

Reviews in
Computational Biology

Inferring Lateral Gene Transfer



Christophe Dessimoz

March 28th, 2011

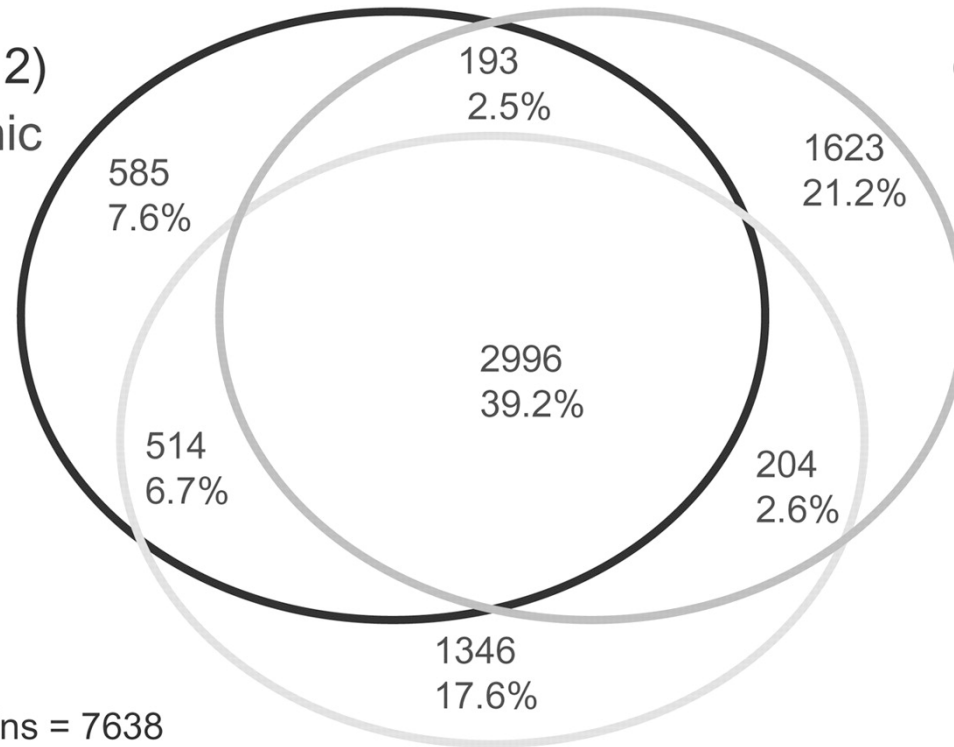
Structure

- **Introduction**
- **Methods to Detect Lateral Gene Transfer**
 - Parametric
 - Phylogenetic
 - Explicit
 - Implicit
- **Outlook**

Shared E. coli proteins.

MG1655 (K-12)
non-pathogenic

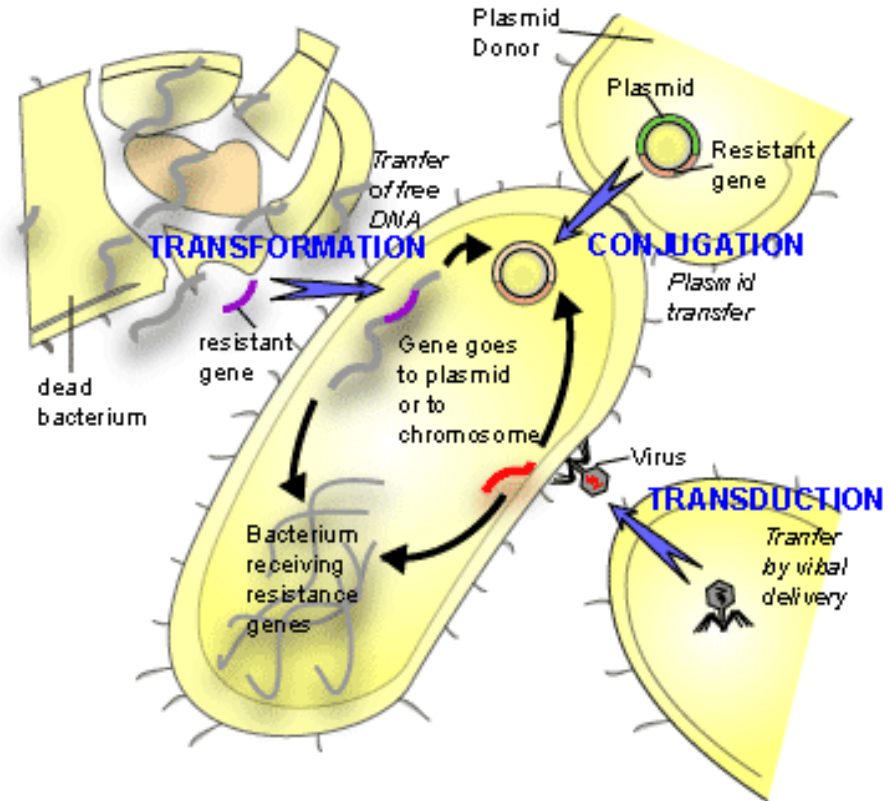
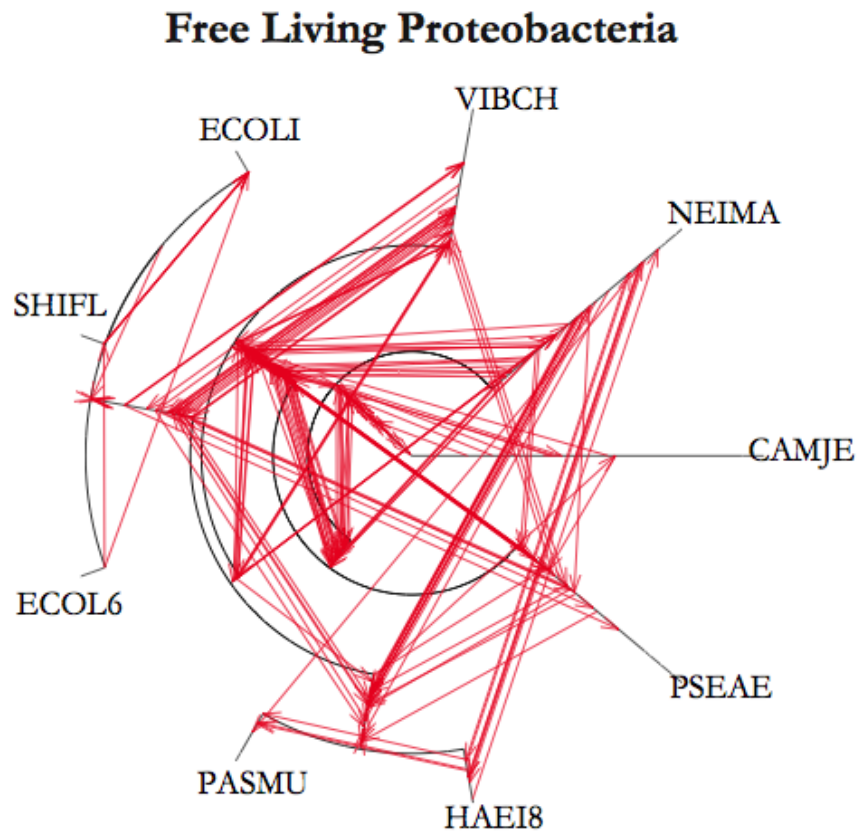
CFT073
uropathogenic



Total proteins = 7638
2996 (39.2%) in all 3
911 (11.9%) in 2 out of 3
3554 (46.5%) in 1 out of 3

EDL933 (O157:H7)
enterohaemorrhagic

Lateral Gene Transfer (LGT)

