# Topics of today:

1. Small Parsimony / Big Parsimony

2. Fitch's algorithm

3. Small Parsimony with SCJ

4. Quiz-review

# Small (and big) parsimony problems

$d_M(\mathbb{G}_1, \mathbb{G}_2)$: distance between two genomes under some model $M$

Binary tree T $\begin{cases} \text{topology of T is known} \\ \text{each leaf } u \text{ of T corresponds to a given genome } \mathbb{G}_u \end{cases}$

Function $\mathcal{A}$: assigns a genome $\mathcal{A}(u)$ to each node $u$ of T

$\qquad\qquad$ (if $u$ is a leaf of T, then $\mathcal{A}(u) = \mathbb{G}_u$)

Weight $w(uv)$ of a branch $uv$ of T under $\begin{cases} \text{assignment } \mathcal{A} \\ \text{model } M \end{cases}$ $\quad : w(uv) = d_M(\mathcal{A}(u), \mathcal{A}(v))$

**Small parsimony problem under the model M**:

$\quad$ Find an assignment $\mathcal{A}$ that minimizes the total branch weight of T:

$$\mathcal{W}(T) = \min_{\mathcal{A} \in \mathfrak{A}} \sum_{uv\, \in E(T)} d_M(\mathcal{A}(u), \mathcal{A}(v))$$

Big parsimony problem under the model M:

$\quad$ Given a set of $k$ genomes $\mathbb{G}_1, \mathbb{G}_2, \dots, \mathbb{G}_k$,

$\qquad$ find $\begin{cases} \text{tree T whose } k \text{ leaves are in one-to-one correspondence with the genomes } \mathbb{G}_1, \mathbb{G}_2, \dots, \mathbb{G}_k \\ \text{assignment } \mathcal{A} \text{ of genomes to the nodes of T} \end{cases}$

$\qquad$ minimizing the total branch weight $\mathcal{W}(T)$

# Small parsimony with Fitch's algorithm

General model:

$$\begin{cases} \text{binary rooted tree } \mathtt{T}, \text{ with } n \text{ nodes} \\ \text{genomes are represented by sequences of a fixed length } \ell, \text{ over a finite alphabet } \Sigma \text{ (with } |\Sigma| = m) \\ \text{hamming distance (hd) gives the weight of the branches} \end{cases}$$

Assuming positions being mutually independent, the problem can be solved for each position separately:

The algorithm determines an optimal value for each position $p$ of each node $\mathtt{v}$, denoted by $s_{\mathtt{v}}[p]$

**Bottom-up phase:** defines set $B(\mathtt{v}, p)$ of possible values for each $s_{\mathtt{v}}[p]$, based on $\mathtt{v}$'s children $\mathtt{x}_1$ and $\mathtt{x}_2$

$$\begin{cases} \text{if } \mathtt{v} \text{ is a leaf}: B(\mathtt{v}, p) = \{s_{\mathtt{v}}[p]\} \\ \text{else} \begin{cases} \text{compute } B(\mathtt{v}, p) = B(\mathtt{x}_1, p) \cap B(\mathtt{x}_2, p) \\ \text{if } B(\mathtt{v}, p) = \emptyset : B(\mathtt{v}, p) = B(\mathtt{x}_1, p) \cup B(\mathtt{x}_2, p) \end{cases} \end{cases}$$

Complexity: $O(mn)$ (bottom-up traversal takes $O(n)$, computation of each node takes $O(m)$)

**Top-down phase:** defines final value $s_{\mathtt{v}}[p]$, based on set $B(\mathtt{v}, p)$ and $\mathtt{v}$'s parent $\mathtt{u}$
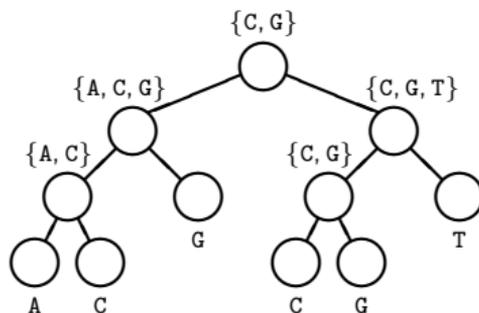
$$\begin{cases} \text{if } \mathtt{v} \text{ has a parent } \mathtt{u} \text{ and } s_{\mathtt{u}}[p] \in B(\mathtt{v}, p) : s_{\mathtt{v}}[p] = s_{\mathtt{u}}[p] \\ \text{otherwise (including the root), arbitrarily assign any value from } B(\mathtt{v}, p) \text{ to } s_{\mathtt{v}}[p] \end{cases}$$

