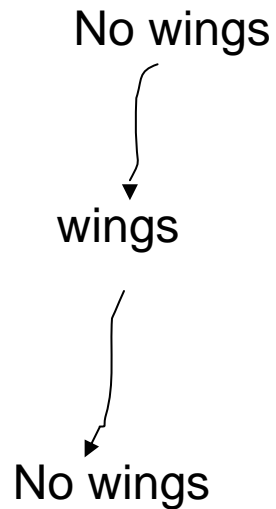
# Convexity: example



| $x$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| $f(x)$ | 🔴 | 🔴 | 🔵 | 🔵 | 🔵 | 🟢 | 🟢 |

# Biological significance of convexity

- **Lemma**: A character $\chi$ is convex on a phylogenetic tree $T$ if and only if $\chi$ could have evolved on $T$ (from any root vertex) without any **reversals** or **convergent evolution.**

No wings

wings

No wings
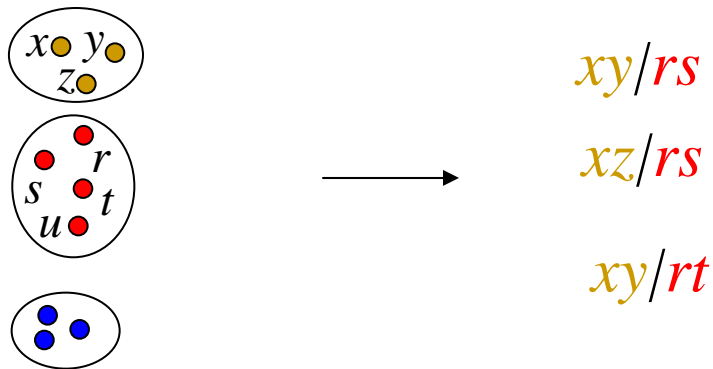
No wings

wings        wings

# Relevance to genomics

- Eg. gene order rearrangements (*n* species, *L* genes, random inversion model)

$$g_1 g_2 \boxed{g_3 g_4 g_5} g_6 g_7 \cdots \quad \Longrightarrow \quad g_1 g_2 g_5 g_4 g_3 g_6 g_7 \cdots$$

$$P[h = 0] \geq 1 - \frac{2(2n - 3)(n - 1)}{L(L - 1)}$$

# Equivalence of character and quartet compatibility

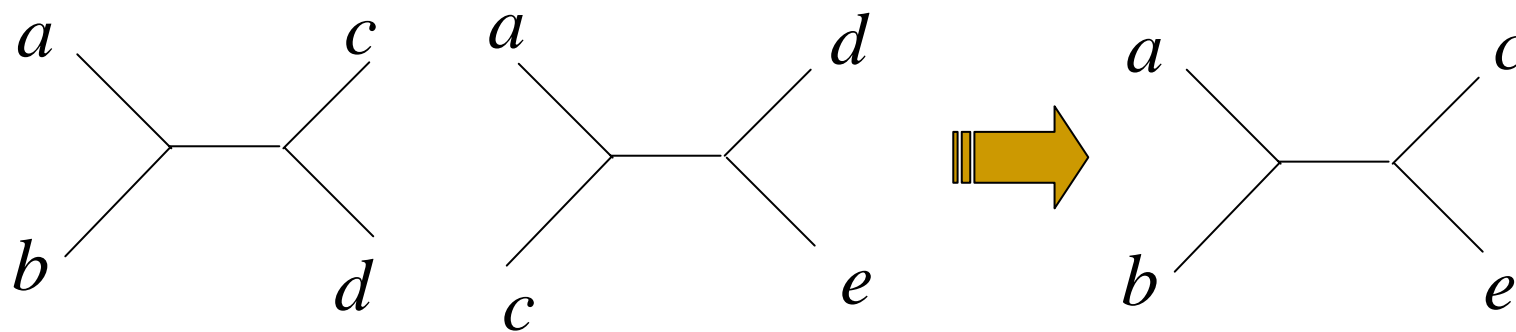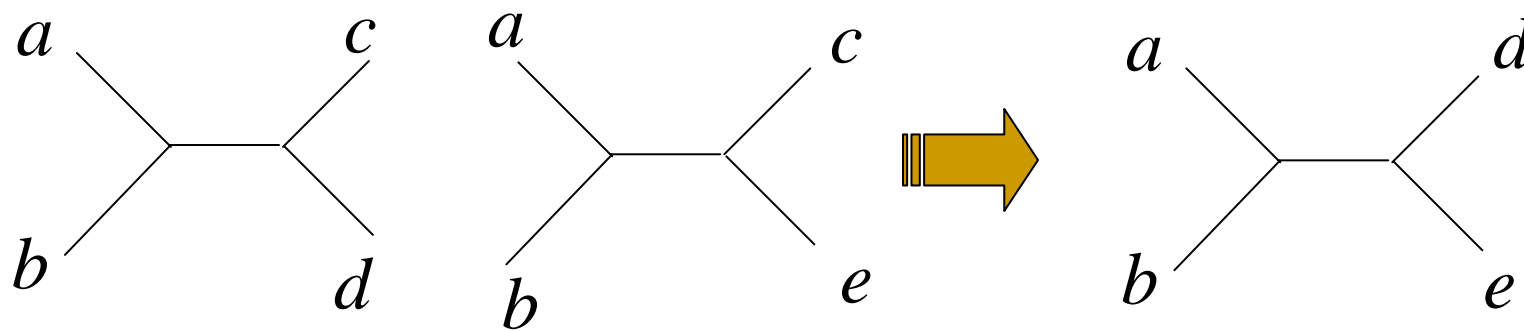$$C \rightarrow Q(C)$$



$xy/rs$

$xz/rs$

$xy/rt$

**Lemma:** Each character in $C$ is convex on $T$ if and only if $T$ displays all the quartets in $Q(C)$.

[$C$ is "compatible", $C$ "defines" $T$ iff $Q(C)$ does]

# New quartet trees from old ones

# Dyadic rules for quartet trees

## (Colonius and Schulze; Dekker)

($\mathbf{Q1}$): $\{ab|cd, ab|ce\} \vdash ab|de$

($\mathbf{Q2}$): $\{ab|cd, ac|de\} \vdash ab|ce, ab|de, bc|de.$

Any phylogenetic $X-$tree that displays the quartet trees on the left of ($\mathbf{Q1}$) or ($\mathbf{Q2}$) also displays the corresponding quartet tree(s) on the right.

