Phylogenetics with distances

Daniel Huson, 2008

# Distance corrections

**Problem**

Hidden mutations (e.g. $\mathbf{A} \to \mathbf{C} \to \mathbf{A}$) mean that we **can't directly observe** the number of mutations between two sequences.

**Solution**

Assume the sequences evolve according to a Markov process and use probability theory to **estimate** the number of hidden mutations.

$$D = -tr(\Pi \log(\Pi^{-1}F))$$

GTR "General time-reversible model"

# Correction formulas

- Most of the standard distance correction formulas can be derived directly from the GTR
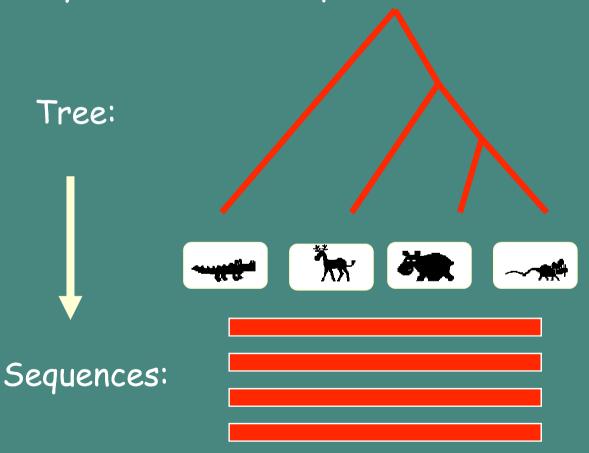- E.g., we obtain the Jukes-Cantor correction as:

$$D = -\frac{3}{4} \log \left( 1 - \frac{4}{3}p \right)$$

# Sequence evolution along one tree

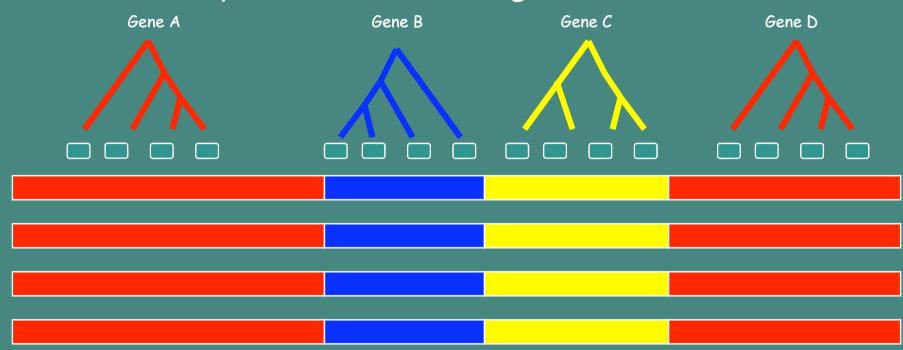We usually assume that sequences evolve on a fixed tree:

Tree:

Sequences:

Standard distance corrections apply to such sequences

# Mosaic sequence evolution

- Can we safely apply standard distance corrections to mosaic sequences?

Trees: $T_1$ $T_2$ ... $T_k$

Sequence proportions: $q_1$ $q_2$ $q_k$

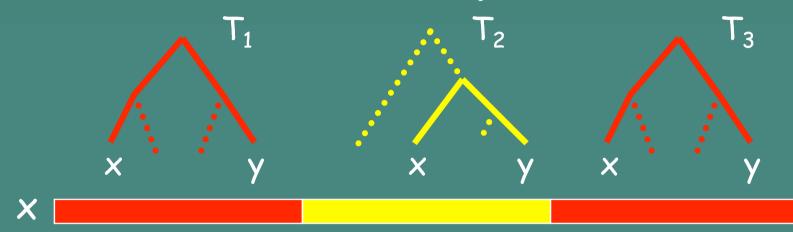Number of mutations between two fixed sequences x and y : $d_1$ $d_2$ $d_k$

True distance between sequences x and y:

$$E[d]=\sum_i q_i d_i \qquad Var[d]=\sum_i q_i (d_i-E[d])^2$$

expected number of mutations

8

# Example

# Main result

Given mosaic sequences. Apply standard correction.

$$E[d] - K \cdot Var[d] \leq \text{corrected distance} \leq E[d]$$

- Corrected distance underestimates true distance.

