

#### The minimum evolution problem

Daniele Catanzaro Service Graphes et Optimisation Mathématique (G.O.M.) Université Libre de Bruxelles





#### From phylogenetics to molecular phylogenetics





#### From phylogenetics to molecular phylogenetics





# HIV-1 phylogeny





# HIV-1 phylogeny







# Applications

medical research - epidemiology population dynamics - drug discovery

Species	Molecular Sequence
Macaca (A)	AAGCTTCATAGGAGCAACCATTCTAATAATCGCACATGGCCTTACATCATCC
Homo sapiens (B)	AAGCTTCACCGGCGCAGTCATTCTCATAATCGCCCACGGGCTTACATCCTCA
Pan (C)	AAGCTTCACCGGCGCAATTATCCTCATAATCGCCCACGGACTTACATCCTCA
Gorilla (D)	AAGCTTCACCGGCGCAGTTGTTCTTATAATTGCCCACGGACTTACATCATCA
Pongo (E)	AAGCTTCACCGGCGCAACCACCCTCATGATTGCCCATGGACTCACATCCTCC



# Phylogenies





#### Phylogenetic estimation models



The direct use of molecular sequences may lead to under estimation problems.

AAAATCTCTCTCGGTCTCACGG AAATGTGTGTGTGC---CATTTTC ATTTTCTCTCTCC---CTCACGG CCCTGTGTGTGCGGTCATTTCC AAAAT---CTCGGTCTCACGG



Hence, some models of molecular evolution have to be taken into account in order to avoid such problems.

#### Models of molecular evolution



The sequence of a gene can be altered in a number of ways. Gene mutations have varying effects on health depending on where they occur and whether they alter the function of essential proteins. Structurally, mutations can be classified as:

Small-scale mutations, such those as affecting a small gene in one or a few nucleotides, including:

Point mutations, Insertions / Deletions

Large-scale mutations in <u>chromosomal</u> structure, including:

- Gene duplications
- Deletions of large chromosomal regions.
- Chromosomal translocations
- Chromosomal inversions
- Loss of heterozygosity

# Models of molecular evolution



#### Neutrality hypothesis



#### Models of molecular evolution



#### Conservative Hypothesis



#### Models of molecular evolution



#### Superposition principle



# $p_{ij}(t+dt)=\Sigma p_{ik}(t)p_{kj}(dt)$

#### Models of molecular evolution





#### AAATCGGC

R

Time

#### CAATCGGT

.

#### CCATCGGT

#### CCATCGTT

Neutral Selection
 Superposition Principle
 Conservative Hypothesis
 Constant instantaneous rates

# Models of molecular evolution



$$f_{ir} \underbrace{f_{c_{1}}}_{r_{r_{1}}} \underbrace{f_{c_{2}}}_{r_{r_{1}}} \underbrace{f_{c_{2}}}_{r_{r_{2}}} f_{c_{2}}}_{r_{r_{2}}} f_{c_{2}}}_{r_{r_{2}}} f_{c_{2}}} f_{c_{2}}$$

$$f_{r_{1}} \underbrace{f_{c_{2}}}_{r_{r_{2}}} \underbrace{f_{c_{2}}}_{r_{r_{2}}}}_{r_{r_{2}}} \underbrace{f_{c_{2}}}_{r_{r_{2}}} f_{c_{2}}}_{r_{r_{2}}} f_{c_{2}}} f_{c_{2}}$$

$$f_{r_{1}} \underbrace{f_{c_{2}}}_{r_{r_{2}}} \underbrace{f_{c_{2}}}_{r_{r_{2}}}}_{r_{r_{2}}} \underbrace{f_{c_{2}}}_{r_{r_{2}}} \underbrace{f_{c_{2}}}_{r_{r_{2}}}}_{r_{r_{2}}} \underbrace{f_{c_{2}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}} \underbrace{f_{c_{2}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}} \underbrace{f_{c_{2}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}} \underbrace{f_{c_{2}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}} \underbrace{f_{c_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}} \underbrace{f_{c_{2}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}} \underbrace{f_{c_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}} \underbrace{f_{c_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}} \underbrace{f_{c_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}} \underbrace{f_{c_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}\underbrace{f_{c_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{r_{2}}}}_{r_{2}}}}_{r_{2}}}}_{r_{r_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}\underbrace{f_{c_{2}}}}_{r_{2}}}$$

# Models of molecular evolution



Models of molecular evolution are quite important. For example they can be used for computing distances.