# Dyadic quartet closure

$$\mathcal{Q} = \mathcal{Q}_1 \subseteq \mathcal{Q}_2 \subseteq \cdots \subseteq \mathcal{Q}_m = \mathrm{qcl}_\theta(\mathcal{Q})$$

where $\mathcal{Q}_{i+1}$ consists of $\mathcal{Q}_i$ together with all additional quartets that can be obtained from a pair of quartets in $\mathcal{Q}_i$ by applying the rule(s) allowed by $\theta$.
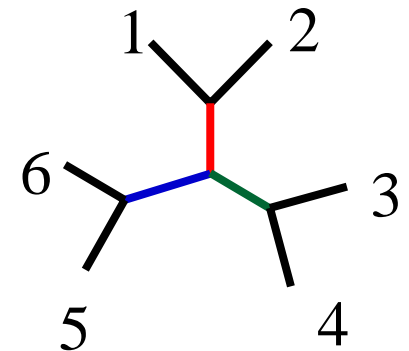
For $\theta \subseteq \{1, 2\}$, let the dyadic quartet closure under rule $\theta$, $\mathrm{qcl}_\theta(\mathcal{Q})$, denote the minimal set of quartet trees that contains $\mathcal{Q}$ and is closed under rule **(Qi)** for each $i \in \theta$.

We denote these closures with: $\mathrm{qcl}_1(\mathcal{Q}), \mathrm{qcl}_2(\mathcal{Q}), \mathrm{qcl}_{1,2}(\mathcal{Q})$.

# Example 1: $qcl_2$

**Definition:** If $Q$ distinguishes every interior edge of a binary phylogenetic tree $T$ and we can order $Q$ so that each quartet tree in the ordering introduces precisely one new leaf label, we say $Q$ has a **tight ordering** for $T$.

Example: $\{12|35, 13|56, 15|34\}$.



**Proposition:**

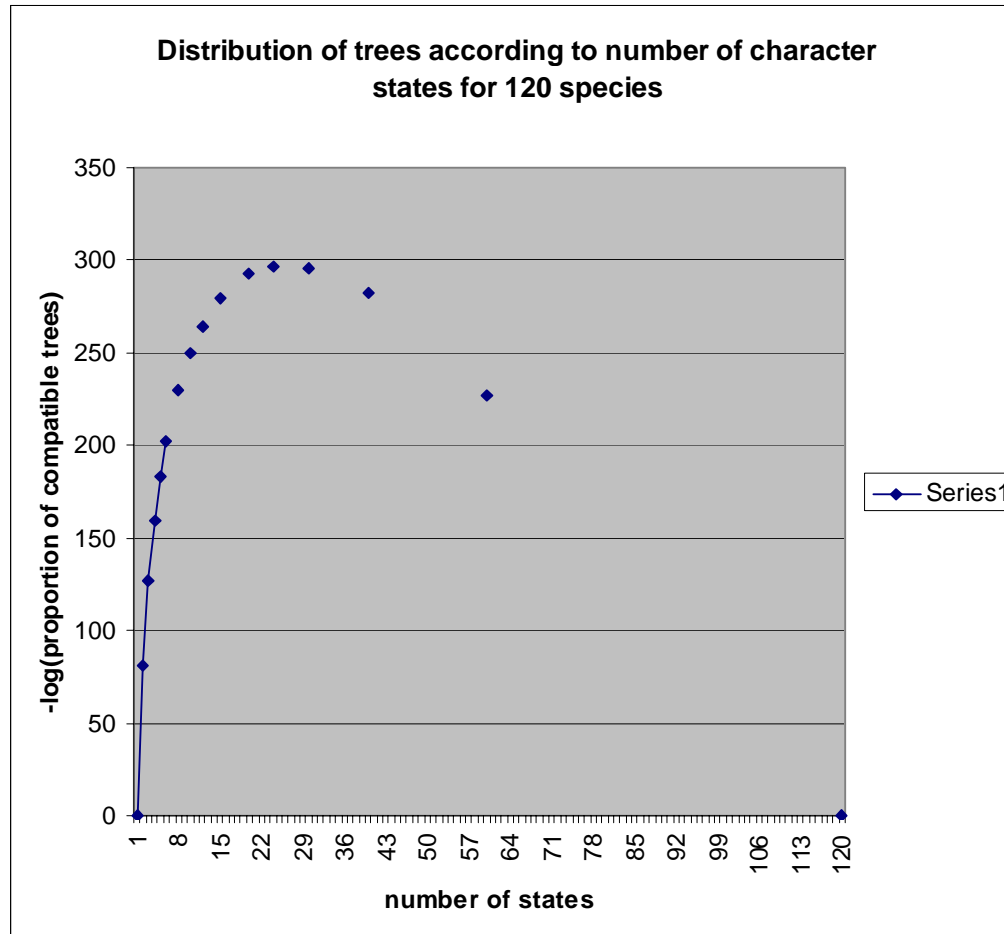If $Q$ has a tight ordering for $T$, then $qcl_2(Q) = Q(T)$

In particular $Q$ defines $T$.

## Application: How many characters are needed to define a binary phylogenetic $X$-tree?
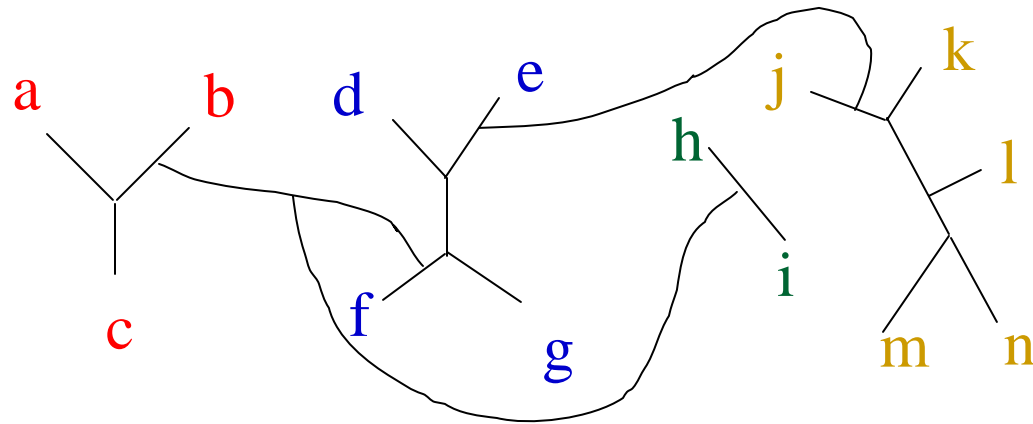
- For binary characters we need $n$-3 ($n=|X|$).
- For $r$-state characters ($r$ fixed) we need at least

    $(n$-3$)/(r$-1$)$

- What if $r$ is not fixed?

(it is not useful to make $r$ too large!)