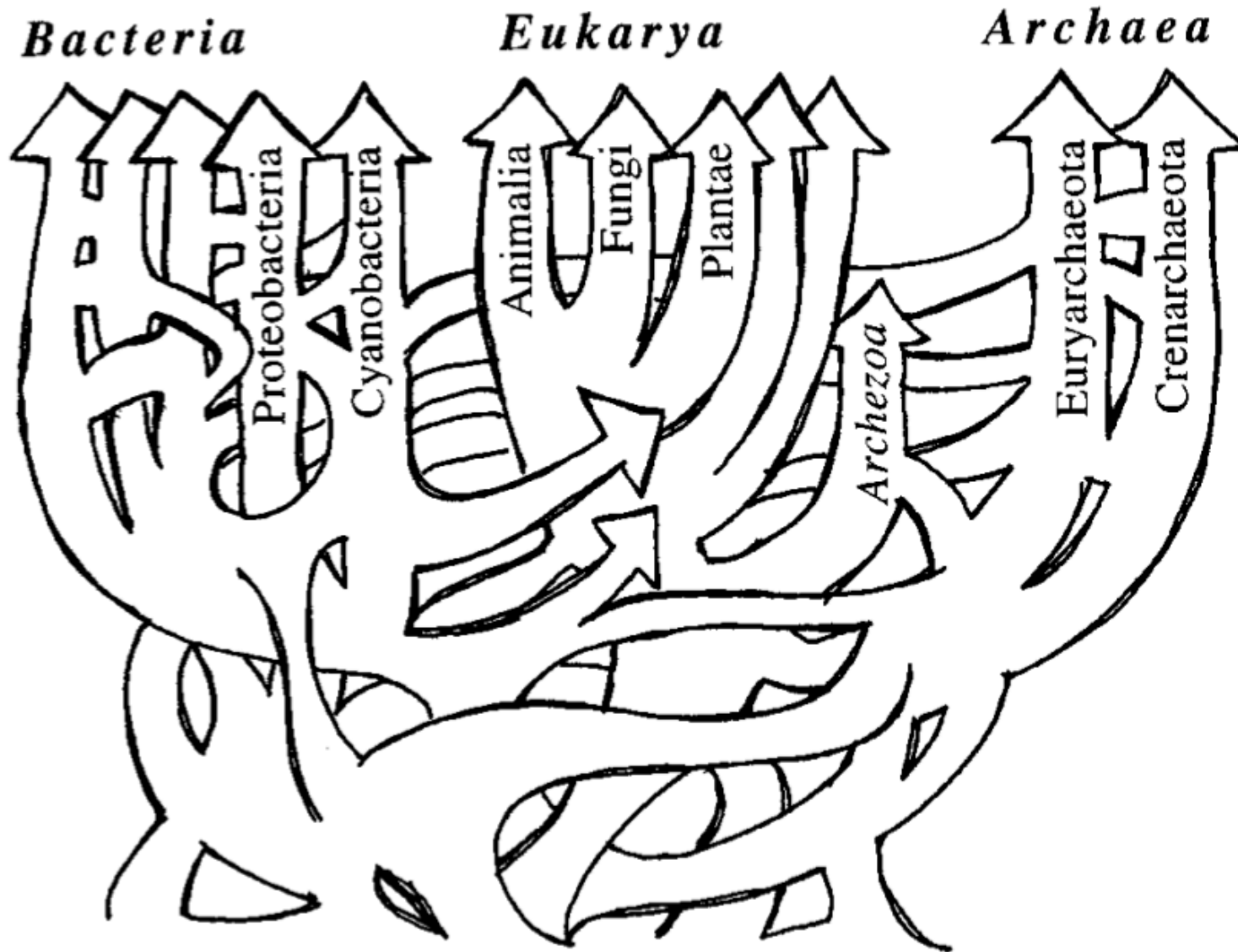


<http://www.scq.ubc.ca/attack-of-the-superbugs-antibiotic-resistance/>

Types of LGTs

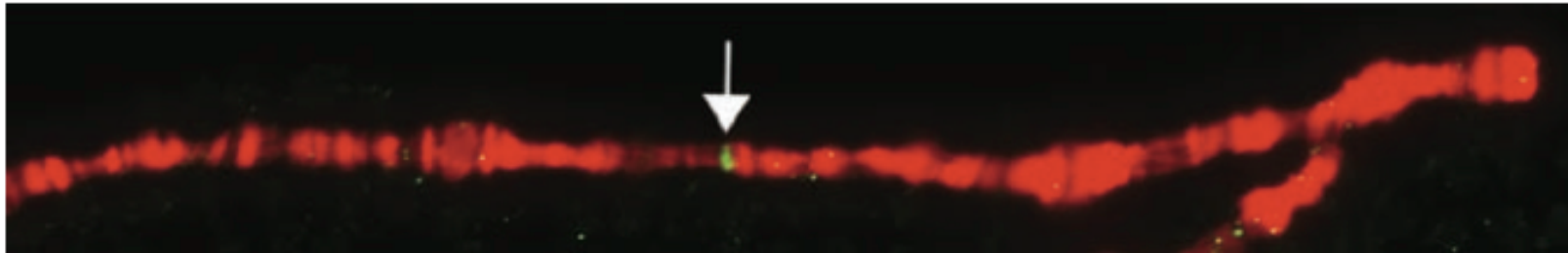
- novel gene acquisition
- orthologous gene replacement

Phylogenetic Classification and the Universal Tree
W. Ford Doolittle, *et al.*
Science 284, 2124 (1999);



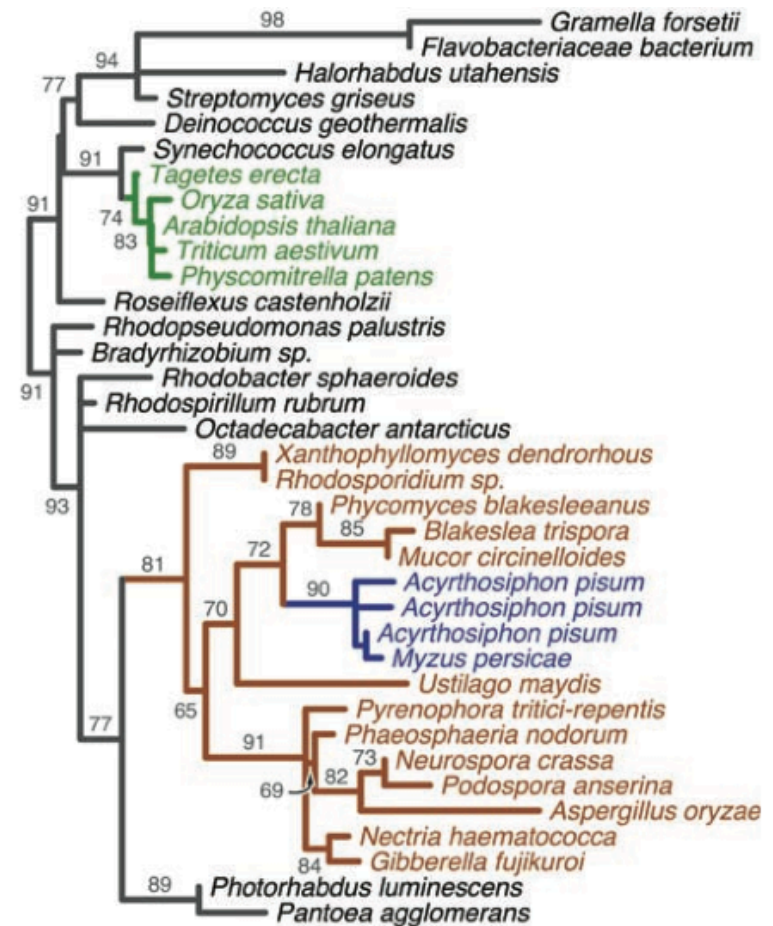
Widespread Lateral Gene Transfer from Intracellular Bacteria to Multicellular Eukaryotes

Julie C. Dunning Hotopp,^{1*†‡} Michael E. Clark,^{2*} Deodoro C. S. G. Oliveira,² Jeremy M. Foster,³
Peter Fischer,⁴ Mónica C. Muñoz Torres,⁵ Jonathan D. Giebel,² Nikhil Kumar,^{1‡}
Nadeeza Ishmael,^{1‡} Shiliang Wang,¹ Jessica Ingram,³ Rahul V. Nene,^{1§} Jessica Shepard,^{1||}
Jeffrey Tomkins,⁵ Stephen Richards,⁶ David J. Spiro,¹ Elodie Ghedin,^{1,7} Barton E. Slatko,³
Hervé Tettelin,^{1‡¶} John H. Werren^{2¶}



Lateral Transfer of Genes from Fungi Underlies Carotenoid Production in Aphids

Nancy A. Moran^{1*} and Tyler Jarvik²



— Bacteria — Plants — Fungi — Aphids

Overview

Parametric methods (based on genome signatures)

- CG-content
- Codon bias
- k-nucleotide frequencies

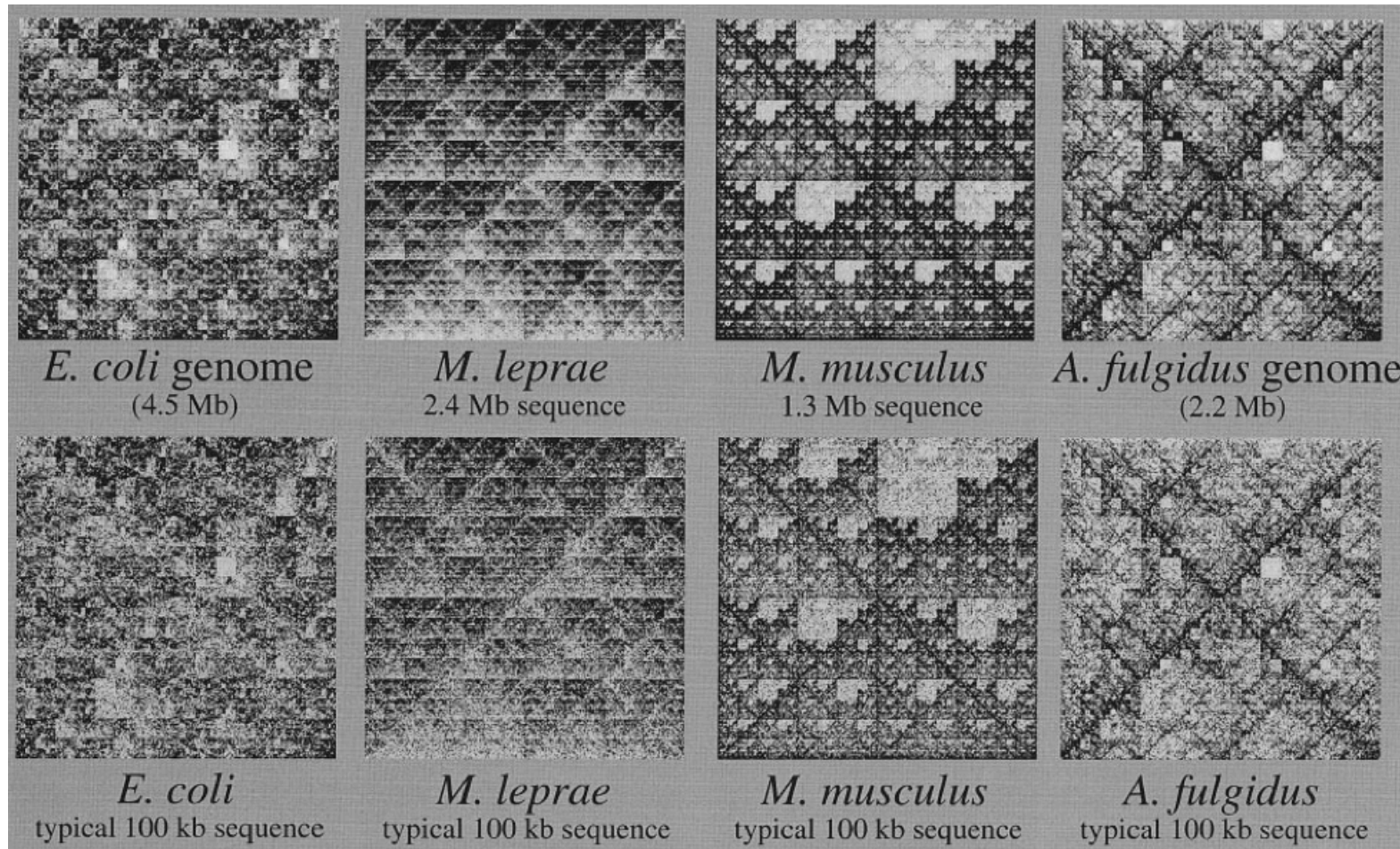
Phylogenetic methods

- Explicit: tree reconciliation
- Implicit: based on underlying tree, but no tree inference/reconciliation

Parametric Methods

Genome signatures

Representation of 7-nucleotide frequencies



Mol. Biol. Evol. 16(10):1391–1399. 1999

Deschavanne et al.

