

## MCB5472 Computer methods in molecular evolution

Lecture 4/7/2014

### input and output

#### Input and output files

For most of the PHYLIP programs, information comes from a series of input files, and ends up in a series of output files:



The programs interact with the user by presenting a menu. Aside from the user's choices from the menu, they read all other input from files. These files have default names. The program will try to find a file of that name - if it does not, it will ask the user to supply the name of that file. Input data such as DNA sequences comes from a file whose default name is `infile`. If the user supplies a tree, this is in a file whose default name is `intree`. Values of weights for the characters are in `weights`, and the tree plotting program needs some digitized fonts which are supplied in `fontfile` (all these are default names).

### Old Assignment

Write a script that takes all phylip formatted aligned multiple sequence files present in a directory, and performs a bootstrap analyses using maximum parsimony.

Files you might want to use are [A.fa](#), [B.fa](#), [alpha.fa](#), [beta.fa](#) from last week's assignment, and [atp\\_all.phy](#). BUT you first have to **align** them and convert them to **phylip format**! AND you should replace gaps with "?"

(In the end you would be able to answer the question "does the resolution increase if a more related subgroup is analyzed independent from an outgroup?")

- `clustalw2` is one program frequently used to convert formats
- `system("clustalw -infile=$file.fa -convert -output=PHYLIP");`

written and distributed by Joe Felsenstein and collaborators (some of the following is copied from the PHYLIP homepage)

PHYLIP (the *PHY*logeny *IN*ference *PA*ckage) is a package of programs for inferring phylogenies (evolutionary trees).

PHYLIP is the most widely-distributed phylogeny package, and competes with PAUP\* to be the one responsible for the largest number of published trees. PHYLIP has been in distribution since 1980, and has over 15,000 registered users.

Output is written onto special files with names like "outfile" and "outtree". Trees written onto "outtree" are in the [Newick](#) format, an informal standard agreed to in 1986 by authors of a number of major phylogeny packages.

Input is either provided via a file called "infile" or in response to a prompt.

### What's in PHYLIP

Programs in PHYLIP allow to do parsimony, distance matrix, and likelihood methods, including bootstrapping and consensus trees. Data types that can be handled include molecular sequences, gene frequencies, restriction sites and fragments, distance matrices, and discrete characters.

Phylip works well with protein and nucleotide sequences. Many other programs mimic the style of PHYLIP programs. (e.g. TREEPUZZLE, phylml, protml)

Many other packages use PHYLIP programs in their inner workings (e.g., SEAVIEW)

PHYLIP runs under all operating systems

Web interfaces are available

### Programs in PHYLIP are Modular

For example:

SEQBOOT take one set of aligned sequences and writes out a file containing bootstrap samples.

PROTDIST takes a aligned sequences (one or many sets) and calculates distance matrices (one or many)

FITCH (or NEIGHBOR) calculate best fitting or neighbor joining trees from one or many distance matrices

CONSENSE takes many trees and returns a consensus tree

.... modules are available to draw trees as well, but often people use [figtree](#) or [njplot](#).

[The Phylip Manual](#) is an excellent source of information.

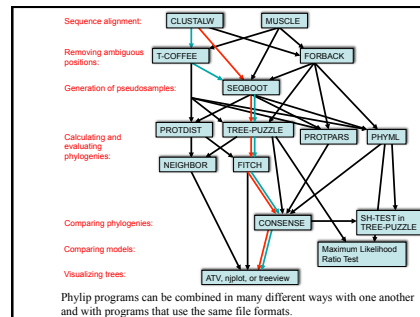
Brief one line descriptions of the programs are [here](#)

The easiest way to run PHYLIP programs is via a command line menu (similar to clustalw). The program is invoked through clicking on an icon, or by typing the program name at the command line.

```
> seqboot
> protpars
> fitch
```

If there is no file called infile the program responds with:

```
[gogarten@carrot gogarten]$ seqboot
seqboot: can't find input file "infile"
Please enter a new file name>
```



Phylip programs can be combined in many different ways with one another and with programs that use the same file formats.

