

MCB 372 #13:
Selection, Data Partitioning
Gene Transfer

J. Peter Gogarten

University of Connecticut
 Dept. of Molecular and Cell Biology

Collaborators:
 Olga Zhaxybayeva (Dalhousie)
 Jinning Huang (ECU)
 Tim Harlow (UConn)
 Pascal Lapiere (UConn)
 Greg Fournier (UConn)



Edward Munch, The Dance of Life (1900)


Funded through the NASA Exobiology and AISR Programs, and NSF Microbial Genetics

Hy-Phy - Hypothesis Testing using Phylogenies.

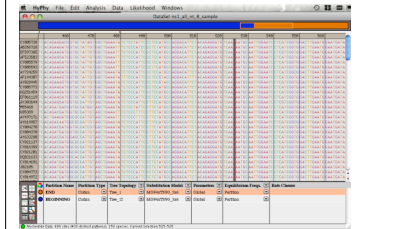
Using Batchfiles or GUI

Information at <http://www.hyphy.org/>

Selected analyses also can be performed online at <http://www.datamonkey.org/>

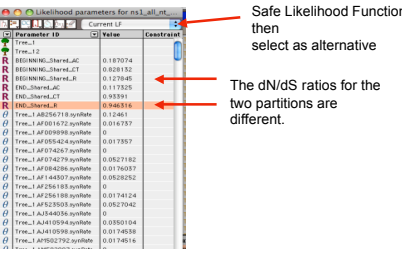


Example testing for dN/dS in two partitions of the data -- John's dataset



Set up two partitions, define model for each, optimize likelihood

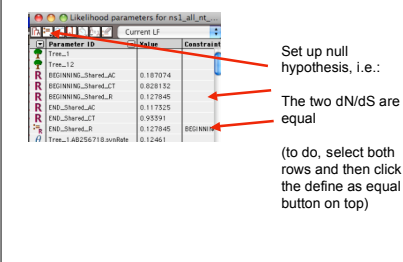
Example testing for dN/dS in two partitions of the data -- John's dataset



Safe Likelihood Function then select as alternative

The dN/dS ratios for the two partitions are different.

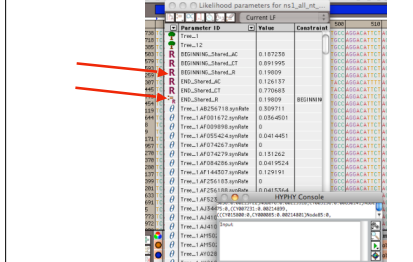
Example testing for dN/dS in two partitions of the data -- John's dataset



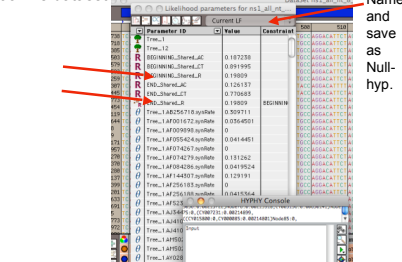
Set up null hypothesis, i.e.: The two dN/dS are equal

(to do, select both rows and then click the define as equal button on top)

Example testing for dN/dS in two partitions of the data -- John's dataset

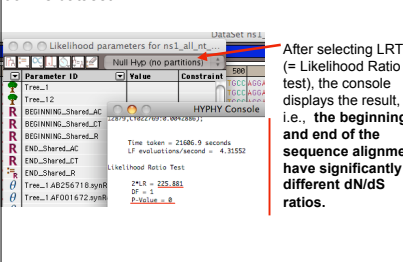


Example testing for dN/dS in two partitions of the data -- John's dataset



Name and save as Null-hyp.

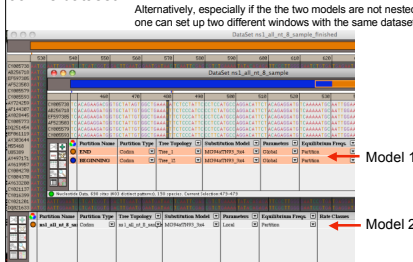
Example testing for dN/dS in two partitions of the data -- John's dataset



After selecting LRT (= Likelihood Ratio Test), the console displays the result, i.e., the beginning and end of the sequence alignment have significantly different dN/dS ratios.