Biased biological functions of horizontally transferred genes in prokaryotic genomes

Yoji Nakamura^{1,5}, Takeshi Itoh^{2,3}, Hideo Matsuda⁴ & Takashi Gojobori^{1,2}

$$P(\text{COD}_1|F) = \frac{P(F|\text{COD}_1) P(\text{COD}_1)}{\sum_{m=1}^6 P(F|\text{COD}_m) P(\text{COD}_m) + P(F|\text{NON}) P(\text{NON})}$$

$\frac{1/12}{1/12} \quad \frac{1/12}{1/2}$

($m = 1, 2, 3, 4, 5, 6$).

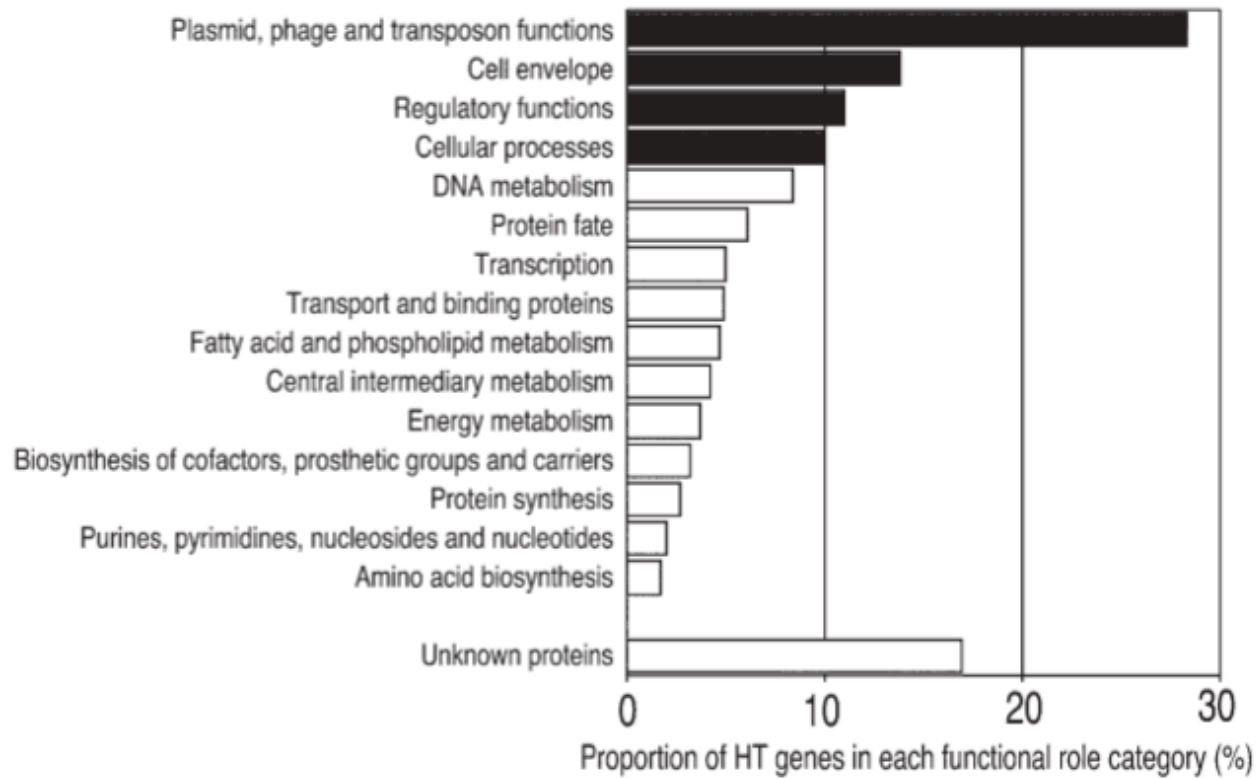
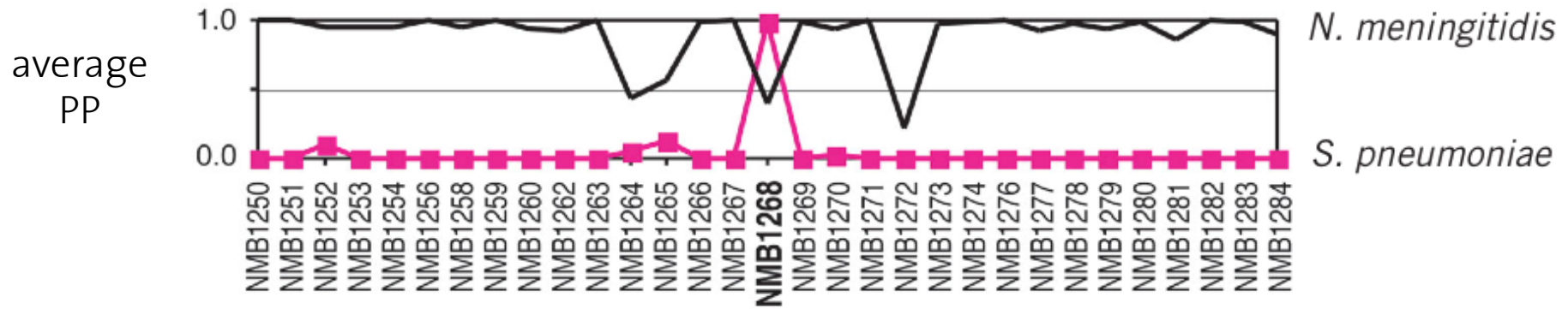
reading frames

F: hexamer

COD_i: coding in frame i

NON: non-coding

The Horizontal Transfer Index is the average $P(\text{COD}_1|F)$ value over a window of size 96 bp, slid on the gene sequence by a step of 12 bp



Parametric Methods: Limitations

- “Amelioration”: adjustment of laterally transferred gene to the nucleotide composition of its new host.
→ restricted to relatively recent transfers
- Donor species are difficult to identify
- Only works for genes transferred from species with significantly different parameters

Amelioration of Bacterial Genomes: Rates of Change and Exchange

Jeffrey G. Lawrence,^{1,*} Howard Ochman²

J Mol Evol (1997) 44:383–397

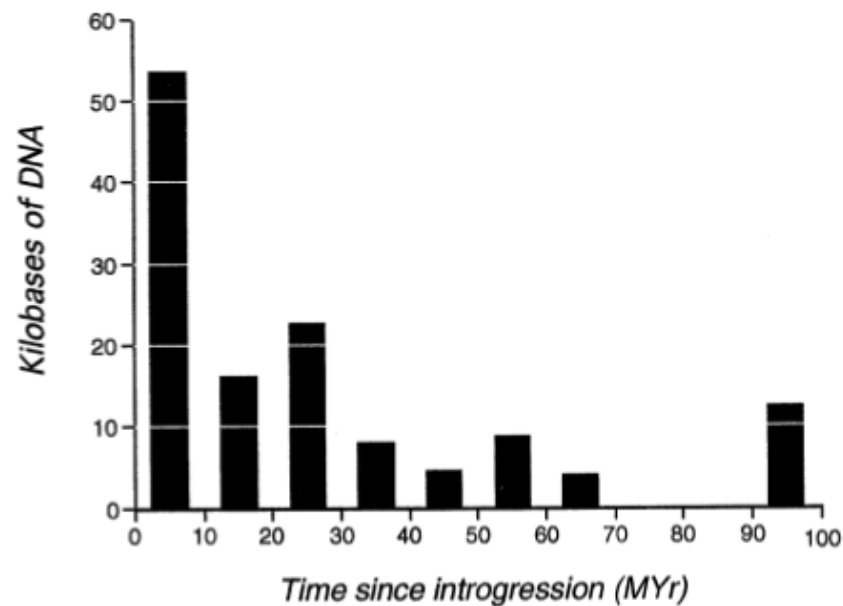
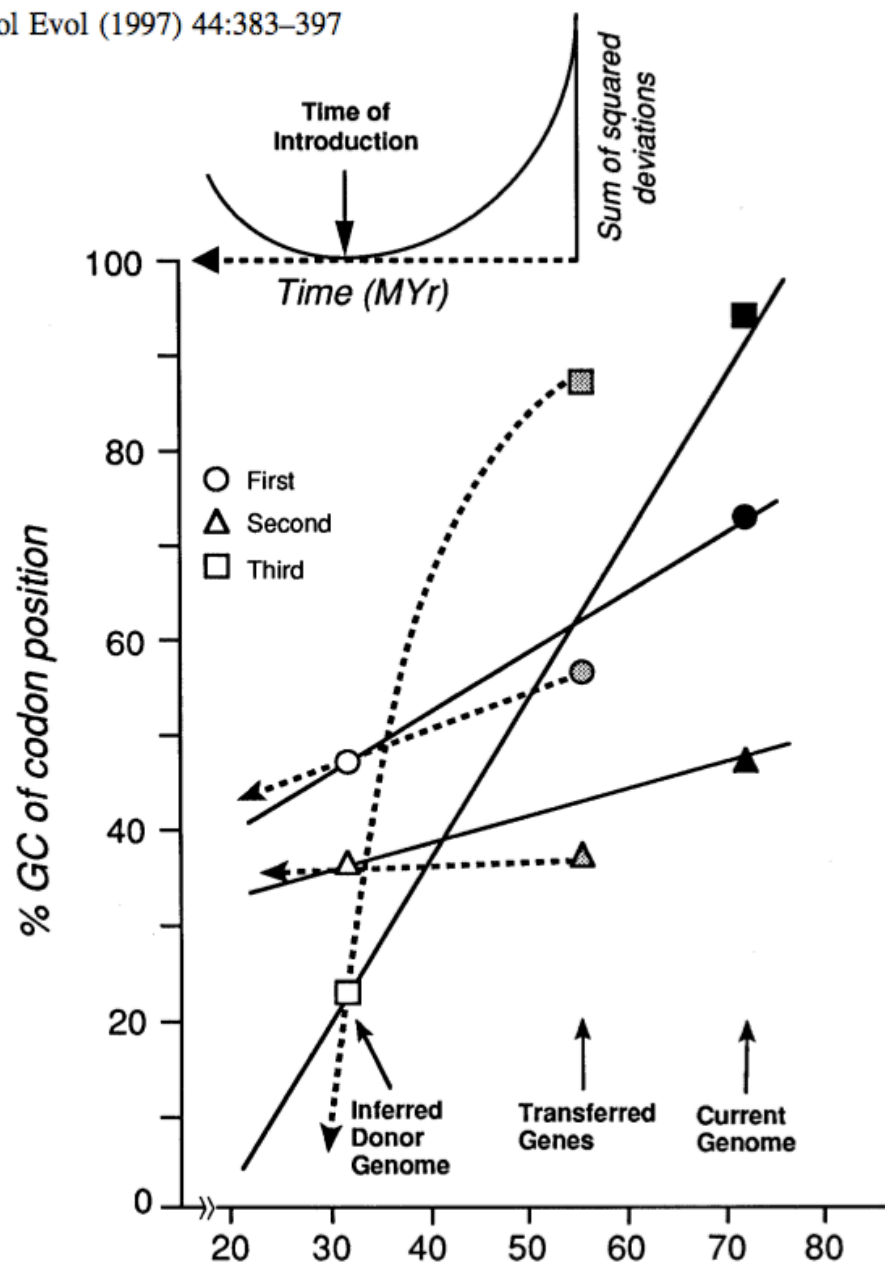