# $\mathrm{I}(\chi) := - \log(\mathbf{Pr}[\chi \text{ is convex on random } T])$



Distribution of trees according to number of character states for 120 species

# Where do these numbers come from?

a   b   d   e   j   k
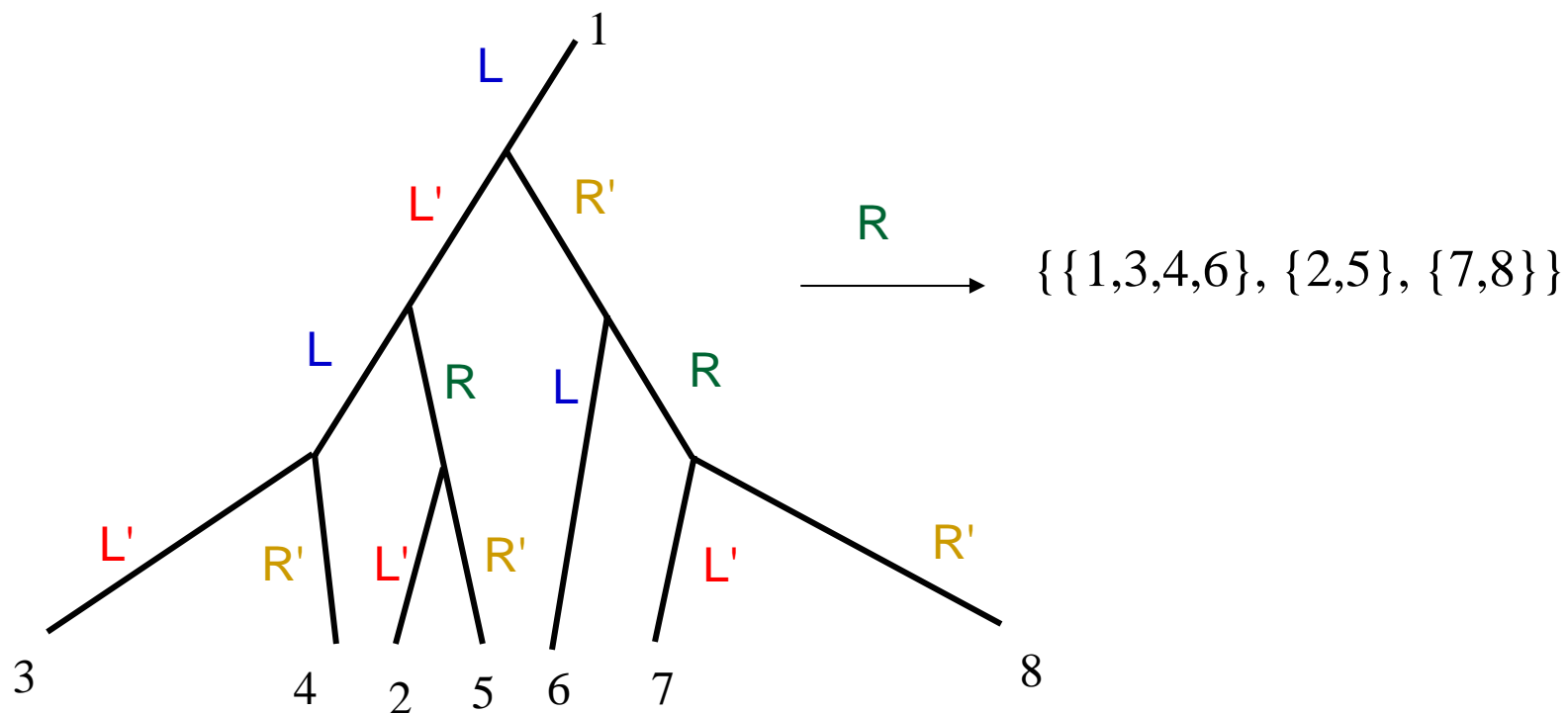
h   l

c   f   g   i   m   n

**Carter *et al.* (1990); Erdös & Székely (1993).**

\# binary phylogenetic trees with n leaves,   $b(n) = 1 \times 3 \times \cdots \times (2n-5)$

\# of these on which $\chi$ is convex =   $b(n) \prod_{i=1}^{r} b(a_i + 1) \Big/ b(n-r+2)$

# Edge-colouring a tree by $Z_2 \times Z_2$



$$\{\{1,3,4,6\}, \{2,5\}, \{7,8\}\}$$

**Theorem** (Huber, Moulton, S, 2003)

$Q(C)$ contains a subset with a tight ordering for $T$.

Thus for any tree there is a set of just *four* characters that defines $T$.

# Distances or characters?

- $d_C(i,j)$ = # characters in C on which $i$ and $j$ differ

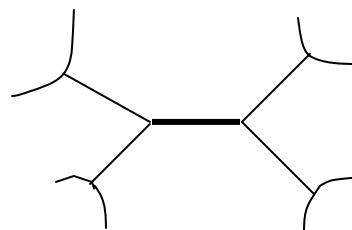If $C$ is compatible is $d_C$ tree-like?

$C$ binary – yes.
$C$ non-binary no.

**Theorem** [Huson and S, 2003]:

For **any** two trees $T_1$, $T_2$ there is a set of multi-state characters $C$ such that
- $C$ defines $T_1$ (i.e. $C$ homoplasy-free only on $T_1$) yet
- $d_C$ is tree-like (and ultrametric!), but only on $T_2$.

# Application 2: "Short" quartets



- $Q_{\mathrm{short}}(T)$

- **Theorem** (Erdös *et al.* 1997)
$Q_{\mathrm{short}}(T)$ contains a subset that has a tight ordering for $T$
(and so $\mathrm{qcl}_2(Q_{\mathrm{short}}(T)) = Q(T)$).

- The number of characters required to reconstruct (wp $>1-\varepsilon$) a binary phylogenetic tree with $n$ leaves from binary characters generated under a finite Markov process is (for almost all trees) at most
$$k \geq \frac{c_\varepsilon (\log(n))^{d(p)}}{a^2}$$
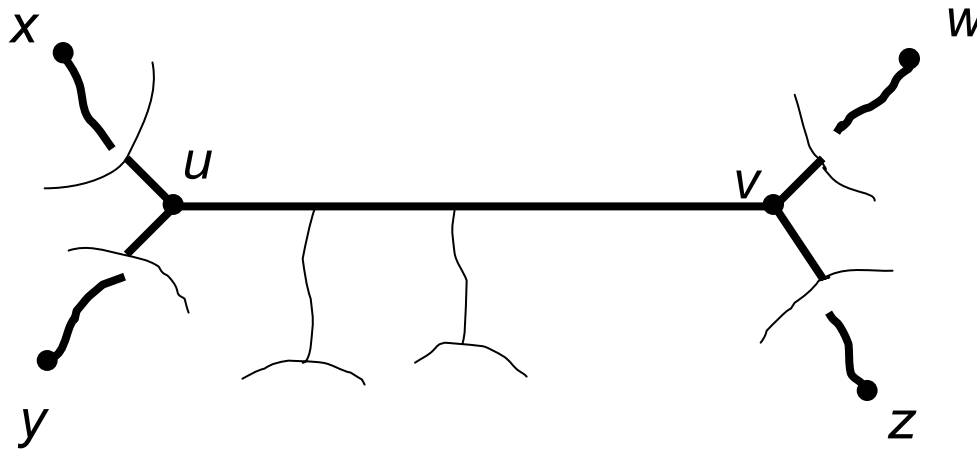
# A further application involving $qcl_2$:

We say $Q$ is **excess-free** if $|L(Q)|-3-|Q| = 0$.