Example testing for dN/dS in two partitions of the data -- John's dataset

Alternatively, especially if the two models are not nested, one can set up two different windows with the same dataset:

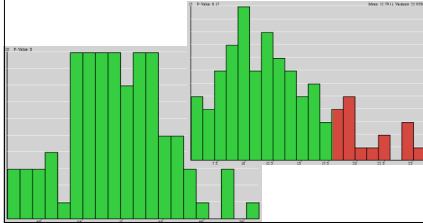


Model 1

Model 2

Example testing for dN/dS in two partitions of the data -- John's dataset

Simulation under model 1, evaluation under model 2, calculate LR
Compare real LR to distribution from simulated LR values. The result might look something like this or this



HGT detection

- **Phylogenetic Incongruence** (conflict between gene and species tree)
- **Phyletic Patterns** (disjunct/spotty distribution)
- **Surrogate Methods** (compositional analyses, violation of clock assumption)

Surrogate Methods - compositional analyses,

Transferred genes often have a different composition compared to the host genome. Especially dinucleotide frequencies provide a useful measure.

Reason A) The transferred gene retains for some time the composition of the donor. (Complete amelioration takes about 100 million years)

<http://www.ncbi.nlm.nih.gov/pubmed/90890782>

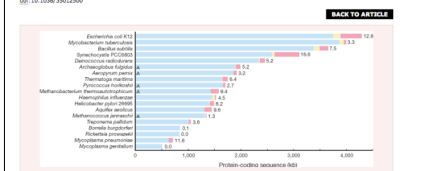
Reason B) The composition reflects the composition of the mobilome, which has a much higher AT content (mutational bias) compared to the genome. (Transferred genes never are AT rich)

<http://www.ncbi.nlm.nih.gov/pubmed/151731102>

Surrogate Methods - compositional analyses,

FIGURE 2. Distribution of horizontally acquired (foreign) DNA in sequenced bacterial genomes.

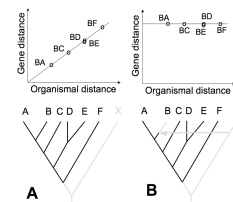
From the following article:
Lateral gene transfer and the nature of bacterial innovation
Howard Ochman, Jeffrey D. Lachance and Eduardo A. Groisman
Nature 405, 299-304 (4 May 2003)
DOI: 10.1038/35912500



Lengths of bars denote the amount of protein-coding DNA. For each bar, the native DNA is blue; foreign DNA identifies its mobile elements, including transposons and bacteriophages, is yellow, and other foreign DNA is red. The percentage of foreign DNA is noted to the right of each bar. 'A' denotes an Archaeal genome.

Surrogate Methods

- clocks



Use of an approximate molecular clock to detect horizontally transferred genes. For each gene, the distance between the gene and its orthologs from closely related genomes is calculated and plotted against the evolutionary distance separating the organisms. The latter can be approximated by ribosomal RNAs or by a genome average. If the gene was inherited vertically, and if the substitution rate remained approximately constant, then the points will fall on a straight line through the origin with a slope depending on the substitution rate of the individual gene (A). If the gene was acquired from outside the organisms considered in the analysis (organism X), then all gene distances will be approximately the same and independent of the distance between the organisms (B). If the transfer occurred to a deeper branch in the tree, part of the points will fall on the diagonal, and part on a parallel line to the abscissa. Modified from Novichkov, P. S., M. V. Ormelchenko, M. S. Gelfand, A. A. Mironov, Y. J. Wolf and E. V. Koonin (2004). Genome-wide molecular clock and horizontal gene transfer in bacterial evolution. *J. Bacteriol.* 186(19): 6575-6585.