AAAATCTCTCTCGGTCTCACGG AAATGTGTGTGCACACATTTTC

$$\mathbf{F}^{\text{#}} = \begin{pmatrix} \mathbf{f}_{AA} & \mathbf{f}_{AC} & \mathbf{f}_{AG} & \mathbf{f}_{AT} \\ \mathbf{f}_{CA} & \mathbf{f}_{CC} & \mathbf{f}_{CG} & \mathbf{f}_{CT} \\ \mathbf{f}_{GA} & \mathbf{f}_{GC} & \mathbf{f}_{GG} & \mathbf{f}_{GT} \\ \mathbf{f}_{TA} & \mathbf{f}_{TC} & \mathbf{f}_{TG} & \mathbf{f}_{TT} \end{pmatrix}$$

 $\mathbf{\Pi}\mathbf{P}(t) = \mathbf{P}(t)^T \mathbf{\Pi}$  $\mathbf{P}(\hat{t}) = \mathbf{P} = \mathbf{\Pi}^{-1}(\mathbf{P}(\hat{t})^T \mathbf{\Pi}) = \mathbf{\Pi}^{-1}\mathbf{F}^{\#}(\hat{t}) = \mathbf{\Pi}^{-1}\mathbf{F}^{\#} - \mathbf{\Pi}^{-1}\mathbf{F}^{\#} - \mathbf{\Pi}^{-1}\mathbf{F}^{\#}(\hat{t}) = \mathbf{\Pi}^{-1}\mathbf{F}^{\#} - \mathbf{\Pi}^{-1}\mathbf{F}^{\#}(\hat{t}) = \mathbf{\Pi}^$ 

 $\mathbf{P}(t) = \mathbf{e}^{\mathbf{R}t} = \mathbf{\Omega}\mathbf{e}^{\mathbf{\Lambda}t}\mathbf{\Omega}^{-1}$ 

 $\hat{t} = -trace[\mathbf{\Pi log}(\mathbf{P})]$   $\log(\mathbf{A}) = \mathbf{\Psi} \log(\mathbf{\Lambda}) \mathbf{\Psi}$ 

Lanave et al. (1984) presented a general model of DNA sequence evolution.

Tavaré (1986), Barry & Hartigan (1987), Rodriguez et al. (1990) gave a different but numerically and algebraically equivalent formulation. Gillespie (1986), Zharkikh (1994), Waddell (1995) noted the <u>time-reversibility</u> (TR) of the Lanave's model; Swofford and Lewis (1997) provided a proof. Waddell & Steel (1997) summarized and extended results on the GTR model and, starting from Rodriguez's formulation, provided algorithms currently implemented in PAUP\*.