- **Proposition:** Suppose a subset $Q$ of $Q(T)$ contains an excess-free subset $Q_0$ that defines $T$. Then $qcl_2(Q)=Q(T)$.

- Why? Let us say a set $Q$ of quartet trees is "good" if (i) $Q$ defines a phylogenetic tree, and (ii) $exc(Q)=0$.

  Theorem [Bocker, Dress 1999] Any good set of ($\geq 2$) quartets is the disjoint union of precisely two good sets.

# Example 2: $qcl_1$
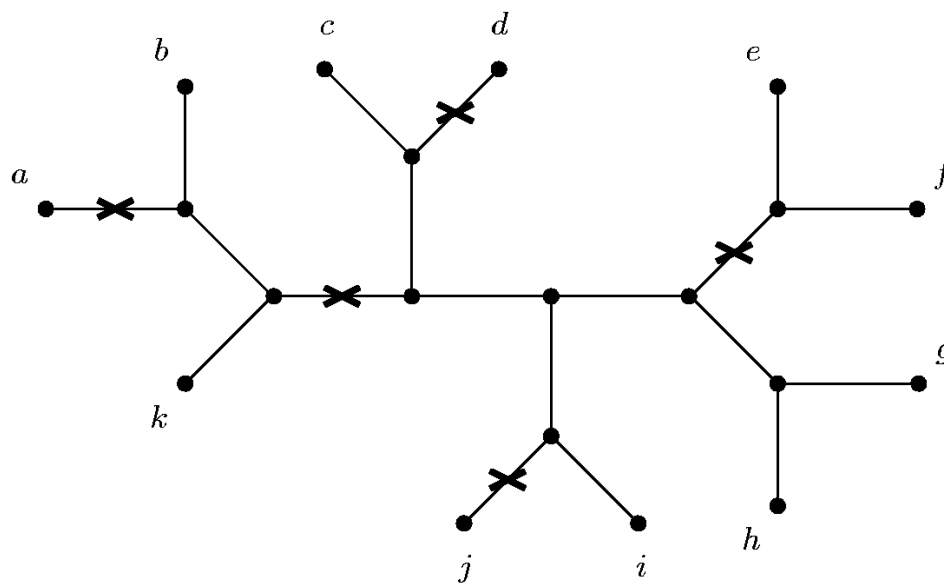
- **Definition:** For a binary phylogenetic tree $T$, a collection $Q$ of displayed quartet trees is a *generous cover* for $T$ if for all pairs $u,v$ of interior vertices of $T$, we have a quartet $xy|wz$ in $Q$ that looks like this:



**Theorem** (Dezulian + S, 2003): If $Q$ is a generous cover for $T$, then $qcl_1(T) = Q(T)$. Thus $Q$ defines T.

# Application: the random cluster model

Random process on a phylogenetic tree $\mathcal{T}$. Independently cutting edges with probability $p(e)$ generates, by connectivity, random characters on $\mathcal{T}$.



Cutting the marked edges yields the character $\{a|bk|cghi|d|ef|j\}$.

# Reconstructing *T* from *k* independent characters (bounds and phase transition)

- **Theorem** (Mossel and S, 2003) For random cluster model, if
- $0 < a \leq p(e) \leq p < 0.5$, every binary phylogenetic tree with *n* leaves can be reconstructed with probability at least 1-ε from *k* indep. characters if

$$k \geq \frac{c_{p,\varepsilon} \log(n)}{a}$$

- A fast (polynomial-time) algorithm to reconstruct *T* from the characters.
- Proof uses generous cover result. Doesn't require i.i.d.
- Lower bound: log(*n*) needed (not trivial) and polynomial (*n*) if p>0.5.
- *cf.* finite-state $\qquad k \geq \dfrac{c_\varepsilon (\log(n))^{d(p)}}{a^2}$

# Relevance to finite-state space?

**Corollary:**
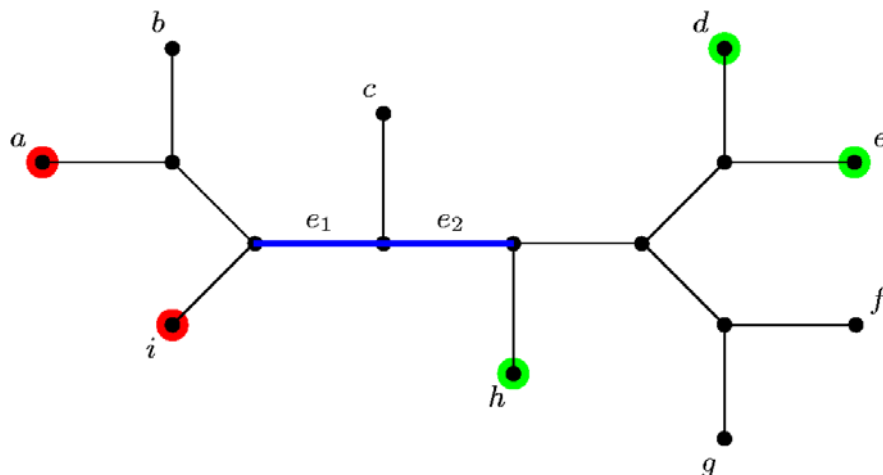
Random walk on group with generating set of size

$$d \geq c_\varepsilon n^2 \log(n)$$

then $T$ can be reconstructed w.p. $> 1-\varepsilon$ with a
$\Theta(\log(n))$ number of characters

# Application 3: $qcl_1$, $qcl_2$, $qcl_{1,2}$

[Definition] A **partial X-split** $A|B$ is a partition of a subset into two non-empty sets, $A,B$. $A|B$ is **displayed** by $T$ if we can remove an edge from $T$ to separate $A$ from $B$.