- If v$\left( K = \dfrac{1}{2} \dfrac{tr(\Pi Q^2)}{t_Q^2} , \right.$
appro

- Corr $\qquad \Pi \qquad$ equilibrium frequencies,
distar where $\quad Q \quad$ rate matrix, and

$\qquad\qquad r_Q \qquad$ mutation rate.$\left. \right)$
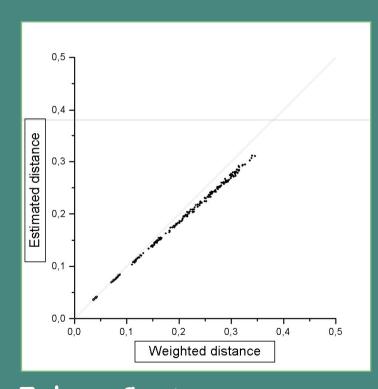
- Proo

# Applications of this result

1) How does undetected recombination effect standard phylogenetic analysis?

2) Consequences for rate variation models?

3) Do network methods explicitly represent recombination?

Daniel Huson, 2003

# Undetected recombination

- Experiments suggest the effect is small:



**Jukes-Cantor
Two trees**

**Distances**

**Kimura 2-p.
Five trees**
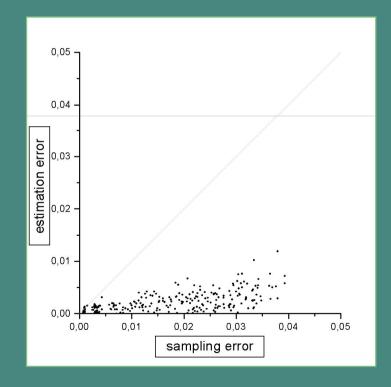
# Undetected recombination

- Experiments suggest the effect is small:



Jukes-Cantor
Two trees

Error

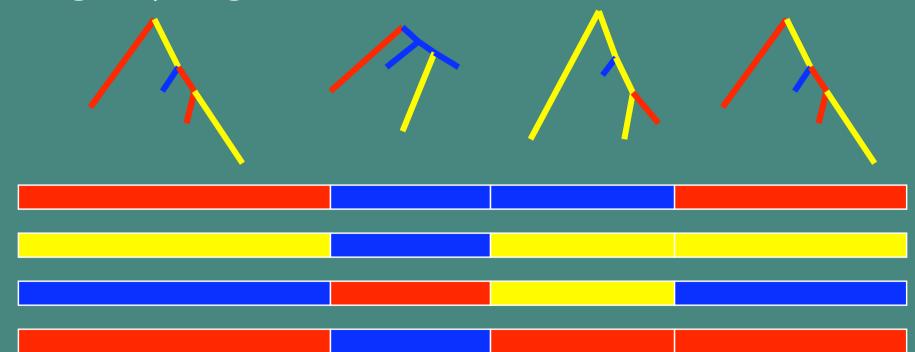Kimura 2-p.
Five trees

# Rate variation

## Site-by-site rate variation



## Different sites have different rates

# Rate variation

## Edge-by-edge rate variation



Different sites have different rates on different edges

# Consistency under rate variation

- Distance and max. likelihood methods can be "inconsistent" when rates vary across sites [J. Chang 1995]
- Our result implies: If

$$Var[d] < \frac{\epsilon}{2K}$$

(Recall: Method is called "consistent", if it converges to the true tree, as the length of available sequence grows longer and longer)

then Neighbor-Joining applied to corrected distances is consistent.