#### Computing distances







The cell

			2nd ba	ase	
		U	С	Α	G
	U	UUU (Phe/F)Phenylalanine UUC (Phe/F)Phenylalanine UUA (Leu/L)Leucine UUG (Leu/L)Leucine	UCU (Ser/S)Serine UCC (Ser/S)Serine UCA (Ser/S)Serine UCG (Ser/S)Serine	UAU (Tyr/Y)Tyrosine UAC (Tyr/Y)Tyrosine UAA Ochre ( <i>Stop</i> ) UAG Amber ( <i>Stop</i> )	UGU (Cys/C)Cysteine UGC (Cys/C)Cysteine UGA Opal ( <i>Stop</i> ) UGG (Trp/W)Tryptophan
1st	с	CUU (Leu/L)Leucine CUC (Leu/L)Leucine CUA (Leu/L)Leucine CUG (Leu/L)Leucine	CCU (Pro/P)Proline CCC (Pro/P)Proline CCA (Pro/P)Proline CCG (Pro/P)Proline	CAU (His/H)Histidine CAC (His/H)Histidine CAA (GIn/Q)Glutamine CAG (GIn/Q)Glutamine	CGU (Arg/R)Arginine CGC (Arg/R)Arginine CGA (Arg/R)Arginine CGG (Arg/R)Arginine
base	A	AUU (IIe/I)Isoleucine AUC (IIe/I)Isoleucine AUA (IIe/I)Isoleucine AUG (Met/M)Methionine, Start <sup>[1]</sup>	ACU (Thr/T)Threonine ACC (Thr/T)Threonine ACA (Thr/T)Threonine ACG (Thr/T)Threonine	AAU (Asn/N)Asparagine AAC (Asn/N)Asparagine AAA (Lys/K)Lysine AAG (Lys/K)Lysine	AGU (Ser/S)Serine AGC (Ser/S)Serine AGA (Arg/R)Arginine AGG (Arg/R)Arginine
	G	GUU (Val/V)Valine GUC (Val/V)Valine GUA (Val/V)Valine GUG (Val/V)Valine	GCU (Ala/A) <mark>Alanine</mark> GCC (Ala/A)Alanine GCA (Ala/A)Alanine GCG (Ala/A)Alanine	GAU (Asp/D)Aspartic acid GAC (Asp/D)Aspartic acid GAA (Glu/E)Glutamic acid GAG (Glu/E)Glutamic acid	GGU (Gly/G)Glycine GGC (Gly/G)Glycine GGA (Gly/G)Glycine GGG (Gly/G)Glycine

#### Models of molecular evolution



Amino acid	Hydrophobic effect []	cal-mol]	Molecular	weight [	Dalton]	Surface area	[Å <sup>2</sup> ]	Side-chain volume [Å <sup>3</sup> ]
Alanine	1.0		71			115		88.6
Arginine	1.1		156			225		173.4
Asparagine	-0.1		114			160		114.1
Aspartic acid	-0.1		115			150		111.1
Cysteine	0.0		103			135		108.5
Glutamic acid	0.5		129			190		138.4
Glutamine	0.5		128			180		143.8
Glycine	0.0		57			75		60.1
Histidine	1.3		137			195		153.2
Isoleucine	2.7		113			175		166.7
Leucine	2.9		113			170		166.7
Lysine	1.9		128			200		168.6
Methionine	2.3		131			185		162.9
Phenylalanine	2.3		147			210		189.9
Proline	1.9		97			145		112.7
Serine	0.2		87			115		89
Threonine	1.1		101			140		116.1
Tryptophan	2.9		186			255		227.8
Tyrosine	1.6		163			230		193.6
Valine	2.2		99			155		140

for synonymous substitutions

1/6 for replacement substitutions involving positions (1 and 2) or (2 and 3)

1/8 for replacement substitutions involving positions 1, 2, 3

$$a_{ij}=r_{ij}h_{ij}e^{-\sum_k\omega_k|c_i^k-c_j^k}$$

$$\dot{\mathbf{P}}(t) = \mathbf{P}(t)\mathbf{A} = \mathbf{A}\mathbf{P}(t)$$

1/10 for replacement substitutions to stop codons

#### Models of molecular evolution

Daniele Catanzaro The minimum evolution problem



1

 $h_{ij} =$ 



#### Phylogenetic estimation models





#### Cavalli-Sforza and Edwards: A first model of evolution



Measure of the dissimilarity between species A and B MA1'+W1'B=dAB+eABWA1'+W1'2'+W2'C=dAC+eACWA1'+W1'2'+W2'D=dAD+eADWB1'+W1'2'+W2'C=dBC+eBCWB1'+W1'2'+W2'C=dBC+eBCWB1'+W1'2'+W2'D=dBD+eBDWC2'+W2'D=dCD+eCD



#### Cavalli-Sforza and Edwards: A first model of evolution





Find the best **X** s.t.