**Example:** $\{a,i\}|\{d,e,h\}$ is displayed by $T$.

## Meacham's dyadic rules for splits (1983)

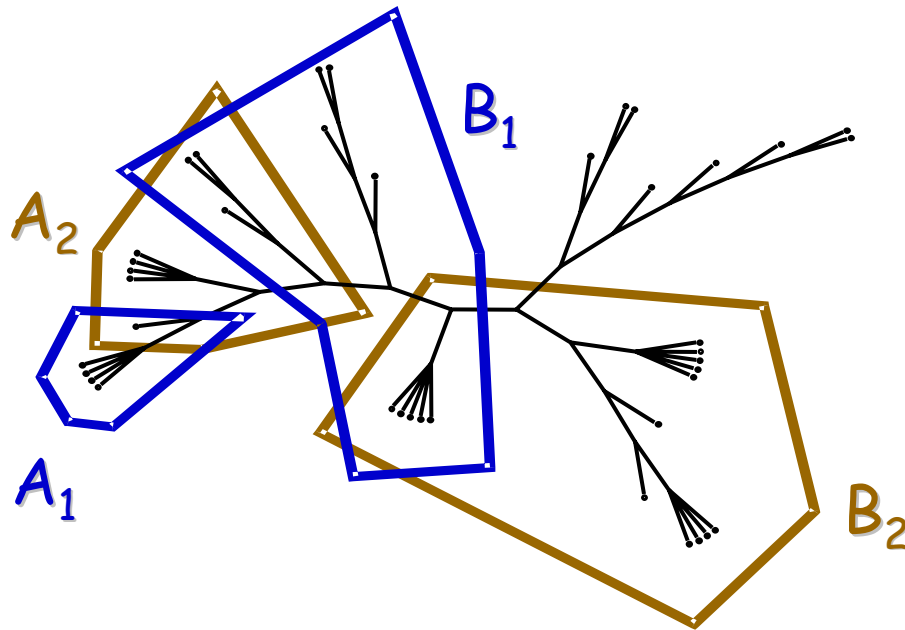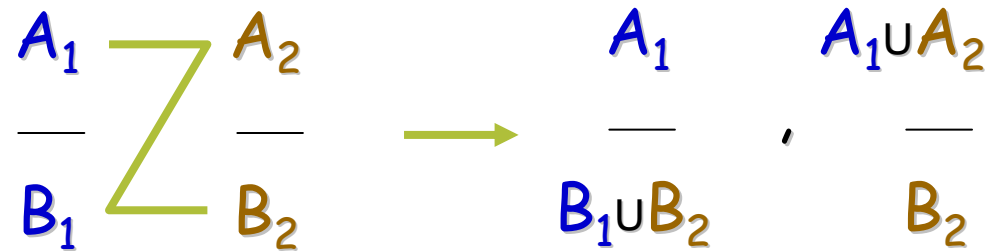(**M1**): If $A_1 \cap A_2 \neq \emptyset$ and $B_1 \cap B_2 \neq \emptyset$ then

$$\{A_1|B_1, A_2|B_2\} \vdash A_1 \cap A_2|B_1 \cup B_2, A_1 \cup A_2|B_1 \cap B_2.$$

(**M2**): If $A_1 \cap A_2 \neq \emptyset$ and $B_1 \cap B_2 \neq \emptyset$ and $A_1 \cap B_2 \neq \emptyset$ then

$$\{A_1|B_1, A_2|B_2\} \vdash A_2|B_1 \cup B_2, A_1 \cup A_2|B_1.$$

> Any phylogenetic $X$−tree that displays the partial $X$−splits on the left of (**M1**) or (**M2**) also displays the corresponding partial $X$−splits on the right.

$$\frac{A_1}{B_1} \diagdown \diagup \frac{A_2}{B_2} \longrightarrow \frac{A_1}{B_1 \cup B_2} , \frac{A_1 \cup A_2}{B_2}$$

# Dyadic split closure

$$\Sigma = \Sigma_1 \subseteq \Sigma_2 \subseteq \cdots \subseteq \Sigma_m = \mathrm{spcl}_\theta(\Sigma)$$

where $\Sigma_{i+1}$ consists of $\Sigma_i$ together with all additional splits that can be obtained from a pair of splits in $\Sigma_i$ by applying the rule(s) allowed by $\theta$.

For $\theta \subseteq \{1, 2\}$, let the dyadic split closure under rule $\theta$, $\mathrm{spcl}_\theta(\Sigma)$, denote the minimal set of splits that contains $\Sigma$ and is closed under rule (**Mi**) for each $i \in \theta$.

We denote these closures with: $\mathrm{spcl}_1(\Sigma), \mathrm{spcl}_2(\Sigma), \mathrm{spcl}_{1,2}(\Sigma)$.

# The (almost) happy marriage

$$
\begin{array}{ccc}
\Sigma & \xrightarrow{\;\;\mathcal{Q}\;\;} & \mathcal{Q}(\Sigma) \\
{\scriptstyle \mathrm{spcl}_\theta} \Big\downarrow & & \Big\downarrow {\scriptstyle \mathrm{qcl}_\theta} \\
\mathrm{spcl}_\theta(\Sigma) & \xrightarrow{\;\;\mathcal{Q}\;\;} & (*)
\end{array}
$$

?

# The (almost) happy marriage

$$\Sigma \xrightarrow{\;\mathcal{Q}\;} \mathcal{Q}(\Sigma)$$

$$\text{spcl}_\theta \Big\downarrow \qquad\qquad \Big\downarrow \text{qcl}_\theta$$

$$\text{spcl}_\theta(\Sigma) \xrightarrow{\;\mathcal{Q}\;} (*)$$

**Theorem 2.1.** *Let $\Sigma$ be a collection of partial $X$–splits. Then,*

$$\text{qcl}_\theta(\mathcal{Q}(\Sigma)) = \mathcal{Q}(\text{spcl}_\theta(\Sigma))$$

*for $\theta = \{1\}$ and $\theta = \{1,2\}$. For $\theta = \{2\}$ we have*

$$\text{qcl}_\theta(\mathcal{Q}(\Sigma)) \subseteq \mathcal{Q}(\text{spcl}_\theta(\Sigma))$$

*and containment can be strict.*

# The closure of a set of quartets

For a compatible set $\mathcal{Q}$ of quartet trees, the *closure* cl($\mathcal{Q}$) is defined as

$$\text{cl}(\mathcal{Q}) = \bigcap_{\mathcal{T} \in \text{co}(\mathcal{Q})} \mathcal{Q}(\mathcal{T})$$