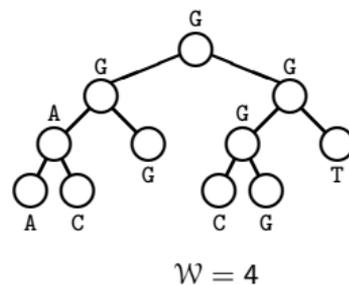Complexity: $O(mn)$ (top-down traversal takes $O(n)$, computation of each node takes $O(m)$)
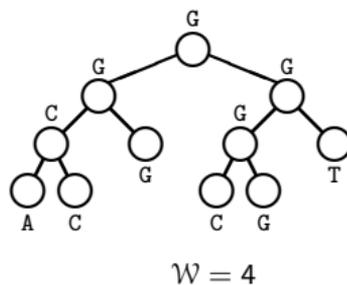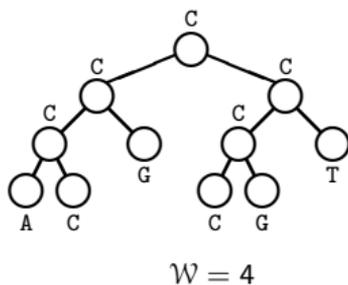
Total complexity: $O(mn)$ per position; with $\ell$ positions: $O(\ell mn)$

# Small parsimony with Fitch's algorithm

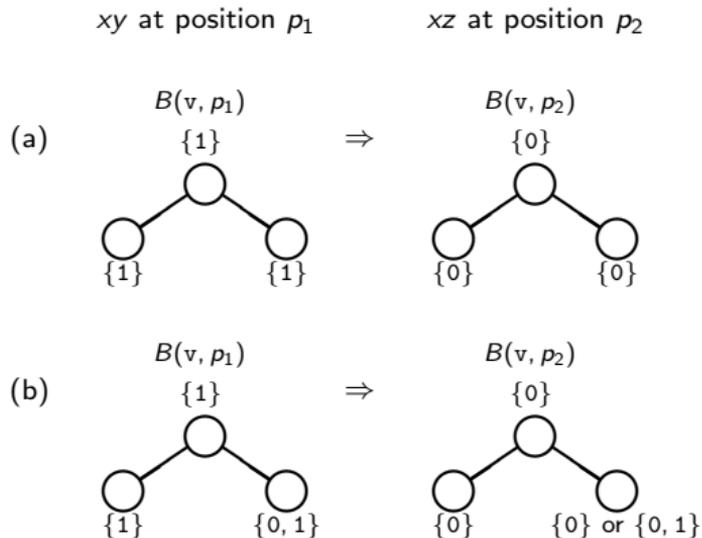Bottom-up phase:



Top-down phase:



$\mathcal{W} = 4$     $\mathcal{W} = 4$     $\mathcal{W} = 4$

# Small parsimony under SCJ with Fitch's algorithm

$$\begin{cases} \text{set of gene families } \mathcal{F} \\[1mm] \text{set of extremities } \xi(\mathcal{F}) = \{f^t : f \in \mathcal{F}\} \cup \{f^h : f \in \mathcal{F}\} \quad (|\xi(\mathcal{F})| = 2|\mathcal{F}|) \\[1mm] \text{set of all possible adjacencies } \widehat{\alpha}(\mathcal{F}) = \{\text{subsets of } \xi(\mathcal{F}) \text{ with size } 2\} \end{cases}$$

$$\Rightarrow \quad |\widehat{\alpha}(\mathcal{F})| = \binom{|\xi(\mathcal{F})|}{2} = \binom{2|\mathcal{F}|}{2} = \frac{2|\mathcal{F}|(2|\mathcal{F}|-1)}{2} = O(|\mathcal{F}|^2)$$