Fig. 7. Times of introgression of DNA fragments into the *E. coli* chromosome. The amount of DNA reflects cumulative pooled sequences (see Table 3).

Phylogenetic Methods

Phylogenetic Methods:

Overview

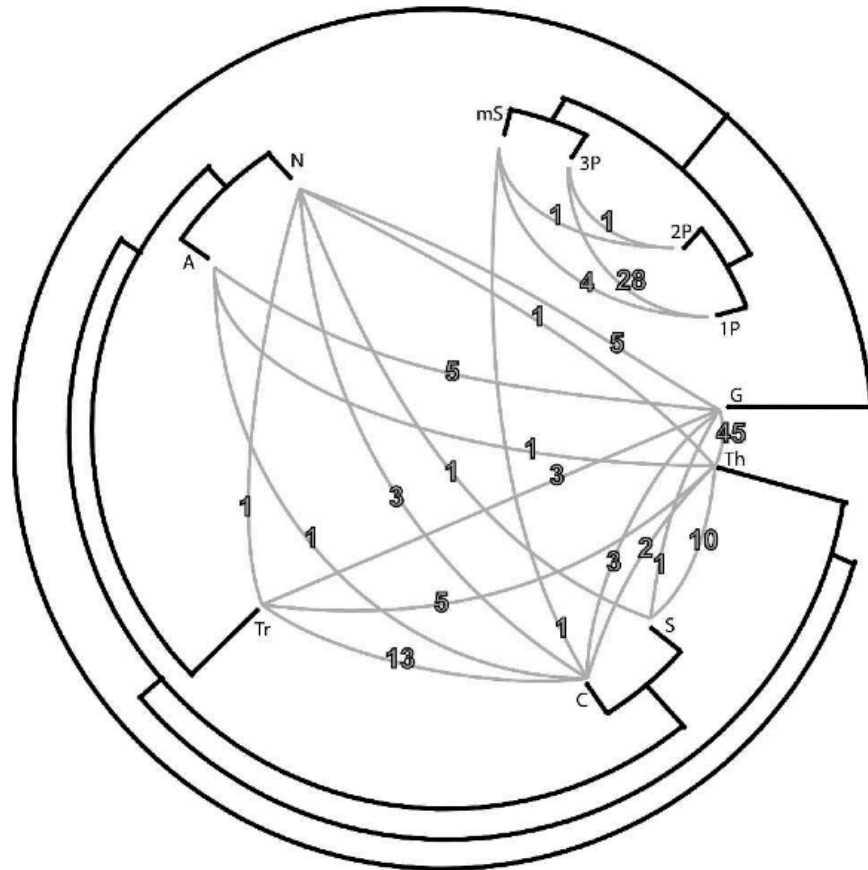
- **Explicit (reconstruct trees)**
 - Discordance
 - Model-based
- **Implicit (rely indirectly on tree)**

Discordance

Phylogenetic analyses of cyanobacterial genomes: Quantification of horizontal gene transfer events

Olga Zhaxybayeva, J. Peter Gogarten, Robert L. Charlebois, et al.

(plurality signal = “species tree”)



But (5-30% FP on simulation w.o. LGT)

are very probable candidates for HGT. (At the same time, our simulation study shows how frequently false positives arise, and casts a shadow on the reliability of phylogenetic reconstruction in general.)

Simulations

We performed simulations of genome evolution using EvolSimulator (<http://bioinformatics.org.au/evolsim/>). Seven hundred genes were simulated for 10,000 generations in a dynamic population of genomes, with speciation balancing extinction (at a nominal rate of 0.015 events per generation) to maintain ~50 extant lineages at any given time following a brief initial phase of population growth. Disabling paralogous duplication as well as gene loss, each genome maintained exactly 700 genes whose orthology could thus be perfectly tracked. Parameters affecting se-

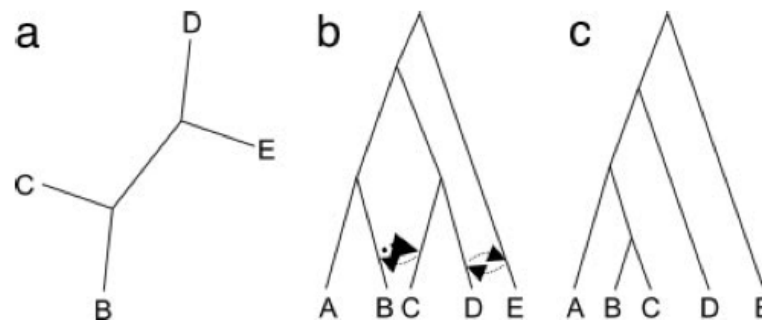
Highways of gene sharing in prokaryotes

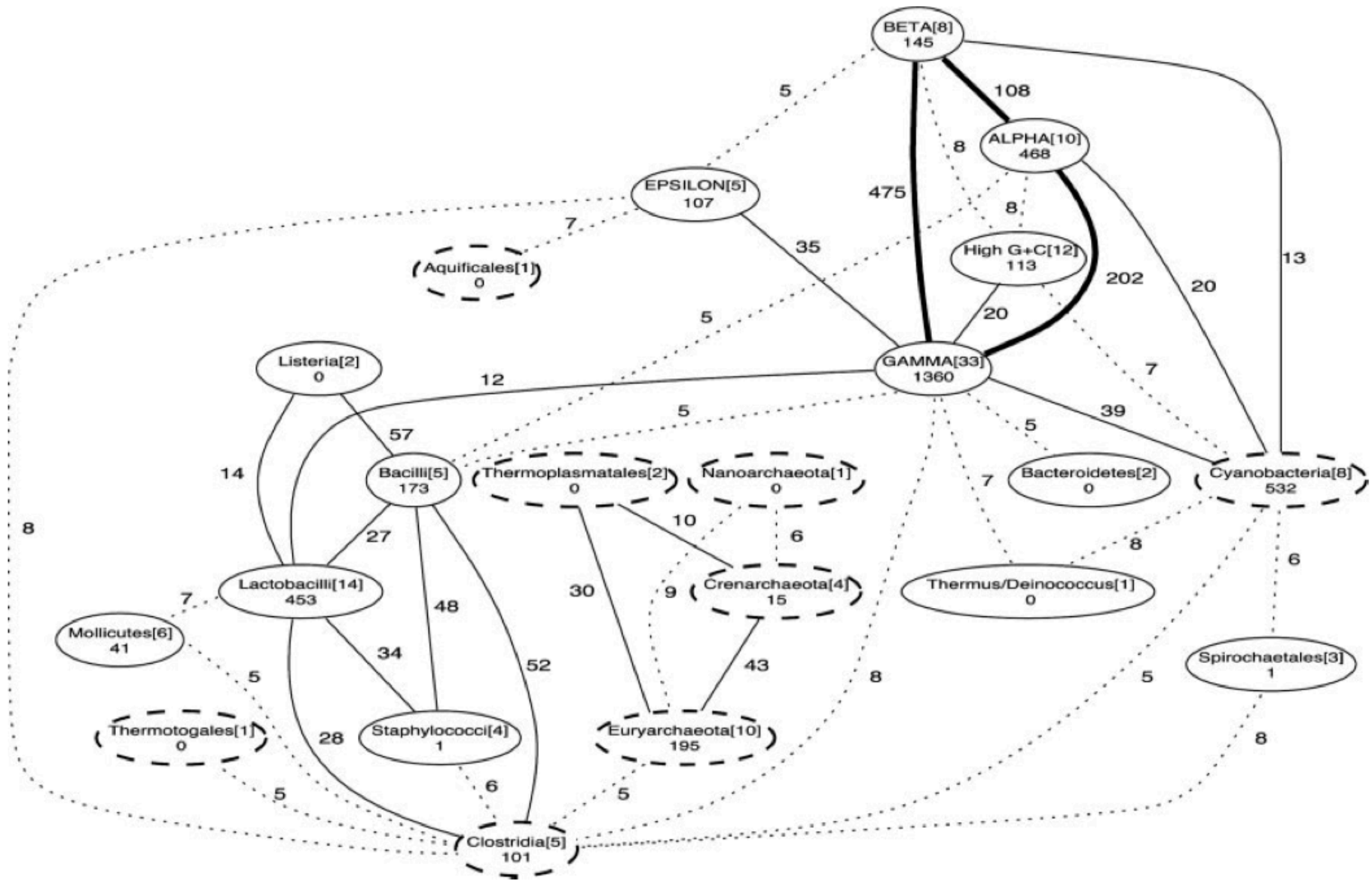
Robert G. Beiko, Timothy J. Harlow, and Mark A. Ragan*