**||e||**2 is minimized



#### Cavalli-Sforza and Edwards: A first model of evolution



In other words, Cavalli-Sforza and Edwards proposed the use of the Least-Squares (LS) method, i.e., to find a phylogeny having the lowest distortion from an additive phylogeny.

# $w=(X^{t}X)^{-1}X^{t}d=X^{\dagger}d$

This estimation model is characterized by several drawbacks:

- Species generally do not evolve independently from each others.
- The rate of evolution may not be the same for each species.
- The additive model may provide phylogenies having negative edge weights which is a nonsense.





In other words, Cavalli-Sforza and Edwards proposed the use of the Least-Squares (LS) method, i.e., to find a phylogeny having the lowest distortion from an additive phylogeny.

# $w=(x^{t}x)^{-1}x^{t}d=x^{\dagger}d$

Some authors proposed to

• Consider the evolutionary dependencies between species (Weighted Least-Squares (WLS) and Generalized Least-Squares (GLS)).

• Consider models which allow different evolutionary rates (Minimum Evolution models, Maximum Likelihood (ML) models, Bayesian models (BM)).

• Impose the positivity constraint in order to remove negative edge weights (Projective Algorithms (PA), Minimum Distortion Algorithms (MDA), Balanced Least-Squares (BLS)).

#### Possible solutions





# The Minimum Evolution (ME) model





# The minimum evolution problem











#### Minimum Evolution Problem Under Linear Programming (MELP)

Given a phylogenetic graph and a distance matrix among species, find a phylogeny whose length is minimum.











Assume the set of species  $\Gamma$  is lexicographically ordered, and assume without loss of generality that the rows of the EPT matrix are always ordered lexicographically on the basis of the order in  $\Gamma$ . Assume also that the first n columns of X correspond to the external edges of a phylogeny T and that they are sorted according to the order of the taxa at one of their extremes.



	$e_A$	$e_B$	$e_C$	$e_D$	$e_E$	$e_1$	$e_2$
AB	1	1	0	0	0	0	0
AC	1	0	1	0	0	1	1
AD	1	0	0	1	0	1	1
AE	1	0	0	0	1	1	0
BC	0	1	1	0	0	1	1
BD	0	1	0	1	0	1	1
BE	0	1	0	0	1	1	0
CD	0	0	1	1	0	0	0
CE	0	0	1	0	1	0	1
DE	0	0	0	1	1	0	1

#### The EPT model



Finally, we assume that the remaining (n-3) columns of X, corresponding to the internal edges of T are sorted according to a relation defined in the following way. Given a generic external edge e, define dist(e) as the topological distance of e from the leaf associated with taxon A. In addition, define path(e) the first path, from a lexicographical point of view, to which e belongs. Then, we impose that the column associated with the internal edge e1 precedes the column associated with the internal edge e2 in X if one of the following two conditions holds: path(e1) lexicographically precedes path(e2) or path(e1) = path(e2) and dist(e1) < dist(e2). This order relation is complete as the lexicographic order is complete in the path set, and in a tree we cannot have path(e1) = path(e2) and dist(e1) = dist(e2).