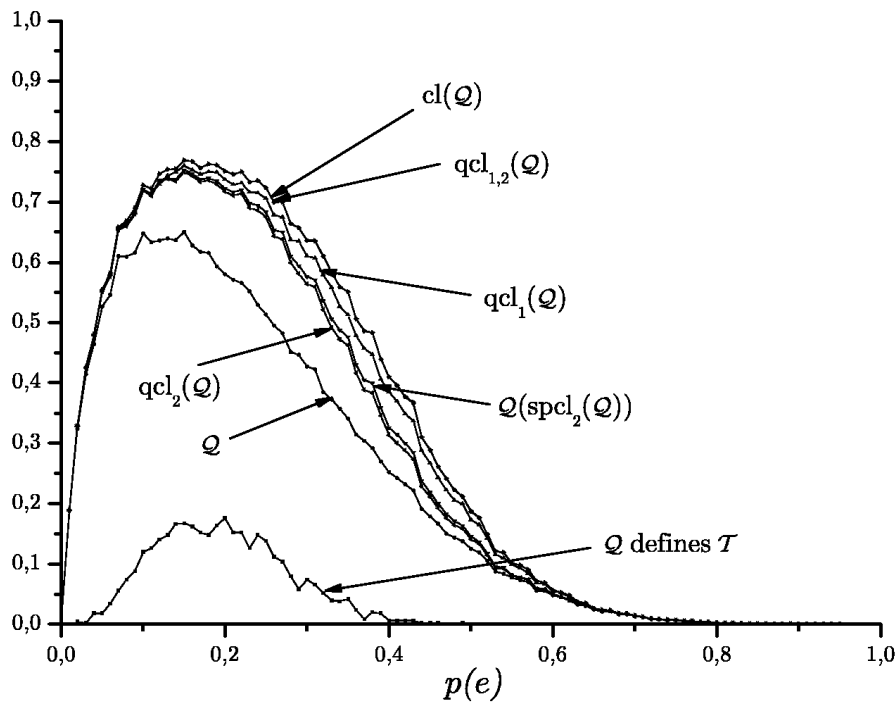where co($\mathcal{Q}$) is the set of phylogenetic trees that display each of the trees in $\mathcal{Q}$.
   Thus cl($\mathcal{Q}$) consists of precisely those quartet trees that are displayed by every phylogenetic tree that displays $\mathcal{Q}$.

• Rules of order $<p$ (for any fixed $p$) do not suffice compute cl($Q$).

• There is a set $Q$ that is incompatible but every strict subset $Q$' is compatible and satisfies cl($Q$')=$Q$'

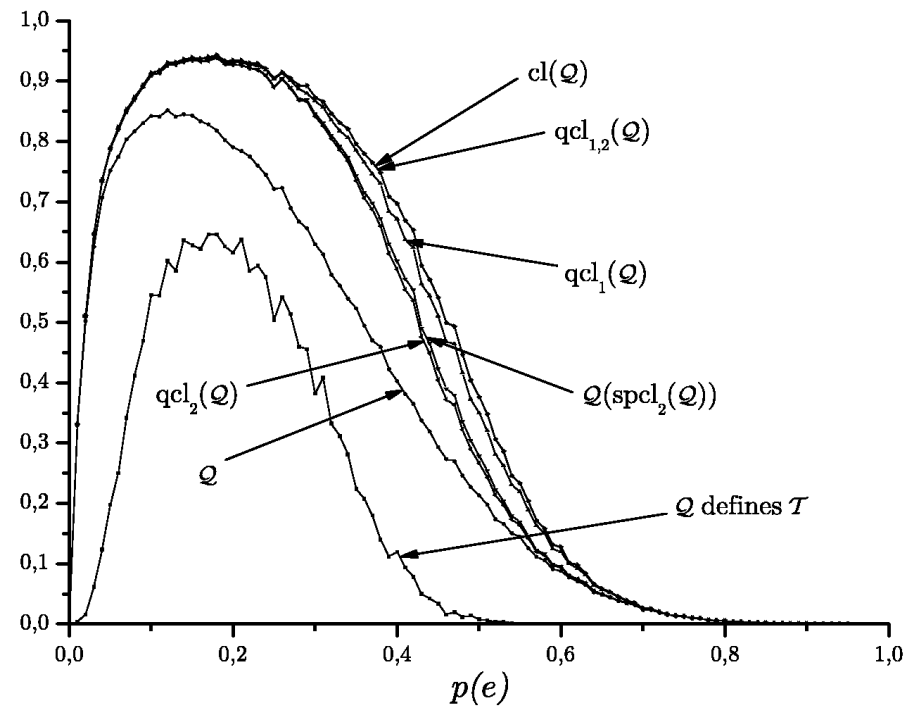# Simulation study

- Random cluster model for binary trees on n=8 leaves.

- *Main question:* How much of cl($Q$) does qcl$_\theta$($Q$) provide? (for $\theta=\{1\}, \{2\}, \{1,2\}$).
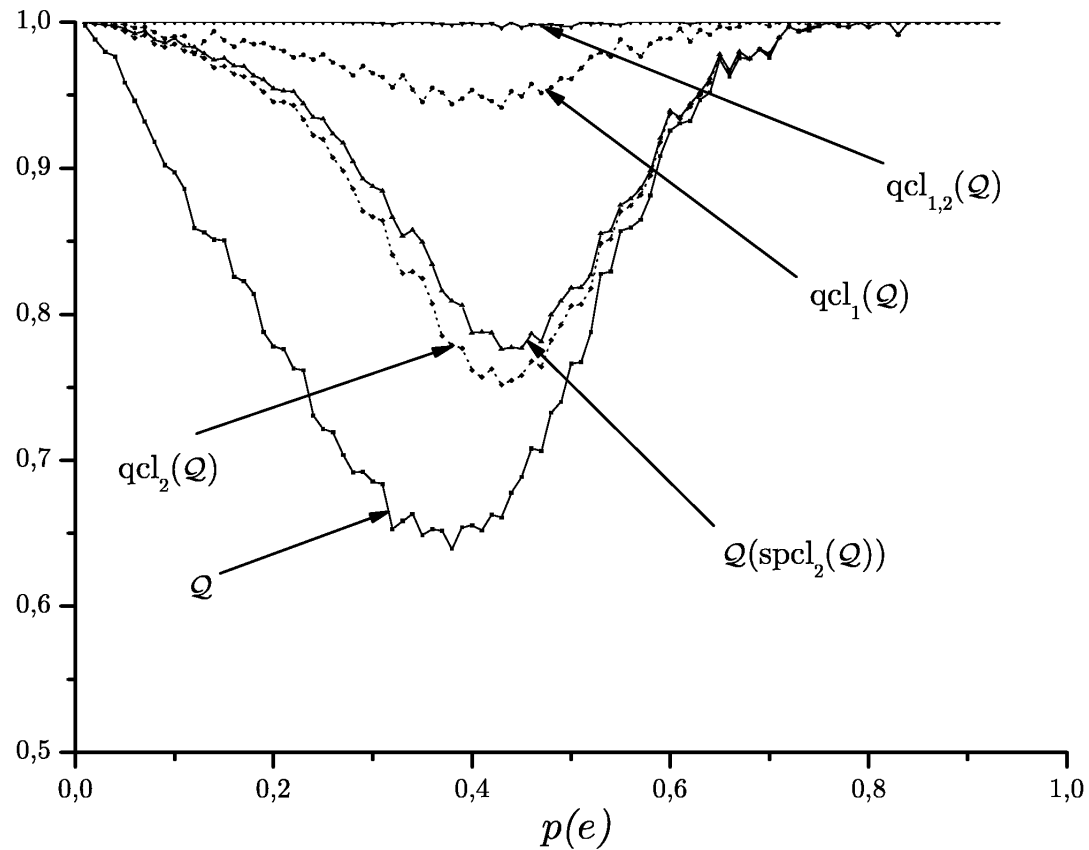
# Simulations 1 (n=8, k=16,32): absolute quartet closure gains