Note: $\widehat{\alpha}(\mathcal{F})$ has many pairs of **conflicting adjacencies** $xy$, $xz$ with $y \neq z$

genome assigned to vertex $u$
$$\begin{cases} \text{represented by a sequence } s_u \text{ of length } \ell = |\widehat{\alpha}(\mathcal{F})| \text{ over the binary alphabet } \{0,1\} \\[2mm] \text{cannot contain conflicting adjacencies:} \\ \quad \text{for each pair of conflicting adj. in } \widehat{\alpha}(\mathcal{F}), \text{ at most one of the two can be set to 1 in } s_u \end{cases}$$

# Small parsimony under SCJ with Fitch's algorithm



Possible bottom-up
scenarios involving
two conflicting adjacencies

Conflicts can be avoided by assigning the value 0 to each ambiguous position of the root genome:

$r$ is the root of T: for each position $p$ of $s_r$ $\begin{cases} s_r[p] = 1 & \text{if } B(r, p) = \{1\} \\ s_r[p] = 0 & \text{otherwise} \end{cases}$

It is easy to verify that: $d_{\text{SCJ}}(u, v) = \text{hd}(s_u, s_v)$

# Close-related problems are NP-hard...

Small parsimony under breakpoint distance

Big parsimony under SCJ

# References

Toward Defining the Course of Evolution: Minimum Change for a Specific Tree Topology

(Walter M. Fitch)

Systematic Zoology, vol. 20, pp. 406–416 (1971)


SCJ: A Breakpoint-Like Distance that Simplifies Several Rearrangement Problems

(Pedro Feijão and João Meidanis)

TCBB volume 8 Number: 5 (2011)

Quiz - Review

# Inversion model 1

1 In the DCJ model any operation reconstructing a target adjacency is optimal, but the same is not true for the inversion model because...

    A a target adjacency can be bad

    B a target adjacency can be already present in the genome

    C reconstructing a target adjacency can be unsafe

2 A cycle is bad when...

    A it cannot be sorted by inversions

    B it interleaves another bad cycle

    C it contains only bad target adjacencies

3 Which data structure helps finding safe inversions?

    A relational diagram

    B overlap graph

    C component tree

4 A bad component can be fixed...

    A with a neutral inversion

    B with a split inversion

    C with a safe inversion

# Inversion model 2

1 Each leaf of the component tree represents...

    A  a bad component

    B  a hurdle

    C  a fortress

2 The cost of covering a component tree can be expressed in terms of...

    A  the number of bad nodes

    B  the length of the longest traversal of the tree

    C  the number of leaves

3 Fixing a super hurdle with a neutral inversion

    A  is a good strategy

    B  creates a new hurdle

    C  destroys a good component

4 Merging two good (or trivial) components..

    A  can merge bad components into a good one

    B  creates a new bad component

    C  is never recommended

# DCJ-indel model 1

1  The indel-potential is defined as...

    A  the number of runs in a component

    B  the smallest number of runs that can be obtained after sorting with internal gaining DCJs

    C  the number of indel-edges in a componenta bad component

2  The indel-potential of a component depends on..

    A  its number of runs

    B  its number of indel-edges

    C  its length

3  The number of runs in a cycle can be...

    A  0,1,2,4,6,8,...

    B  any non-negative integer

    C  any positive integer

4  The number of runs in a path can be...

    A  0,1,2,4,6,8,...

    B  any non-negative integer

    C  any positive integer

# DCJ-indel model 2

1 A recombination can reduce the overall number of runs by at most...

    A 1          B 2          C 3

2 A recombination can reduce the overall overall indel-potential by at most...

    A 1          B 2          C 3

3 A recombination involving a cycle is...

    A gaining

    B neutral

    C losing

4 A recombination involving a cycle can be...

    A deducting

    B part of an optimal sorting sequence

    C none of those

# DCJ-indel model - Path recombinations

**With respect to the endpoints:**

$$\mathbb{A}\!-\!\mathbb{A} \;+\; \mathbb{B}\!-\!\mathbb{B} \left\{ \begin{array}{ll} \mathbb{A}\!-\!\mathbb{B} \;+\; \mathbb{A}\!-\!\mathbb{B} & \text{(gaining)} \\ \mathbb{A}\!-\!\mathbb{B} \;+\; \mathbb{A}\!-\!\mathbb{B} & \text{(gaining)} \end{array} \right.$$