	$e_A$	$e_B$	$e_C$	$e_D$	$e_E$	$e_1$	$e_2$
AB	1	1	0	0	0	0	0
AC	1	0	1	0	0	1	1
AD	1	0	0	1	0	1	1
AE	1	0	0	0	1	1	0
BC	0	1	1	0	0	1	1
BD	0	1	0	1	0	1	1
BE	0	1	0	0	1	1	0
CD	0	0	1	1	0	0	0
CE	0	0	1	0	1	0	1
DE	0	0	0	1	1	0	1

Given this order, any EPT matrix of a phylogeny can be decomposed in blocks:

	$e_A$	$e_B$	$e_C$	$e_D$	$e_E$	$e_1$	$e_2$
AB	1	1	0	0	0	0	0
AC	1	0	1	0	0	1	1
AD	1	0	0	1	0	1	1
AE	1	0	0	0	1	1	0
BC	0	1	1	0	0	1	1
BD	0	1	0	1	0	1	1
BE	0	1	0	0	1	1	0
CD	0	0	1	1	0	0	0
CE	0	0	1	0	1	0	1
DE	0	0	0	1	1	0	1

It is easily seen that only the red block (hereafter indicated as F) is necessary to describe a phylogeny. In fact the bleu and the green blocks can be obtained as xor of the relative yellow and red blocks.



#### The EPT model

In order to represent a phylogeny the entries of matrix F must obey to the following theorem:

Theorem 2.1. F is feasible if and only if all the following conditions hold on its entries:

1. It does not include any of the following  $2 \times 2$  submatrices

$$\mathbf{M}_1 = \begin{pmatrix} 1 & 1\\ 0 & 1 \end{pmatrix} \tag{2.7}$$

or

$$\mathbf{M}_2 = \begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix} \tag{2.8}$$

where the columns and the rows whose intersection defines  $M_1$  and  $M_2$  are not necessarily adjacent in F.

2. Any of its column has at least two entries equal to 1.

3. If columns r and s, with r < s, have an entry equal to 1 in the same row then the number of entries equal to 1 in column r is larger than the number of entries equal to 1 in s.

#### The EPT model



The xor conditions can be expressed as follows:

$$egin{aligned} x_{ij,e} &\leq x_{Ai,e} + x_{Aj,e} \ x_{ij,e} &\leq 2 - x_{Ai,e} - x_{Aj,e} \ x_{ij,e} &\geq x_{Ai,e} - x_{Aj,e} \ x_{ij,e} &\geq -x_{Ai,e} + x_{Aj,e}. \end{aligned}$$

The first condition of Theorem 1 can be expressed as follows:

$$\begin{aligned} x_{Ai,f} + x_{Aj,f} + x_{Ai,e} - x_{Aj,e} &\leq 2\\ x_{Ai,f} + x_{Aj,f} - x_{Ai,e} + x_{Aj,e} &\leq 2. \end{aligned}$$

The second condition of Theorem 1 can be expressed as follows:

$$\sum_{\in V_e \setminus A} x_{Ai,e} \ge 2$$

The third condition of Theorem 1 can be expressed as follows:

$$\sum_{i \in V_e \backslash A} x_{Ai,e} \le n-2$$

$$\sum_{i \in V_e \setminus A} x_{Ai,s} \leq \sum_{i \in V_e \setminus A} x_{Ai,r} - 1 + (n-1)(2 - x_{Aj,r} - x_{Aj,s})$$







$$\min \ z = \sum_e w_e$$

linearizing constraints  $v_{ij,e} \leq d_{ij} x_{ij,e}$  $v_{ij,e} \le w_e$ 

**Biological constraints** 

$$\sum_{e} v_{ije} \ge d_{ij}$$

XOR constraints  $x_{ij,e} \leq x_{Ai,e} + x_{Aj,e}$   $x_{Ai,f} + x_{Aj,f} + x_{Ai,e} - x_{Aj,e} \leq 2$  $x_{ij,e} \ \leq \ 2 - x_{Ai,e} - x_{Aj,e} \quad x_{Ai,f} + x_{Aj,f} - x_{Ai,e} + x_{Aj,e} \ \leq \ 2.$  $x_{ij,e} \geq x_{Ai,e} - x_{Aj,e}$  $x_{ij,e} \geq -x_{Ai,e} + x_{Aj,e}.$ 

Degree constraints

$$\sum_{i \in V_e \setminus A} x_{Ai,e} \leq (n-2)$$

 $\sum_{i \in V_e \setminus A} x_{Ai,s} \leq \sum_{i \in V_e \setminus A} x_{Ai,r} - 1 + (n-1)(2 - x_{Aj,r} - x_{Aj,s})$ 

Anti-cycle constraints

#### The EPT model





. . .