Phyletic Methods

- taxonomic distribution of blast hits
- taxonomic position of best blast hit

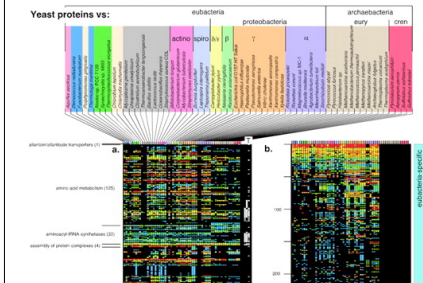
Any non-taxonomic distribution of gene presence and absence can be explained either by

Gene transfer,

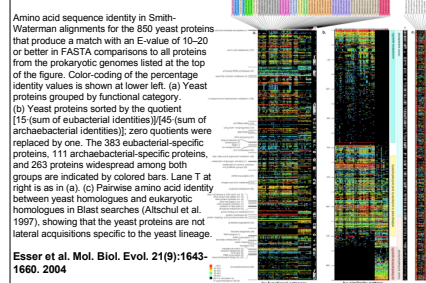
or by

Gene loss. Under the assumption of gene loss any gene present in at least one archaeon and one bacterium would have to be assumed present in the ancestral "Garden of Eden" genome. (Doolittle, W. F., Y. Boucher, C. L. Nesbo, C. J. Douady, J. O. Andersson and A. J. Roger (2003). How big is the iceberg of which organellar genes in nuclear genomes are but the tip? *Philosophical Transactions of the Royal Society B. Biological Sciences* 358(1429): 39-58.)

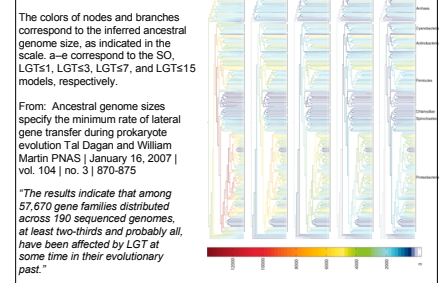
Phyletic Patterns (aside)



Phyletic Patterns (aside)