- Considered 144 prokaryotes
- Built Bayesian trees for 22,437 families
- Combined them into a supertree
- Incongruence between bipartitions with high posterior probability and supertree suggests LGT:

support a regime of vertical inheritance at that node. Protein trees that are strongly incongruent with a supertree node are discordant, and provide *prima facie* evidence of LGT. Of the

- Heuristic to find SPR-moves to explain incongruences





Beiko et al. Highways of gene sharing in prokaryotes. PNAS (2005)

Explicit, Model-Based

- Usually based on subtree pruning and regraft (SPR)
- If only interested in the “pruning” -> **Maximum Agreement Forest Problem** (smallest number of edges to cut in 2 trees to yield two identical forests of rooted subtrees)
- **SPR distance between two unrooted binary trees is NP-hard**
 - practical algorithms are approximations

Simultaneous Identification of Duplications and Lateral Transfers

[Extended Abstract]

Mike Hallett
McGill Centre for
Bioinformatics
3775 University Ave.
Montréal, Québec, Canada
hallett@mcb.mcgill.ca

Jens Lagergren
Stockholm Bioinformatics
Centre
Dept. of Numerical Analysis
and Computer Science
Stockholm, Sweden
jensl@nada.kth.se

Integrating Sequence and Topology for Efficient and Accurate Detection of Horizontal Gene Transfer

Cuong Than, Guohua Jin, and Luay Nakhleh*

Inferring and Validating Horizontal Gene Transfer Events Using Bipartition Dissimilarity

ALIX BOC¹, HERVÉ PHILIPPE², AND VLADIMIR MAKARENKO^{1,*}

Detecting Highways of Horizontal Gene Transfer

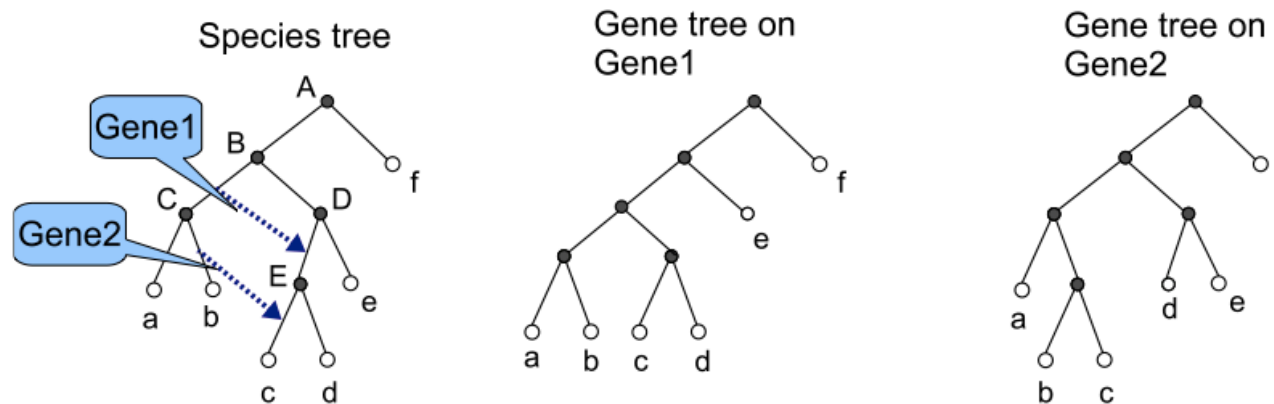
Mukul S. Bansal¹, J. Peter Gogarten², and Ron Shamir¹

Detecting lateral gene transfers by statistical reconciliation of phylogenetic forests

Sophie S Abby, Eric Tannier, Manolo Gouy and Vincent Daubin*

Detecting Highways of Horizontal Gene Transfer

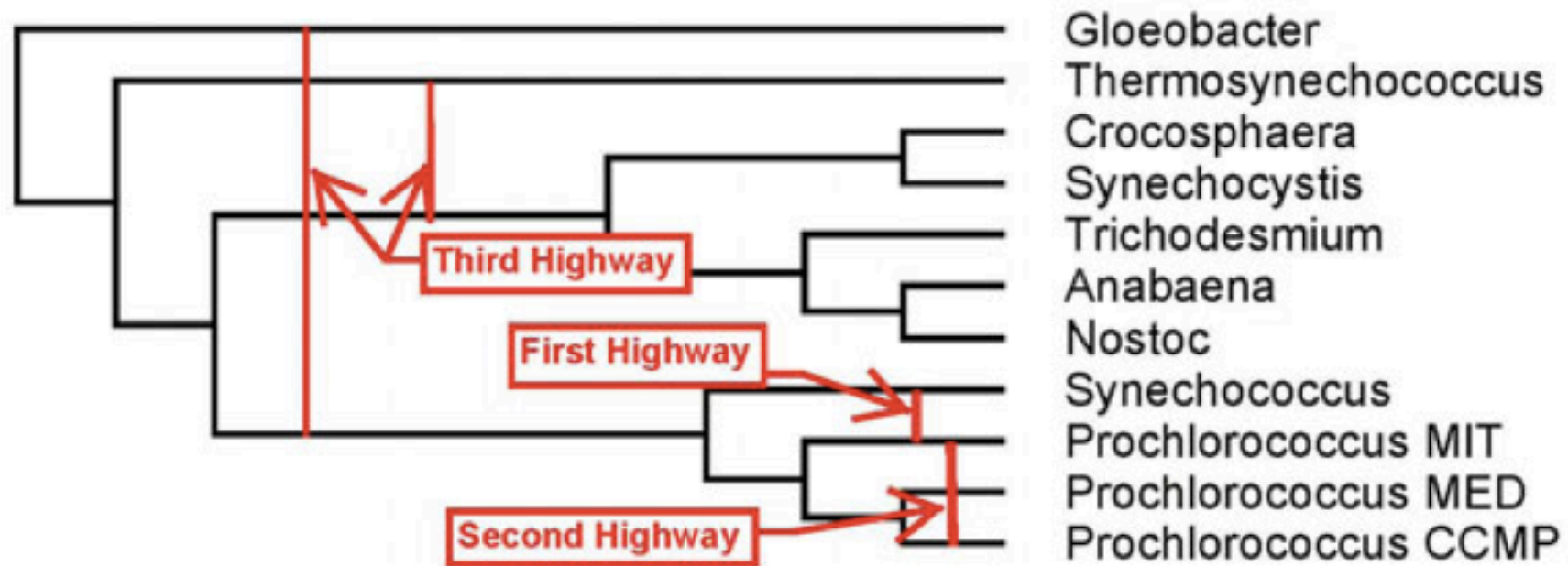
Mukul S. Bansal¹, J. Peter Gogarten², and Ron Shamir¹ (2011)



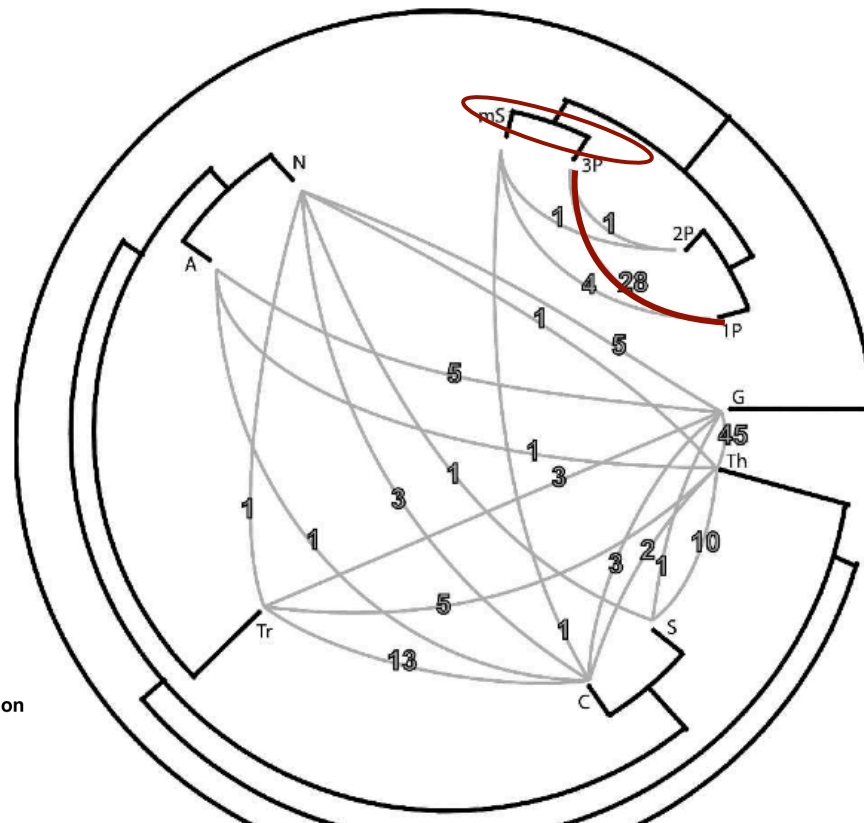
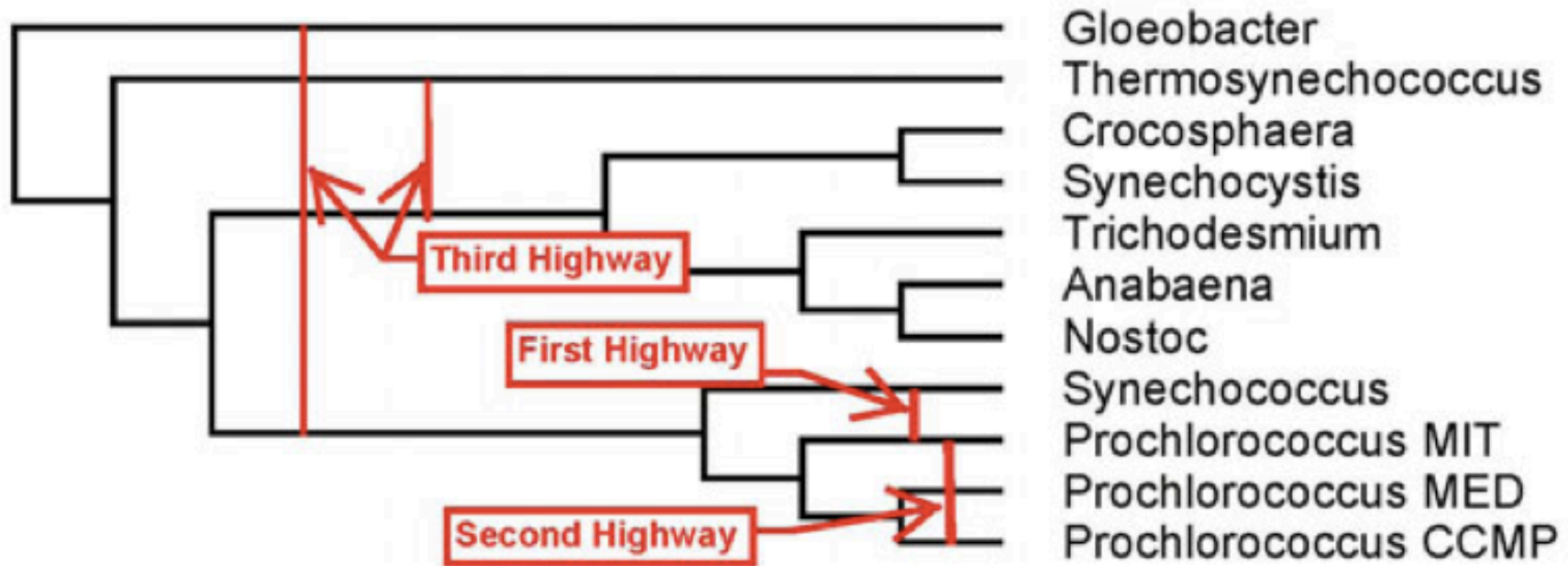
- Step 1:** Decompose each input gene tree T into its constituent set of $\binom{|Le(T)|}{4}$ quartet trees, and combine the quartet trees from the different gene trees into a single weighted set, Φ , of quartet trees.
- Step 2:** Remove from Φ all those quartet trees that are consistent with S .
- Step 3:** Compute the HGT score of each edge in $H(S)$. This HGT score for an edge is computed based on Φ , and is explained in detail below.
- Step 4:** Select the highest scoring horizontal edge as a highway.
- Step 5:** Remove from Φ all those quartet trees that are explained by the proposed highway, and go to Step 3 to start the next iteration.