Measure of the dissimilarity between species A and B MAu+ WBu=dAB+eAB WAu+ Wur+ Wrw+ Wwc=dAC+eAC WAu+ Wur+ Wrw+ Wwv+ WvD=dAD+eAD WAu+ Wur+ Wrw+ Wwv+ WvE=dAE+eAE WAu+ Wur+ Wrw+ Wwv+ WvE=dAE+eAE

# A different version of MEP: BME





#### The balanced minimum evolution criterion of phylogenetic estimation



A minimal length phylogeny provides a lower bound on the overall amount of mutation events occurred along evolution of the set of species analyzed.

The balanced minimum evolution criterion is a variation of ME in which the length of a phylogeny is computed as:



#### Fundamentals of the balanced minimum evolution criterion



The phylogeny length under BME is equivalent to the average of the circular orders associated to a given phylogeny.



#### Combinatorial interpretation of BME



The problem of finding a phylogeny which satisfies the balanced minimum evolution criterion is known as Balanced Minimum Evolution Problem (BME) and consists of minimizing the function

$$\mathbf{L} = \sum_{i=1}^{n} \sum_{j=1}^{n} \frac{\mathbf{d}_{ij}}{\mathbf{2}^{\tau_{ij}}}$$

with the constraint that  $\{\tau ij\}$  form a phylogeny.

BME is in P if

$$d_{zi} + d_{kj} \le \max\{d_{zj} + d_{ik}, d_{kz} + d_{ij}\}$$
  
 $d_{ij} \le \max\{d_{ik}, d_{kj}\}$ 

However in the most general case the complexity of BME is unknown.

#### The Balanced Minimum Evolution Problem (BME)





# Approaches to solution: Heuristics





# Approximate algorithm for MEP













#### Ant philosophy







The total number of possible trees with n leaves is (2n-3)!! The number of non-isomorphic shapes increases more slowly! A possible strategy to solve the optimization problem could be: (*i*) to enumerate all the possible non-isomorphic shapes (P# problem); (ii) to find an optimal assignment for each nonisomorphic shape (NP-hard problem).

Leaves	NI-Shapes	Shapes
3	1	3
4	1	15
5	1	105
6	2	945
7	2	10395
8	3	135135
9	4	2027025
10	11	34459425
15	265	1012
20	11020	1021
30	14502229	1038
40	11077270355	1057

#### Non-isomorphic generation







$v_1$	$v_2$	$v_3$	$v_4$	$v_5$	$v_6$	$v_7$	
1	1	0	0	0	0	1	
1	0	1	0	0	1	0	
1	0	0	1	0	0	1	
1	0	0	0	1	1	0	
0	1	1	0	0	1	1	= <b>X</b> ····
0	1	0	1	0	0	0	
0	1	0	0	1	1	1	
0	0	1	1	0	1	1	
0	0	1	0	1	0	0	
0	0	0	1	1	1	1	
	$v_1$ 1 1 1 1 0 0 0 0 0 0 0 0 0	$egin{array}{cccc} v_1 & v_2 \\ 1 & 1 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{array}$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$

It is possible to change the assignment of k leaves on the tree by simply swapping the k corresponding rows of **X**.

The k-OPT local search is therefore easy and fast.

Contra: it seems not possible to get an apriori cost of each swap; consequently each swap requires a complete evaluation of the entire tree (heavy).