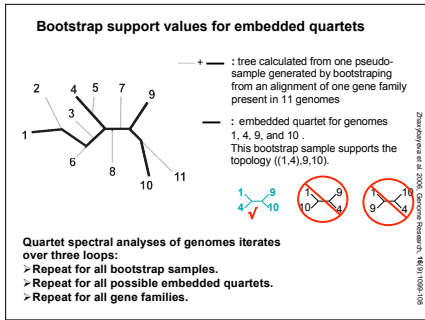
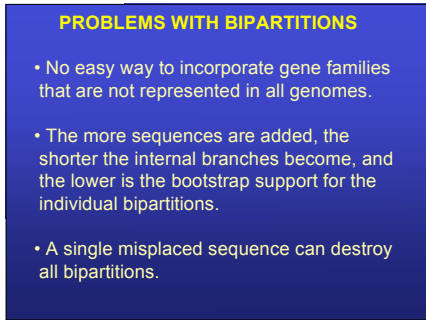
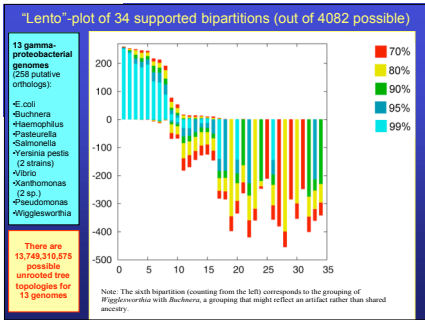
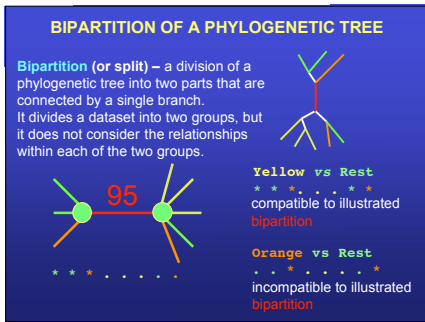
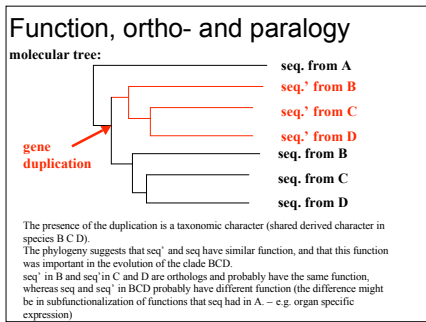
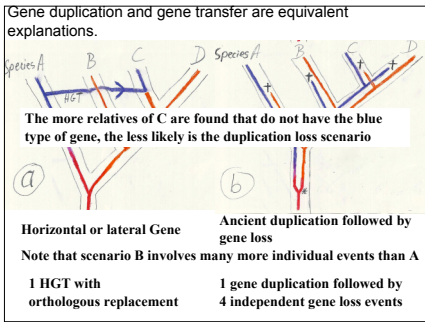
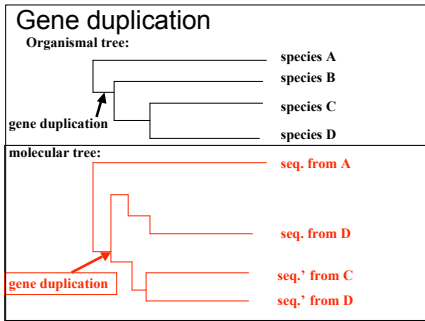
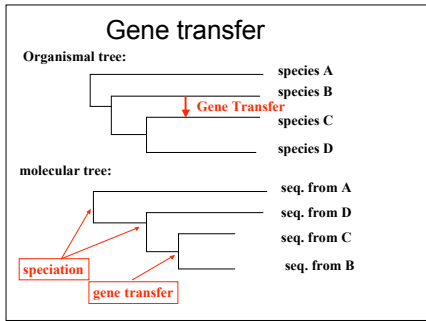
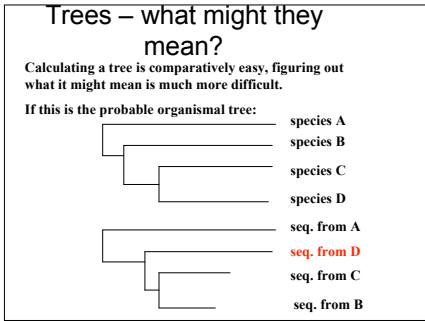
Phyletic Patterns (Garden of Eden Genome)