Detecting Highways of Horizontal Gene Transfer

Mukul S. Bansal¹, J. Peter Gogarten², and Ron Shamir¹



(data from Zhaxybayeva et al. 2006)



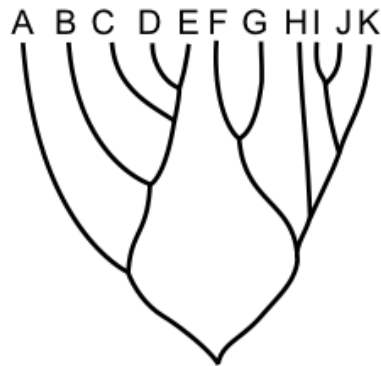
Phylogenetic analyses of cyanobacterial genomes: Quantification of horizontal gene transfer events

Olga Zhaxybayeva, J. Peter Gogarten, Robert L. Charlebois, et al.

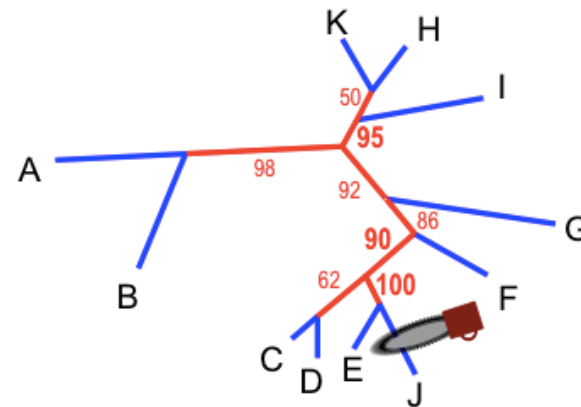
Detecting lateral gene transfers by statistical reconciliation of phylogenetic forests

Sophie S Abby, Eric Tannier, Manolo Gouy and Vincent Daubin* (2010)

Species tree S



Gene tree G



Prunier 1st round

Agree(TS, TG) => NO

Sort candidate transfers (common edges):

- 1) Rm{T(J)} discards: **100** ; 95 ; 90
- 2) Rm{T(E)} discards: 100
- 3) Rm{T(A)} or {T(B)} discards: 98
- 4) Rm{T(C)} or {T(D)} discards: 62
- 5) Rm{T(K)} or {T(H)} discards: 50

Agree{ $TS(J), TG(J)$ } => YES

Prune the highest scoring candidate: T(J)

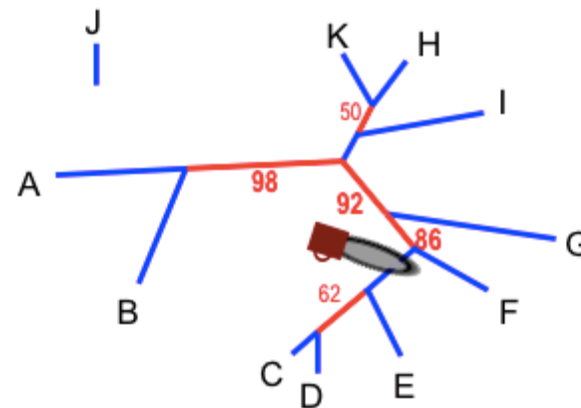
Detecting lateral gene transfers by statistical reconciliation of phylogenetic forests

Sophie S Abby, Eric Tannier, Manolo Gouy and Vincent Daubin*

Species tree S



Gene tree G



Prunier 2nd round

Agree(TS, TG) \Rightarrow NO

Sort candidate transfers (common edges):

- 1) $Rm\{T(CDE)\}$ discards: 98 ; 92 ; 86
- 2) $Rm\{T(A)\}$ or $\{T(B)\}$ discards: 98
- 3) $Rm\{T(C)\}$ or $\{T(D)\}$ or $\{T(E)\}$ discards: 62
- 4) $Rm\{T(K)\}$ or $\{T(H)\}$ or $\{T(I)\}$ discards: 50

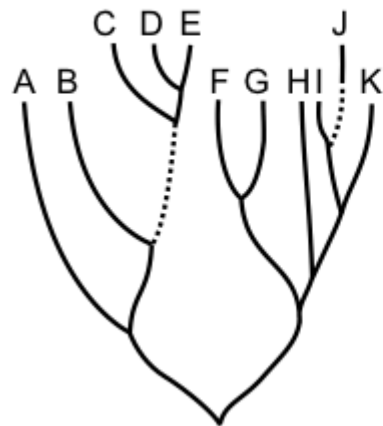
Agree($TS(CDE), TG(CDE)$) \Rightarrow YES

Prune the highest scoring candidate: $T(CDE)$

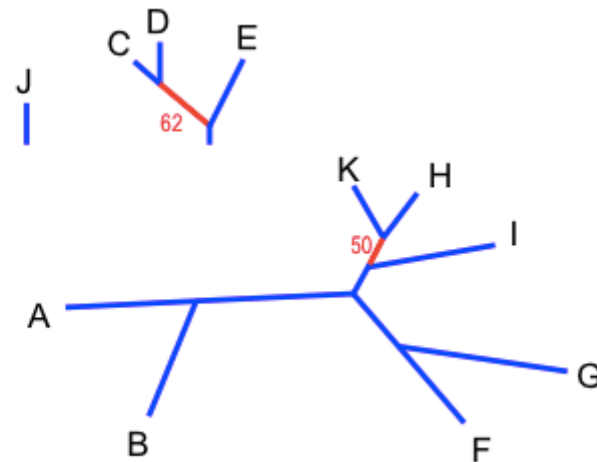
Detecting lateral gene transfers by statistical reconciliation of phylogenetic forests

Sophie S Abby, Eric Tannier, Manolo Gouy and Vincent Daubin*

Species tree S



Gene tree G



Prunier 3rd round

Agree(TS, TG) => **YES**

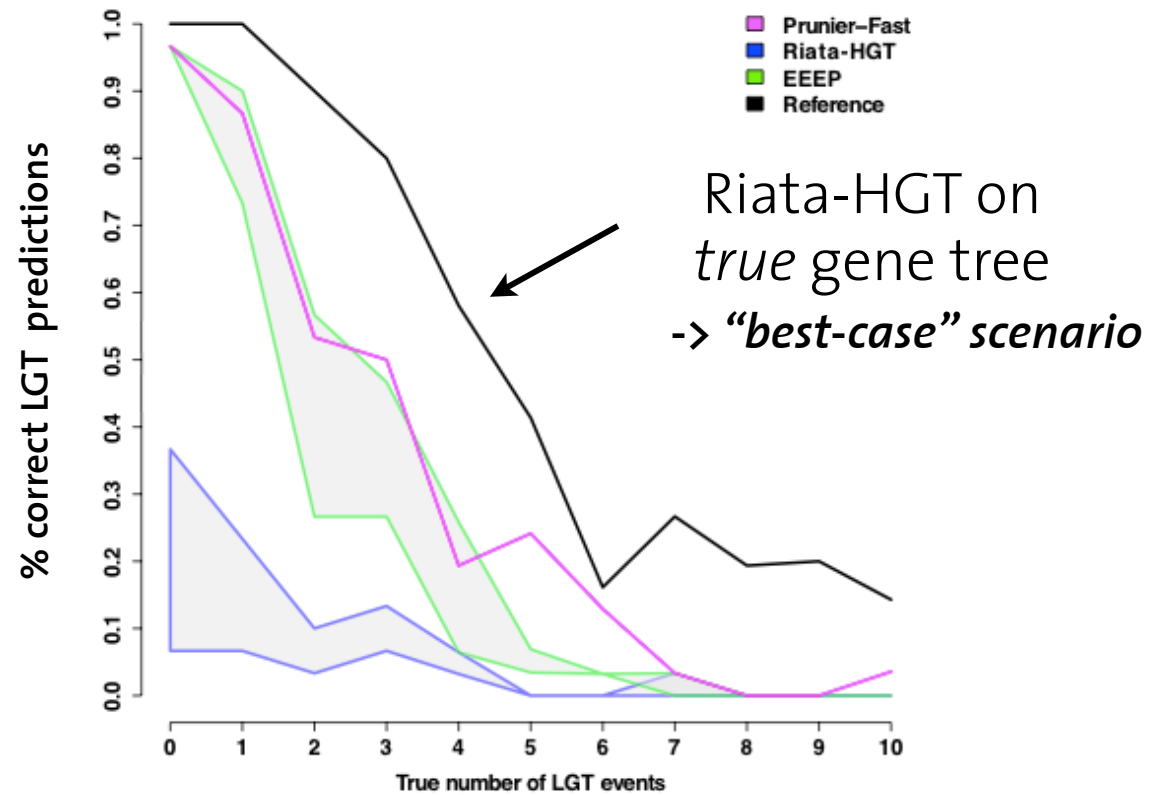
Output: a MSAF decomposition

A reconciled backbone tree:
 $T(ABFGHIK)$

Two reconciled LGT subtrees:
 $T(J)$ and $T(CDE)$

Detecting lateral gene transfers by statistical reconciliation of phylogenetic forests