#### Species assingment









Ant algorithm





$$p_{ij} = \frac{\alpha \tau_{ij} + (1 - \alpha) \eta_{ij}}{\sum_{q \in \Gamma \setminus \Gamma_{\mathcal{G}_k}} \alpha \tau_{qj} + (1 - \alpha) \eta_{qj}}$$

$$\eta_{ij} = (\Delta_{ij} - u_i - u_j)^{-1}$$

 $\tau_{ij} \leftarrow (1-\rho)\tau_{ij} + \epsilon_{ij} \qquad \quad \epsilon_{ij} = \begin{cases} \kappa \rho / l_{\mathcal{G}_{|\Gamma|}}, \text{ if } w_i \text{ is adjacent to } w_j \text{ in solution } s^{best}; \\ 0, \qquad \text{otherwise.} \end{cases}$ 



ACO details



#### Phylogenetic estimation models



The likelihood criterion states that under many plausible explanations of an observed phenomenon, the one with the highest probability of occurring should be preferred. Hence, under the likelihood criterion, a phylogeny is defined to be optimal (or the most likely) if it has the highest probability of explaining the observed taxa.



#### The likelihood criterion of phylogenetic estimation



The likelihood criterion is introduced by Joe Felsenstein in 1981, as an attempt at solving some "distortion" problems of Cavalli-Sforza and Edwards model. In fact, it is possible to prove that in presence of high mutation rates (evidenced e.g., by a high divergence of the molecular sequences) or convergent/divergent evolution, Cavalli-Sforza and Edwards model leads to edge weights having values quite different from the true phylogeny. This phenomenon is known as "long branch attraction".



#### The likelihood criterion of phylogenetic estimation





# Rooted phylogenies







# Rooted phylogenies







#### Likelihood score of a phylogeny



$$\mathbf{L} = \sum_{j=1}^{4} \mathbf{S} (j) \text{ frequency } (j)$$

 $\{0.1, 0.2, 0.4, 0.4\}$ 

# Likelihood score of a phylogeny



$$\mathbf{L} = \prod_{c=1}^{3} \sum_{j=1}^{4} \mathbf{S} (j) \text{ frequency } (j)$$

#### $\{\{0.1, 0.2, 0.4, 0.4\}, \{0.3, 0.12, 0.2, 0.1\}, \{0.5, 0.1, 0, 0.1\}\}$

t₁

# AAC TGG

 $t_2$ 

# Likelihood score of a phylogeny



maximize L (T, t, P)
T,t,P

#### The estimation problem









#### In the literature...





# Very Large-Scale Neighborhood (VLSN) techniques





 $c_{ij,k} = \|S_k - S_j\| + \|S_k - S_i\|$  $c_{ij,k} = L(S_k | S_i, S_j, w_{ki}, w_{kj})$ 

#### VLSN techniques for phylogeny estimation





#### Minimum Cost Assignment Neighborhood





# Minimum Cost Cycle Neighborhood





# The algorithm

Estimating phylogenies under maximum likelihood: A very large-scale neighborhood approach D. Catanzaro, R. Pesenti, and M. C. Milinkovitch - Université Libre de Bruxelles







THE REPART OF THE PERSON OF TH

Neighborhood	Time (sec.)	Number of calls
NNI <sup>(1)</sup>	6.2±1	328±3
VLSN	3.1±0.12	38±1







Neighborhood	Time (sec.)	Number of calls
NNI <sup>(1)</sup>	n.a.	>149±1
VLSN	413±1	116±4







NUNCERE TEACHER

Phylogenetic estimation is one of the most important problem in computational biology. It is a flourishing area of interaction between molecular biology, operations research, computer science, and physics.

The day by day growing amount of molecular data stored in public databases forces to search for:

- Ad hoc models of molecular evolution
- > Optimization algorithms to select phylogenies among possible alternatives

Here we have presented a first introduction to phylogenetics. The most relevant issues used in tackling real-world sized problems have been outlined, as have the most interesting refinements deserving further research effort.



#### Conclusion





#### The minimum evolution problem

Daniele Catanzaro Service Graphes et Optimisation Mathématique (G.O.M.) Université Libre de Bruxelles