( $\epsilon$ =expected number of mutations on shortest branch)

# Networks



Neighbor-net on human mtDNA

# Trees and splits



Edge *e* corresponds to "split" $\{t_1, t_2, t_6, t_7, t_8\}$ vs $\{t_3, t_4, t_5\}$

Daniel Huson, 2003

# Splits and splits graphs



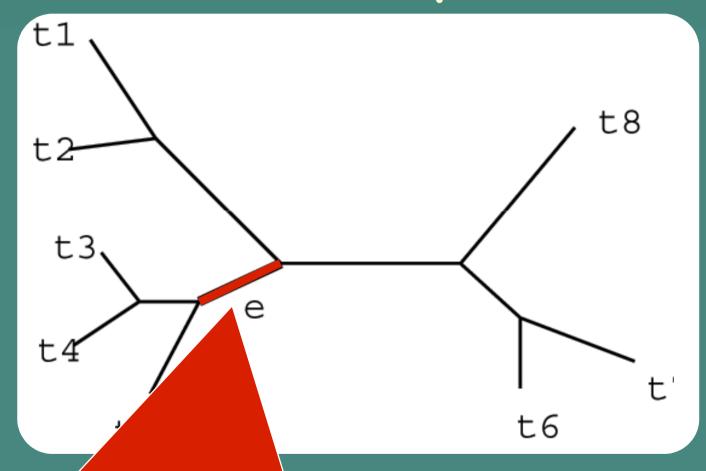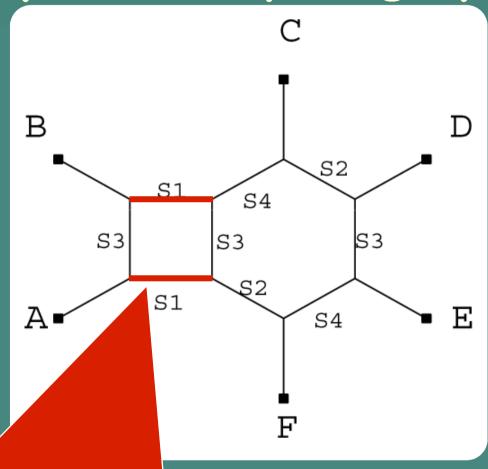Cut-set of parallel edges defines "split" {*A,B*} vs rest

# Mosaic evolution and splits

Trees: $T_1$   $T_2$   ...   $T_k$

Splits sets: $\Sigma_1$   $\Sigma_2$   ...   $\Sigma_k$

Seq. proportions: $q_1$   $q_2$   $q_k$

$$d(x,y) = \sum_{i=1}^{k} q_i \left( \sum_{A|B \in \Sigma_i} b_i(A \mid B)\delta_{A|B}(x,y) \right)$$

$$= \sum_{A|B} \overline{b}(A \mid B)\delta_{A|B}(x,y)$$

where $b_i(A|B)$ is the branch length of $A|B$ in tree $T_i$, and
$\delta_{A|B}(x,y) = 1$, if $A|B$ separates $x,y$, and $0$, else

# Mosaic evolution and splits

Our results imply:

- A splits graph G estimates the set of splits $\Sigma = \cup \Sigma_i$ of the trees $T_1, ..., T_k$
- The lengths of the edges in G estimate the corresponding branch lengths, weighted by the frequencies

# Mosaic evolution and splits

- The split decomposition method [Bandelt & Dress 1992] is consistent when all splits are "weakly compatible"

- The Neighbor-Net method [Bryant & Moulton 2002] is consistent when all splits are "circular"

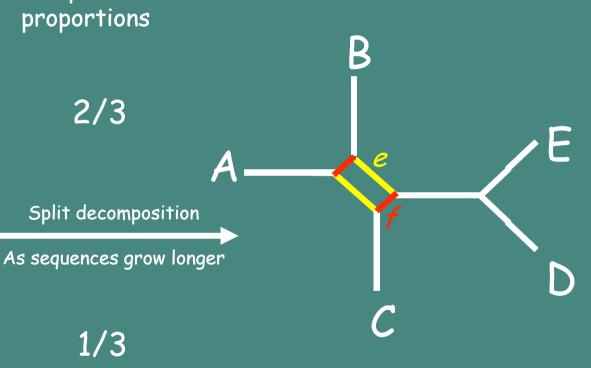$\Rightarrow$ both methods will reconstruct the generating splits and branch lengths, given long enough sequences

# Example

**Two trees**

**Sequence proportions**



2/3

**Split decomposition**

**As sequences grow longer**

1/3

Splits graph containing the splits of both trees

# Summary

- We have established a general result for distance corrections on mosaic sequences

- The effect of using standard distance corrections may not be too bad when the variance is small

- Splits graphs estimate the generating splits and their branch lengths