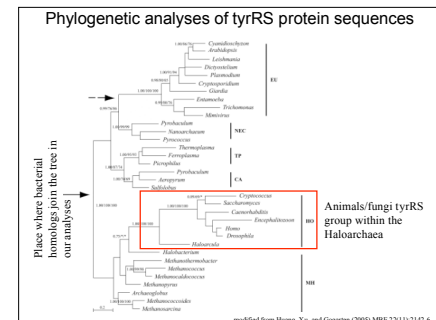
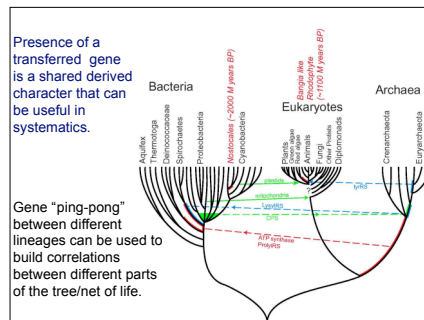
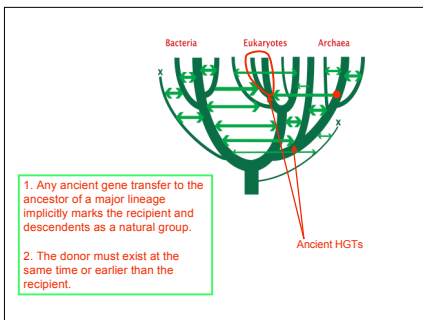
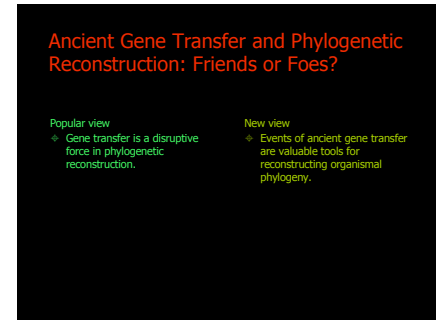
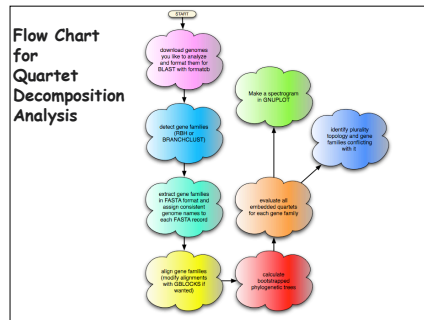
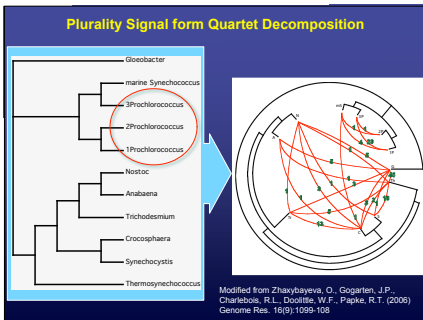
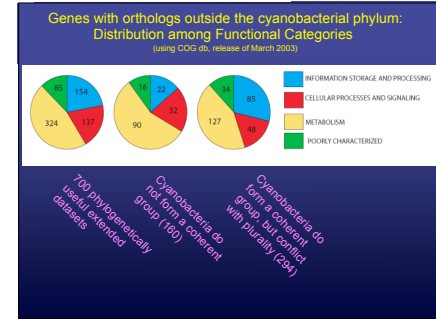
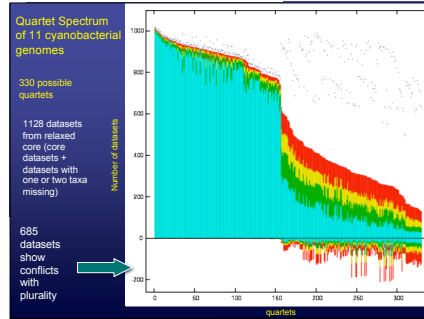
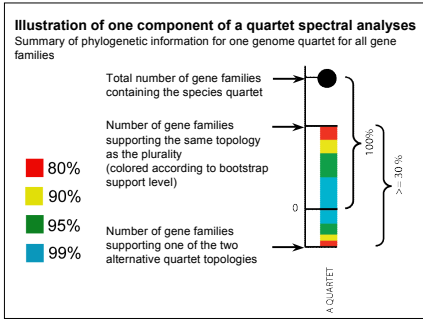


The colors of nodes and branches correspond to the inferred ancestral genome size, as indicated in the scale, a-e correspond to the SO, LGT ≤ 1 , LGT ≤ 3 , LGT ≤ 7 , and LGT ≤ 15 models, respectively.

From: Ancestral genome sizes specify the minimum rate of lateral gene transfer during prokaryote evolution T al Dagan and William Martin PNAS | January 16, 2007 | vol. 104 | no. 3 | 870-875

"The results indicate that among 57,670 gene families distributed across 190 sequenced genomes, at least two-thirds and probably all, have been affected by LGT at some time in their evolutionary past."





Multiple protein sequence alignment for tyrRS.

```

Arabidopsis thaliana |...|
Cyanobacterium |...|
Chlamydomonas reinhardtii |...|
...

```

Signature residues for association of metazoan/fungal/haloarchaeal homologs

Transferred tyrRS supports monophyletic opisthokonts

The same conclusion is reached, if haloarchaeal type tyrRS in opisthokonts is explained by ancient paralogy and differential gene loss.

Monophyly of primary photosynthetic eukaryotes is supported by more than 50 ancient gene transfers from different bacterial phyla to the ancestor of the red algae and green plant lineage.