$$\mathbb{A}\!-\!\mathbb{B} \;+\; \mathbb{A}\!-\!\mathbb{B} \left\{ \begin{array}{ll} \mathbb{A}\!-\!\mathbb{A} \;+\; \mathbb{B}\!-\!\mathbb{B} & \text{(losing)} \\ \mathbb{A}\!-\!\mathbb{B} \;+\; \mathbb{A}\!-\!\mathbb{B} & \text{(neutral)} \end{array} \right.$$

$$\mathbb{A}\!-\!\mathbb{A} \;+\; \mathbb{A}\!-\!\mathbb{A} \left\{ \begin{array}{ll} \mathbb{A}\!-\!\mathbb{A} \;+\; \mathbb{A}\!-\!\mathbb{A} & \text{(neutral)} \\ \mathbb{A}\!-\!\mathbb{A} \;+\; \mathbb{A}\!-\!\mathbb{A} & \text{(neutral)} \end{array} \right.$$

**With respect to the runs:**

$$\mathcal{A}\mathcal{B} \;+\; \mathcal{A}\mathcal{B} \left\{ \begin{array}{ll} \mathcal{A}\mathcal{A} \;+\; \mathcal{B}\mathcal{B} & (\Delta_\lambda = -2) \\ \mathcal{A}\mathcal{B}\mathcal{B}\mathcal{A} \;+\; \varepsilon & (\Delta_\lambda = -2) \end{array} \right.$$

$$\mathcal{A}(\mathcal{B}) \;+\; \mathcal{A} \left\{ \begin{array}{ll} \mathcal{A}\mathcal{A} \;+\; (\mathcal{B}) & (\Delta_\lambda = -1) \\ \mathcal{A}\mathcal{A}(\mathcal{B}) \;+\; \varepsilon & (\Delta_\lambda = -1) \end{array} \right.$$

$$(\mathcal{A})\mathcal{B} \;+\; \mathcal{B} \left\{ \begin{array}{ll} (\mathcal{A}) \;+\; \mathcal{B}\mathcal{B} & (\Delta_\lambda = -1) \\ (\mathcal{A})\mathcal{B}\mathcal{B} \;+\; \varepsilon & (\Delta_\lambda = -1) \end{array} \right.$$

**Deducting path recombinations:** $\left\{ \begin{array}{l} \text{gaining with } \Delta_\lambda = -2 \\ \text{gaining with } \Delta_\lambda = -1 \\ \text{neutral with } \Delta_\lambda = -2 \end{array} \right.$

# DCJ-indel model - Path recombinations

**Putting together (examples):**

$$\mathbb{AB}_{\mathcal{AB}} \quad + \quad \mathbb{AB}_{\mathcal{AB}} \quad = \quad \mathbb{AB}_{\mathcal{AB}} \quad + \quad \mathbb{AB}_{\varepsilon}$$
$$\mathbb{A}-\mathcal{AB}-\mathbb{B} \quad + \quad \mathbb{A}-\mathcal{AB}-\mathbb{B} \quad = \quad \mathbb{A}-\mathcal{ABAB}-\mathbb{B} \quad + \quad \mathbb{A}-\mathbb{B} \qquad \text{(neutral, with } \Delta_\lambda = -1)$$

$$\mathbb{AB}_{\mathcal{AB}} \quad + \quad \mathbb{AB}_{\mathcal{AB}} \quad = \quad \mathbb{AB}_{\mathcal{AB}} \quad + \quad \mathbb{AB}_{\varepsilon}$$
$$\mathbb{A}-\mathcal{AB}-\mathbb{B} \quad + \quad \mathbb{A}-\mathcal{AB}-\mathbb{B} \quad = \quad \mathbb{A}-\mathcal{ABAB}-\mathbb{B} \quad + \quad \mathbb{A}-\mathbb{B} \qquad \text{(neutral, with } \Delta_\lambda = -1)$$