Sophie S Abby, Eric Tannier, Manolo Gouy and Vincent Daubin*



Limitations of Explicit Phylogenetic Methods

- Factors other than LGT can lead to discordance.
- Cannot identify LGT between sister taxa.
- Relatively computationally expensive.
- *Discordance*: does not characterize the LGT: number of events, donor/recipient species, etc.
- *Model-based*: heuristic-based due to computational complexity

Implicit Phylogenetic Methods

Lawrence & Hartl 1992

Clarke et al. 2002

Novichkov et al. 2004

Choi & Kim 2007

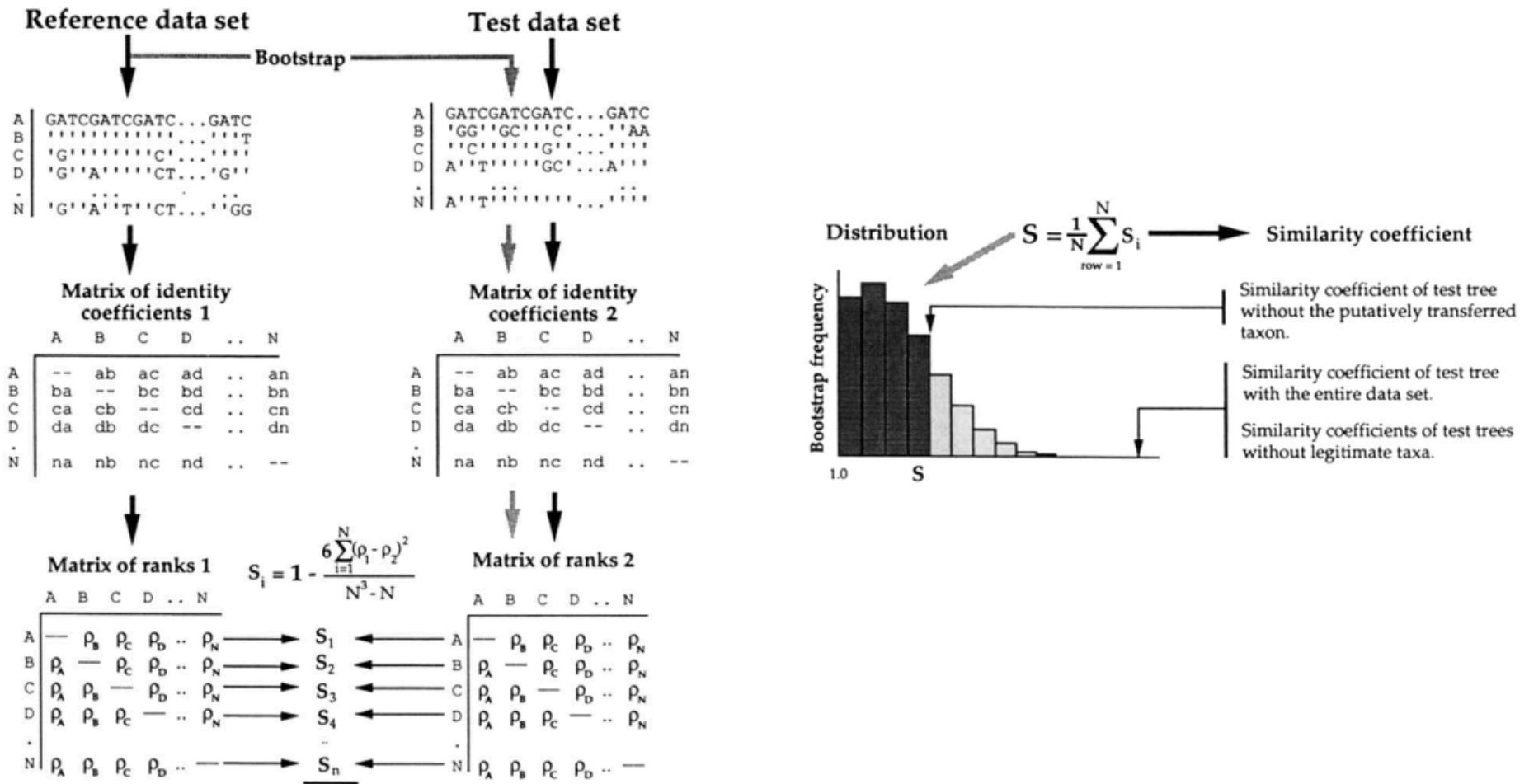
Podell & Gaasterland 2007

Dessimoz et al. 2008

Kanhere & Vingron 2009

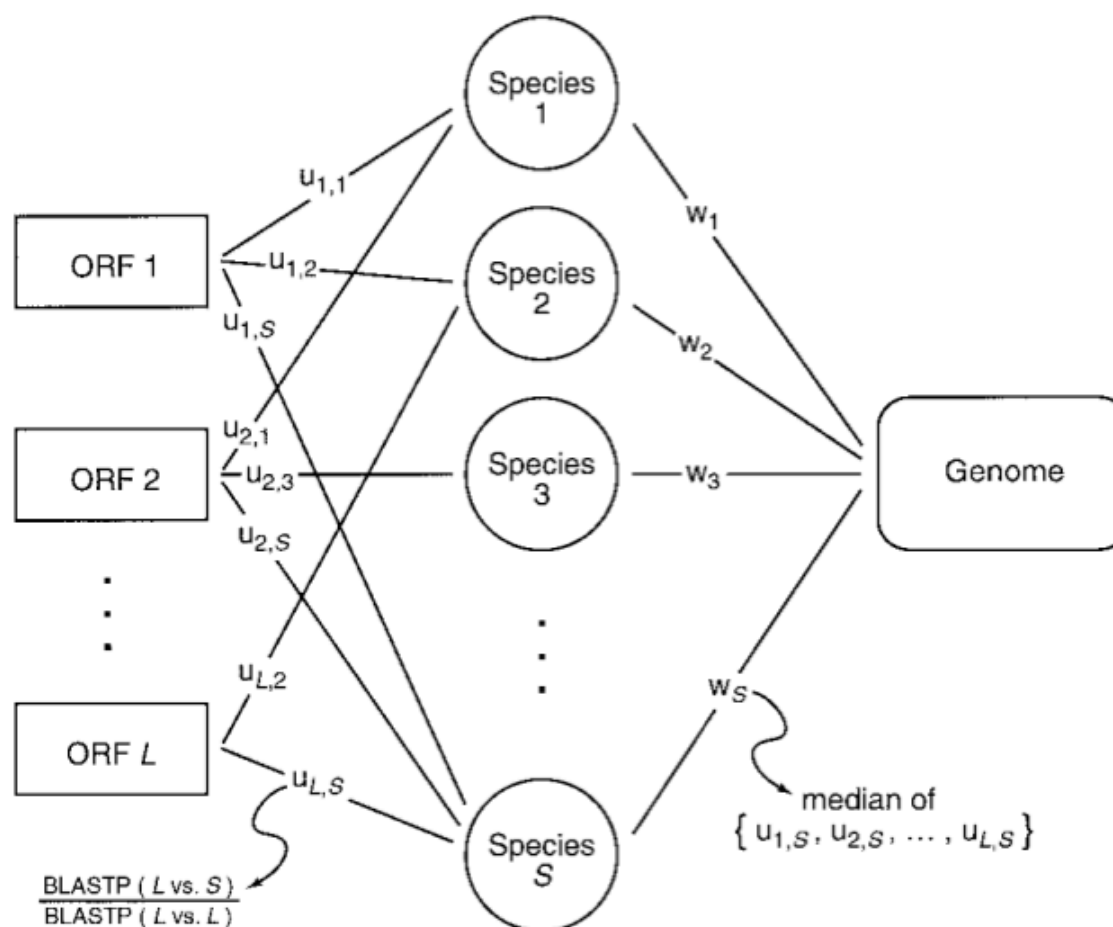
Inference of Horizontal Genetic Transfer From Molecular Data: An Approach Using the Bootstrap