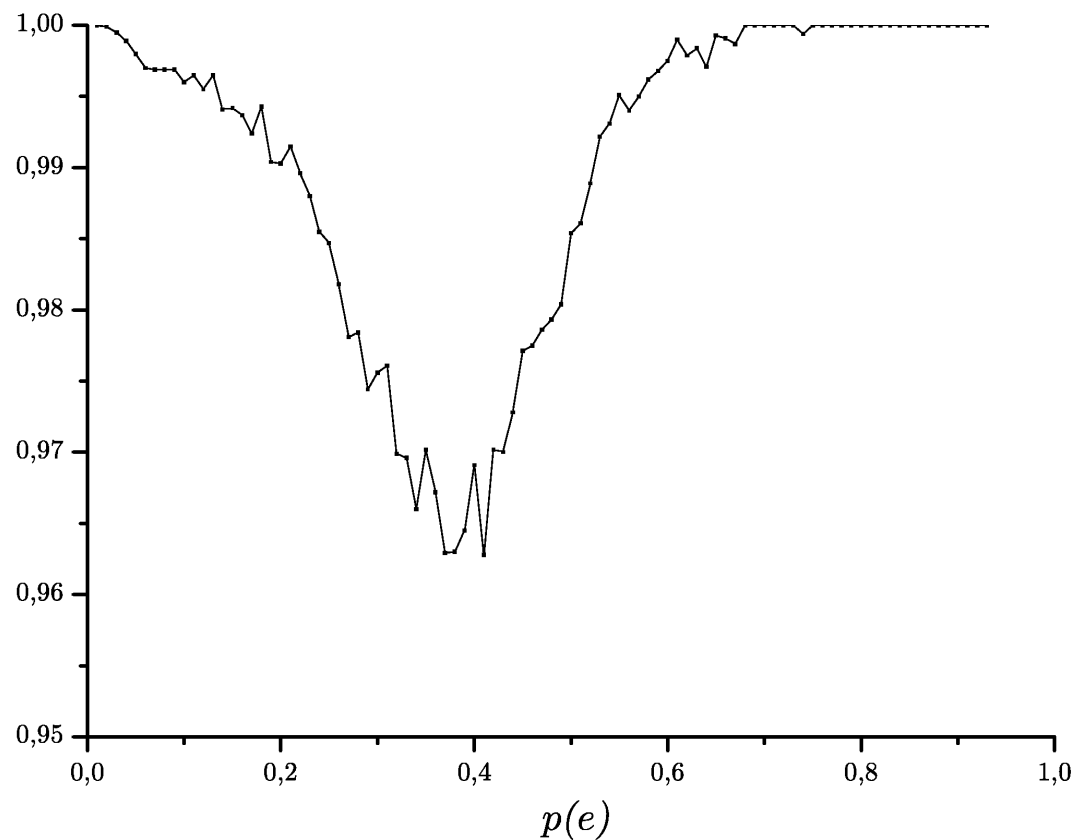(a) The case $k = 16$.

(b) The case $k = 32$.

$$\frac{|\mathcal{Q}|}{\binom{n}{4}}, \frac{|qcl_1(\mathcal{Q})|}{\binom{n}{4}}, \frac{|qcl_2(\mathcal{Q})|}{\binom{n}{4}}, \frac{|\mathcal{Q}(spcl_2(\mathcal{Q}))|}{\binom{n}{4}}, \frac{|qcl_{1,2}(\mathcal{Q})|}{\binom{n}{4}}, \frac{|cl(\mathcal{Q})|}{\binom{n}{4}}$$

# Simulations 2: Relative quartet closure gains



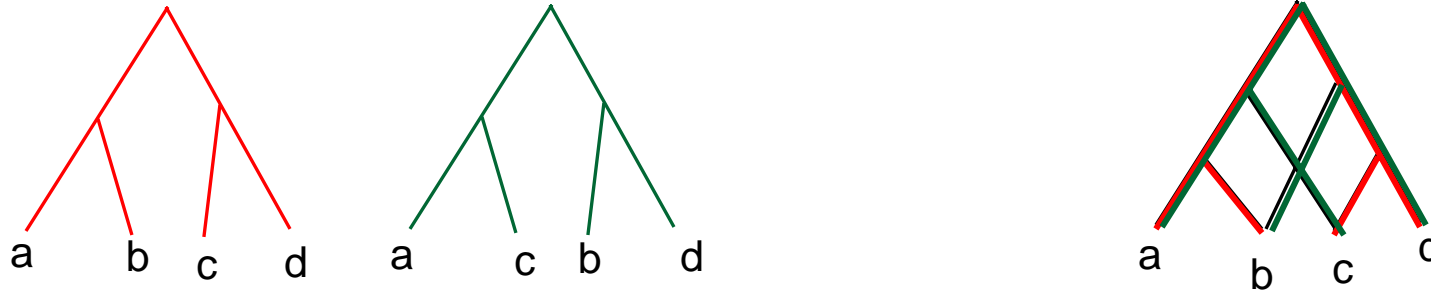$$\frac{|\mathcal{Q}|}{|\mathrm{cl}(\mathcal{Q})|}, \ \frac{|\mathrm{qcl}_1(\mathcal{Q})|}{|\mathrm{cl}(\mathcal{Q})|}, \ \frac{|\mathrm{qcl}_2(\mathcal{Q})|}{|\mathrm{cl}(\mathcal{Q})|}, \ \frac{|\mathcal{Q}(\mathrm{spcl}_2(\mathcal{Q}))|}{|\mathrm{cl}(\mathcal{Q})|}, \ \frac{|\mathrm{qcl}_{1,2}(\mathcal{Q})|}{|\mathrm{cl}(\mathcal{Q})|}$$