$$\mathbb{AB}_{\mathcal{AB}} \quad + \quad \mathbb{AB}_{\mathcal{BA}} \quad = \quad \mathbb{AB}_{\mathcal{A}} \quad + \quad \mathbb{AB}_{\varepsilon}$$
$$\mathbb{A}-\mathcal{AB}-\mathbb{B} \quad + \quad \mathbb{A}-\mathcal{BA}-\mathbb{B} \quad = \quad \mathbb{A}-\mathcal{ABBA}-\mathbb{B} \quad + \quad \mathbb{A}-\mathbb{B} \qquad \text{(neutral, with } \Delta_\lambda = -2)$$

$$\mathbb{AB}_{\mathcal{AB}} \quad + \quad \mathbb{AB}_{\mathcal{BA}} \quad = \quad \mathbb{AB}_{\mathcal{A}} \quad + \quad \mathbb{AB}_{\mathcal{B}}$$
$$\mathbb{A}-\mathcal{AB}-\mathbb{B} \quad + \quad \mathbb{A}-\mathcal{BA}-\mathbb{B} \quad = \quad \mathbb{A}-\mathcal{AA}-\mathbb{B} \quad + \quad \mathbb{A}-\mathcal{BB}-\mathbb{B} \qquad \text{(neutral, with } \Delta_\lambda = -2)$$

$$\mathbb{AB}_{\mathcal{AB}} \quad + \quad \mathbb{AB}_{\mathcal{BA}} \quad = \quad \mathbb{AB}_{\mathcal{B}} \quad + \quad \mathbb{AB}_{\varepsilon}$$
$$\mathbb{A}-\mathcal{AB}-\mathbb{B} \quad + \quad \mathbb{A}-\mathcal{BA}-\mathbb{B} \quad = \quad \mathbb{A}-\mathcal{BAAB}-\mathbb{B} \quad + \quad \mathbb{A}-\mathbb{B} \qquad \text{(neutral, with } \Delta_\lambda = -2)$$

$$\mathbb{AB}_{\mathcal{AB}} \quad + \quad \mathbb{AB}_{\mathcal{BA}} \quad = \quad \mathbb{AB}_{\mathcal{B}} \quad + \quad \mathbb{AB}_{\mathcal{A}}$$
$$\mathbb{A}-\mathcal{AB}-\mathbb{B} \quad + \quad \mathbb{A}-\mathcal{BA}-\mathbb{B} \quad = \quad \mathbb{A}-\mathcal{BB}-\mathbb{B} \quad + \quad \mathbb{A}-\mathcal{AA}-\mathbb{B} \qquad \text{(neutral, with } \Delta_\lambda = -2)$$

| $\Lambda$ | $\lambda$ |
|---|---|
| 0 | 0 |
| 1 | 1 |
| 2 | 2 |
| 3 | 2 |
| 4 | 3 |
| 5 | 3 |
| 6 | 4 |
| 7 | 4 |

# DCJ-indel model - Path recombinations

**Putting together (examples):**

$\mathbb{AB}_{\mathcal{AB}}$  $\mathbb{AB}_{\mathcal{BA}}$  $\mathbb{AB}_{\mathcal{A}}$  $\mathbb{AB}_{\varepsilon}$
$\mathbb{A}-\mathcal{AB}-\mathbb{B}$ $+$ $\mathbb{A}-\mathcal{BA}-\mathbb{B}$ $=$ $\mathbb{A}-\mathcal{ABBA}-\mathbb{B}$ $+$ $\mathbb{A}-\mathbb{B}$ (neutral, with $\Delta_\lambda = -2$)

$\mathbb{AB}_{\mathcal{AB}}$  $\mathbb{AB}_{\mathcal{BA}}$  $\mathbb{AB}_{\mathcal{A}}$  $\mathbb{AB}_{\varepsilon}$
$\mathbb{A}-\mathcal{ABAB}-\mathbb{B}$ $+$ $\mathbb{A}-\mathcal{BA}-\mathbb{B}$ $=$ $\mathbb{A}-\mathcal{ABABBA}-\mathbb{B}$ $+$ $\mathbb{A}-\mathbb{B}$ (neutral, with $\Delta_\lambda = -2$)