Jeffrey G. Lawrence¹ and Daniel L. Hartl



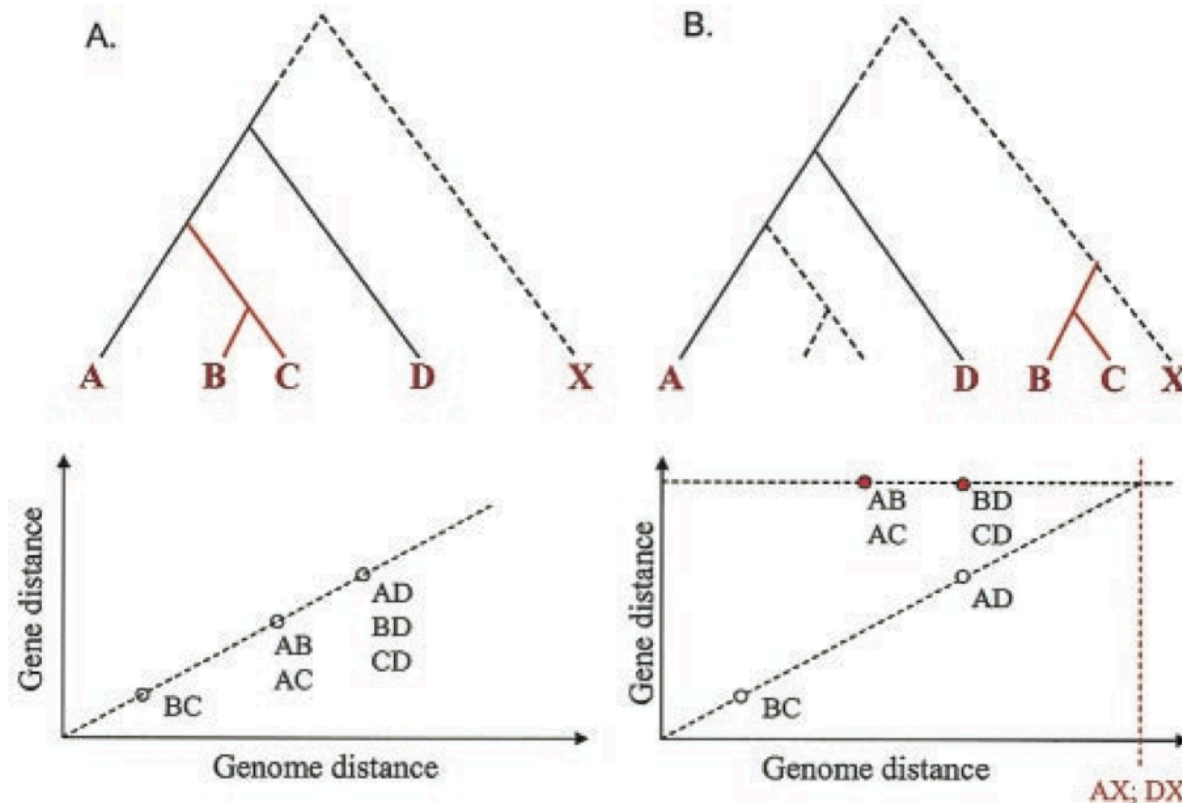
Inferring Genome Trees by Using a Filter To Eliminate Phylogenetically Discordant Sequences and a Distance Matrix Based on Mean Normalized BLASTP Scores

G. D. Paul Clarke,¹ Robert G. Beiko,² Mark A. Ragan,^{3,4*} and Robert L. Charlebois^{1,2,4}



Genome-Wide Molecular Clock and Horizontal Gene Transfer in Bacterial Evolution

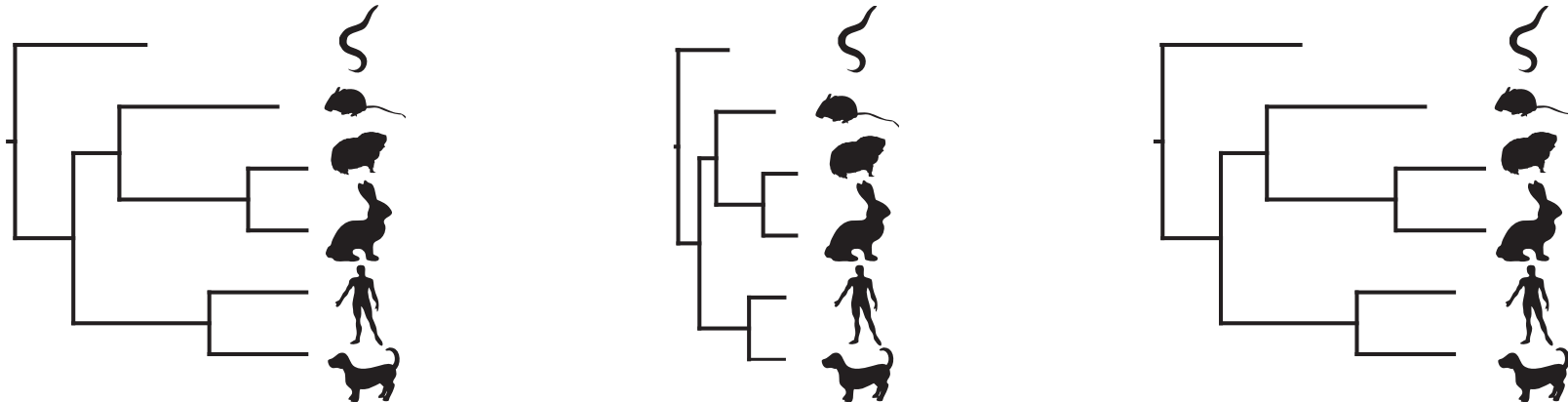
Pavel S. Novichkov,¹ Marina V. Omelchenko,^{2,3} Mikhail S. Gelfand,^{4,5} Andrei A. Mironov,¹
Yuri I. Wolf,³ and Eugene V. Koonin^{3*}



DLIGHT – Lateral Gene Transfer Detection Using Pairwise Evolutionary Distances in a Statistical Framework

Christophe Dessimoz*, Daniel Margadant, and Gaston H. Gonnet

- gene tree = $\text{rate}_f \times \text{species tree}$



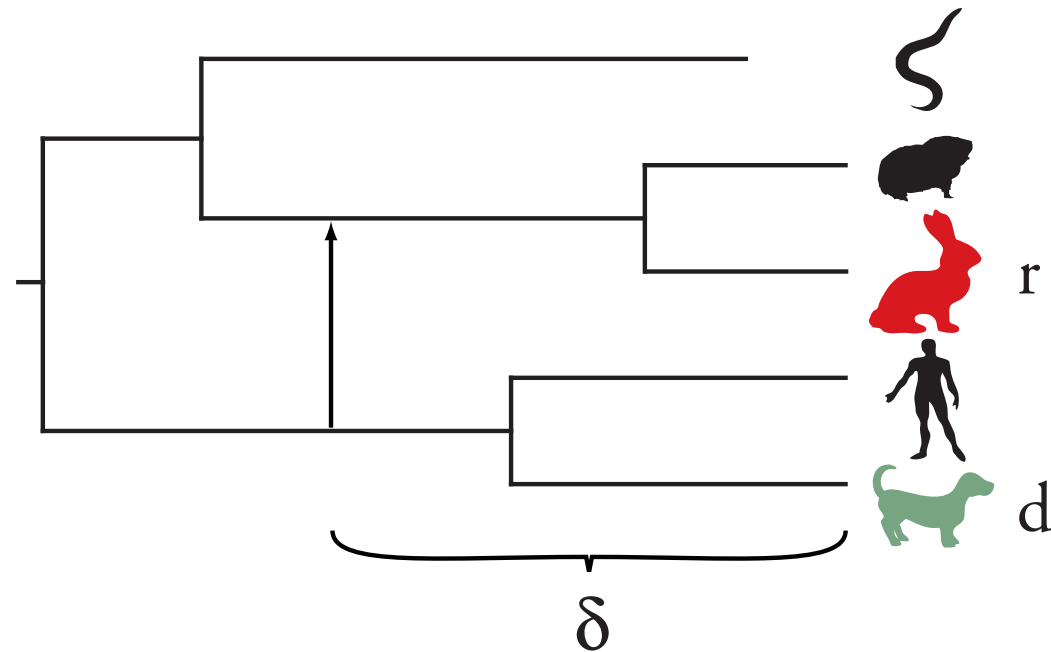
Algorithm

for all orthologous families f do

for all pairs d, r with a seq. in f do

if $2 \ln \frac{l(f, d, r, \delta_{ML})}{l(f, d, r, \delta = \infty)} > \chi^2(\alpha, 1)$ then

the triplet (f, d, r) is a LGT transfer



Likelihood computation: the idea

- Observed distances

$$\mathbf{x} = (x_1, x_2, \dots, x_n)^T \sim \mathcal{N}_n(\mathbf{T}, \Sigma_x)$$

Likelihood computation: the idea

- Observed distances

$$\mathbf{x} = (x_1, x_2, \dots, x_n)^T \sim \mathcal{N}_n(\mathbf{T}, \Sigma_x)$$

- Modeled distances

$$\mathbf{y} = (y_1, y_2, \dots, y_n)^T \sim \mathcal{N}_n(\mathbf{T}, \Sigma_y)$$

where $y_i = g(\text{seqs}_i, d, r, \delta)$

Likelihood computation: the idea

- Observed distances

$$\mathbf{x} = (x_1, x_2, \dots, x_n)^T \sim \mathcal{N}_n(\mathbf{T}, \Sigma_x)$$

- Modeled distances

$$\mathbf{y} = (y_1, y_2, \dots, y_n)^T \sim \mathcal{N}_n(\mathbf{T}, \Sigma_y)$$

$$\text{where } y_i = g(\text{seqs}_i, d, r, \delta)$$

- Thus, assuming that \mathbf{x}, \mathbf{y} are independent,

$$l(f, d, r, \delta) = \text{pdf}(\mathbf{x} - \mathbf{y}) = \frac{\exp(-\frac{1}{2}(\mathbf{x} - \mathbf{y})^T (\Sigma_x + \Sigma_y)^{-1} (\mathbf{x} - \mathbf{y}))}{\sqrt{(2\pi)^n |\Sigma_x + \Sigma_y|}}$$

Limitations of Implicit Phylogenetic Methods

- Cannot identify LGT between sister taxa
- *Most methods* do not characterize the LGTs identified: number of events, donor/recipient species, time of events, etc.

Outlook

- Traditionally parametric vs phylogenetic, but latter is gaining momentum because profits more from more genomes
- Somewhat fragmented community:
 - Microbiologists: parametric methods
 - Theoreticians: model-based phylo.
 - Pragmatic bioinform: discordance
- I see potential for (1) comparing methods; (2) providing large-scale results into DB