# Simulations 3: dyadic closure of splits and quartets comparison



$$\frac{|\mathrm{qcl}_2(\mathcal{Q})|}{|\mathcal{Q}(\mathrm{spcl}_2(\mathcal{Q}))|} \text{ graphed against } p(e)$$
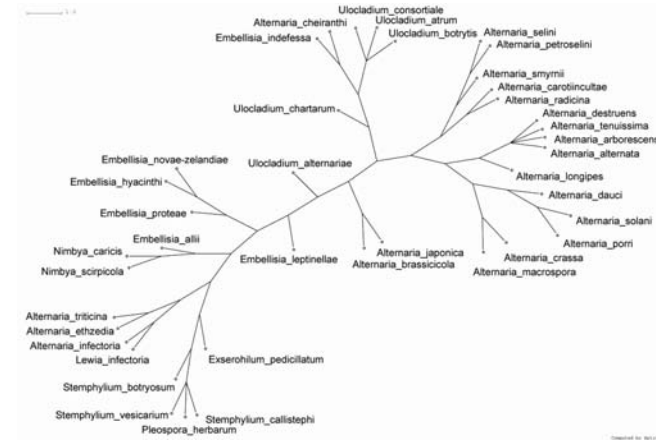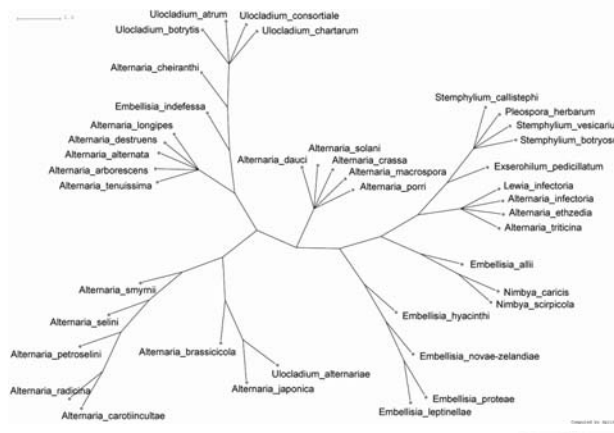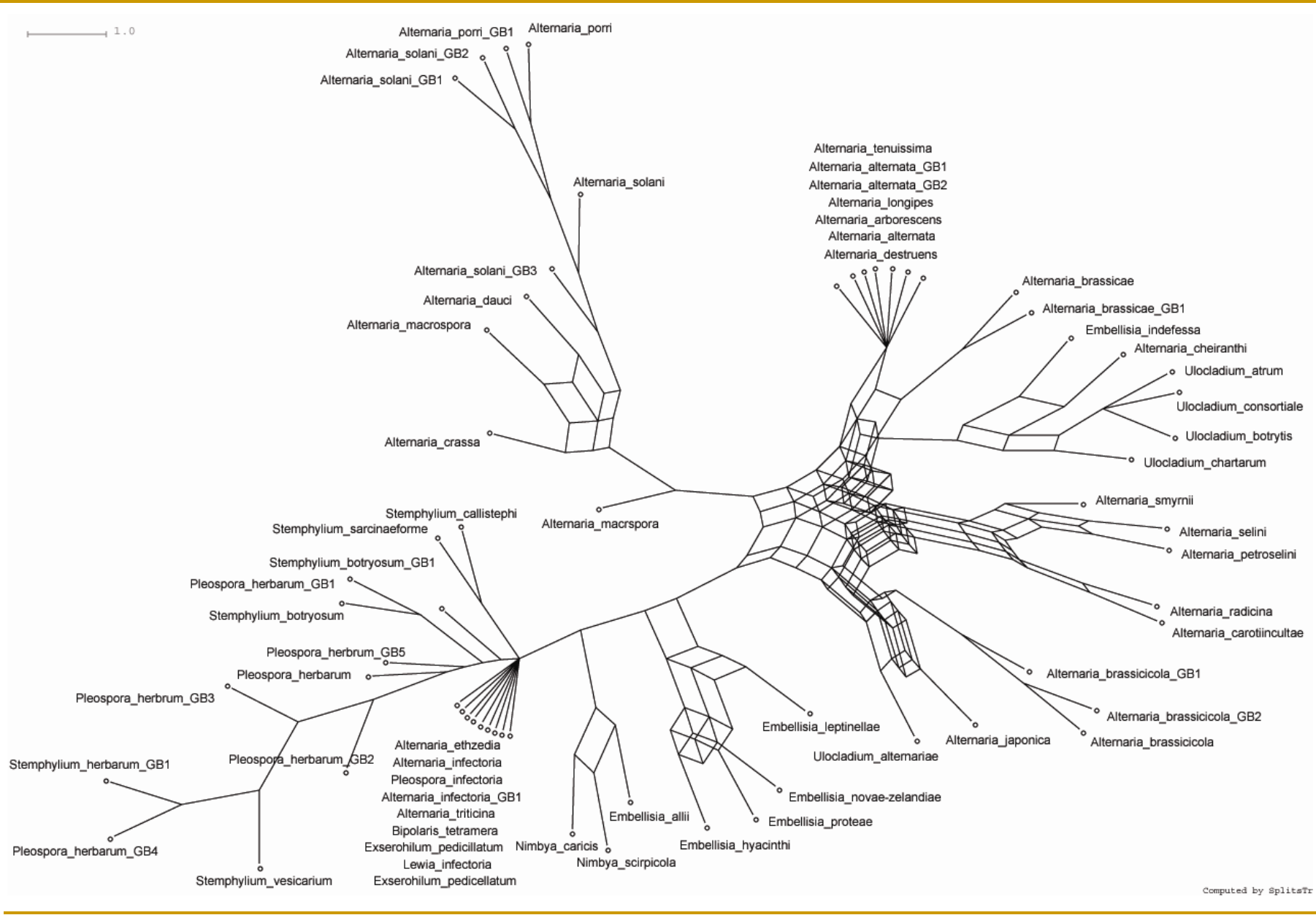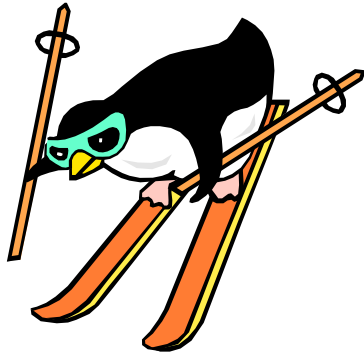
# An application of spcl$_2$



Networks can represent:
- Reticulate evolution (eg. hybrid species)
- Phylogenetic uncertainty (i.e. possible alternative trees)

**Approach:**  Given $T_1, \ldots, T_k$ on overlapping sets of species,
let $\Sigma = \Sigma(T_1) \cup \cdots \cup \Sigma(T_k)$
construct spcl$_2(\Sigma)$ and construct the
'splits graph' of the resulting splits that are 'full'.

# The end

## Further details

•A phase transition for a random cluster model on phylogenetic trees. E. Mossel and M. Steel, *Mathematical Biosciences*, 187 (2004), 189-203.

• Phylogenetic closure operations, and homoplasy-free evolution, T. Dezulian and M. Steel *Proceedings of the International Federation of Classification Societies*, Chicago, 2004.

•Four characters suffice to convexly define a phylogenetic tree. K. Huber, V. Moulton and M. Steel (2003). Submitted.

•How much can evolved characters tell us about the tree that generated them? E. Mossel and M. Steel, Book chapter (Oxford University Press).

•Phylogenetic super-networks from partial trees. D. H. Huson, T. Dezulian, T. Kloepper and M. A. Steel, To appear in *WABI 2004*.