$\mathbb{AB}_{\mathcal{AB}}$  $\mathbb{AB}_{\mathcal{BA}}$  $\mathbb{AB}_{\mathcal{A}}$  $\mathbb{AB}_{\varepsilon}$
$\mathbb{A}-\mathcal{ABAB}-\mathbb{B}$ $+$ $\mathbb{A}-\mathcal{BABA}-\mathbb{B}$ $=$ $\mathbb{A}-\mathcal{ABABBABA}-\mathbb{B}$ $+$ $\mathbb{A}-\mathbb{B}$ (neutral, with $\Delta_\lambda = -2$)

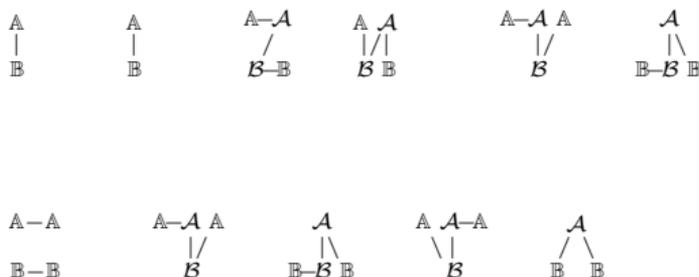| $\Lambda$ | $\lambda$ |
|---|---|
| 0 | 0 |
| 1 | 1 |
| 2 | 2 |
| 3 | 2 |
| 4 | 3 |
| 5 | 3 |
| 6 | 4 |
| 7 | 4 |

# DCJ and DCJ-indel models - Capping

Add caps to close all paths of the graph into cycles, preserving the distance

Canonical capping (no indel edges): maximizes the number of cycles

| paths | linking cycle | Δn | Δc | Δ(2𝔸𝔹) | Δ$_{DCJ}$ |
|---|---|---|---|---|---|
| 𝔸𝔹 | (𝔸𝔹) | +0.5 | +1 | −0.5 | 0 |
| 𝔸𝔸 + 𝔹𝔹 | (𝔸𝔸, 𝔹𝔹) | +1 | +1 | 0 | 0 |
| 𝔸𝔸 | (𝔸𝔸, Γ$_\mathbb{B}$) | +1 | +1 | 0 | 0 |
| 𝔹𝔹 | (𝔹𝔹, Γ$_\mathbb{A}$) | +1 | +1 | 0 | 0 |

Singular capping (with indel edges): optimizes the number of cycles and of runs at the same time



| paths | linking cycle | Δn | Δc | Δ(2𝔸𝔹) | Δλ | Δ$^λ_{DCJ}$ |
|---|---|---|---|---|---|---|
| 𝔸𝔸$_{𝒜�ℬ}$ + 𝔹𝔹$_{𝒜�ℬ}$ | (𝔸𝔸$_{𝒜�ℬ}$, 𝔹𝔹$_{ℬ𝒜}$) | +1 | +1 | 0 | −2 | −2 |
| 2 × 𝔸𝔸$_{𝒜�ℬ}$ + 𝔹𝔹$_{𝒜�ℬ}$ + 𝔹𝔹$_ℬ$ | (𝔸𝔸$_{𝒜�ℬ}$, 𝔹𝔹$_ℬ$, 𝔸𝔸$_{ℬ𝒜}$, 𝔹𝔹$_𝒜$) | +2 | +1 | 0 | −4 | −3 |
| 𝔸𝔹$_{𝒜�ℬ}$ + 𝔸𝔹$_{ℬ𝒜}$ | (𝔸𝔹$_{𝒜�ℬ}$, 𝔸𝔹$_{ℬ𝒜}$) | +1 | +1 | −1 | −2 | −1 |
| 𝔸𝔹 | (𝔸𝔹) | +0.5 | +1 | −0.5 | 0 | 0 |
| 𝔸𝔸 + 𝔹𝔹 | (𝔸𝔸, 𝔹𝔹) | +1 | +1 | 0 | 0 | 0 |