- E.g., ancient gene transfer of fp-gene (fiorfenicol resistance protein)

Gene from Green plants and Red algae groups with delts proteobacteria

modified from Huang and Carpenter (2004), Trends in Genetics 23, 361-366

Red Algal and Green Plant Genes of Chlamydiae Origin

Gene Name or Gene Product	Protein	Paralog Function
ADH/ATP isomerase	R and G	ATP/ADP transport
Phosphite transporter	R and G	Phosphate transport
Sodium hydrogen antiporter	R and G	Ion transport
Cu-ATPase	R and G	Ion transport
4-ketolactyl-3-methylcrotonyl-3-oxo-1-ol diphosphate synthase (gapD)	R and G	Isoprenoid biosynthesis
4-diphosphocetyl-CoA-methyl-D-erythritol kinase (ispE)	R and G	Isoprenoid biosynthesis
2,4-bisubstituted-erythritol diphosphate synthyltransferase (ispD)	R and G	Isoprenoid biosynthesis
Enoyl-ACP reductase (fabF)	R and G	Fatty acid biosynthesis
Beta-ketoadyl-ACP synthase (fabF)	R and G	Fatty acid biosynthesis
Glucosyl-phosphate acyltransferase	R and G	Phospholipid
Polynucleotide phosphorylase	R and G	RNA degradation
Phosphoglycerate mutase	G	Glycerol
Oligonucleotide P	R	Amino acid biosynthesis
Aspartate transaminase	G	Amino acid metabolism
Malate dehydrogenase	G	Energy conversion
Tyrosyl-tRNA synthetase	R and G	Translation
23S rRNA U1 and U5-methyltransferase	R and G	RNA modification
Invertase	R and G	Sacchar biosynthesis
Hydroxylase protein	R	Lipidases
Sugar phosphatase isomerase	G	Sugar interconversion
CMP-K2O synthetase	G	Cell envelope formation

Consistent phylogenetic signal links Chlamydiae, red algae and green plants.

modified from Huang and Carpenter 2007, Genome Biol 8, R90

Chlamydial-type genes in red algae and plants are often specifically associated with *Protochlamydia* (*Parachlamydia*)

Beta-ketoadyl-ACP synthase (fabF)

4-diphosphocetyl-1,2-C-methyl-D-erythritol kinase (ispE)

The chlamydial genes (plant & red algae) group separate from the cyanobacterial homologs

Polynucleotide Phosphorylase

Tyrosyl-tRNA synthetase

Examining possible Hypotheses

- Plants acquired chlamydial genes via insect feeding activities (Everett et al. 2005).
No. The ancestor of red algae and green plants is much older than insects.
- Chlamydiae acquired plant-like genes via *Ascaritamaeba hosta* (Stephens et al. 1999; Wolf et al. 1999; Orulay et al. 2003).
No. All these genes are of bacterial origin. The direction of gene transfer is from bacteria to eukaryotes.
- Chlamydial and plant sequence similarities reflect an ancestral relationship between chlamydiae and cyanobacteria (Brinkmann et al. 2002; Horn et al. 2004).
No. Genes of chloroplast ancestry should still be more similar to cyanobacterial than to chlamydial sequences. In many instances the cyanobacterial homologs form a clearly distinct, and separate clade.

What do we learn from the data??
(Chlamydial genes in red algal and plant genomes)

Unless a stable physical association existed, it is highly unlikely for any single donor to transfer such a number of genes to a single recipient.

Our Hypothesis: An ancient, unappreciated symbiotic association existed between chlamydiae and the ancestor of red algae and green plants.

Genes from this chlamydial symbiont might have been crucial to establish communication between host and the cyanobacterial cytoplasm.

Hypothesis: Chlamydiae and the primary plastids

A) The Host

White: α -proteobacterial (mitochondrial) symbiont

- Gene transfer to the nucleus
- Transport of nuclear encoded proteins to symbiont
- Direction of symbiotic benefit

Hypothesis: Chlamydiae and the primary plastids

B) The Host invaded by a parasite

Yellow: parasitic chlamydial bacterium

Green: cyanobacterium (as food)

- Gene transfer to the nucleus
- Transport of nuclear encoded proteins to symbiont
- Direction of symbiotic benefit