### Example 1 Protpars

example: `seqboot, protpars, consense` on `atp_all.phy`

NOTE the bootstrap majority consensus tree does not necessarily have the same topology as the "best tree" from the original data!

threshold parsimony,  
gap symbols - versus ?  
(in vi you could use :%s/-/?/g to replace all - ?)  
outfile  
outtree compare to distance matrix analysis

create \*.phy files

the easiest (probably) is to run clustalw with the phylip option:  
For example [here](#):

```
#!/usr/bin/perl -w
print "0 This program aligns all multiple sequence files with names *.fa in
# found in its directory using clustalw, and saves them in phylip format.n";
while(defined($file=glob("*.fa"))){
    @parts=split(/./,$file);
    $file=$parts[0];
    system("clustalw -infile=$file -fa -align -output=PHYLIP");
    #if you only want to convert file use:
    system("clustalw -infile=$file -fa -convert -output=PHYLIP");
};
# cleanup:
system("rm *.dnd");
exit;
```

Alternative for entering the commands for the menu:

```
#!/usr/bin/perl -w
system ("cp A.phy infile");
system ("echo -e 'y\n9\n'|seqboot");
exit;

echo returns the string in '\n', i.e., y\n9\n.
The -e options allows the use of \n
The | symbol pipes the output from echo to seqboot
```

## New Assignment

"Given a multiple fasta sequence file", write a script that for each sequence extract the gi number and the species name, and then rewrites the file so that the annotation line starts with the gi number, followed by the species/strain name, followed by a space. (The gi number and the species name should not be separated by or contain any spaces – replace them by \_). This is useful, because many programs will recognize the number and name as handle for the sequence (e.g., clustalw2 and phylml)

Assume that the annotation line follows the NCBI convention and begins with the > followed by the gi number, and ends with the species and strain designation given in []

Example:  
>gi|229240723|ref|ZP\_04365119.1| primary replicative DNA  
helicase; intein [Cellulomonas flavigena DSM 20109]

\*An example multiple sequence file in the unofficial NCBI formatted annotation line is [here](#).

## Bayes' Theorem



Reverend Thomas Bayes  
(1702-1761)

$P(\text{model}|\text{data}, I) = \frac{P(\text{data}|\text{model}, I) \cdot P(\text{model}, I)}{P(\text{data}, I)}$

**Posterior Probability** describes the degree to which we believe a given model accurately describes the situation given the available data and all of our prior information I

**Prior Probability** describes the degree to which we believe the model accurately describes reality based on all of our prior information.

**Likelihood** describes how well the model predicts the data

**Normalizing constant**

## Illustration of a biased random walk

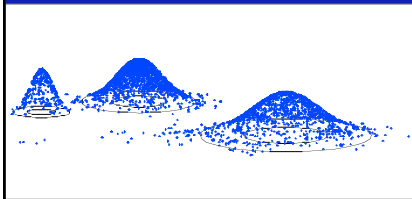


Figure generated using MCRobot program (Paul Lewis, 2001)

## Alternative Approaches to Estimate Posterior Probabilities

Bayesian Posterior Probability Mapping with MrBayes  
(Huelsenbeck and Ronquist, 2001)

**Problem:**

Stimmer's formula  $p_i = \frac{L_i}{L_1 + L_2 + L_3}$  only considers 3 trees (those that maximize the likelihood for the three topologies)

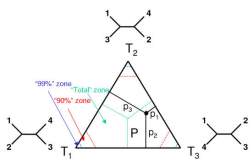
**Solution:**

Exploration of the tree space by sampling trees using a biased random walk (implemented in MrBayes program)

Trees with higher likelihoods will be sampled more often

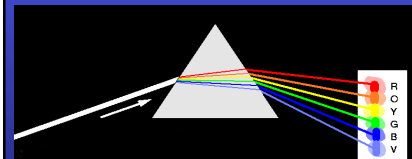
$p_i = \frac{N_i}{N_{\text{total}}}$  where  $N_i$  - number of sampled trees of topology  $i$ ,  $i=1,2,3$   
 $N_{\text{total}}$  - total number of sampled trees (has to be large)

## mi mapping



From: Olga Zhaxybayeva and J Peter Gogarten *BMC Genomics* 2002, 3:4

## Decomposition of Phylogenetic Data



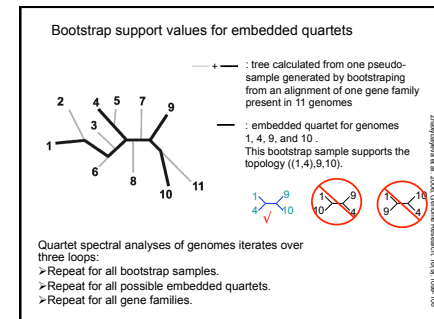
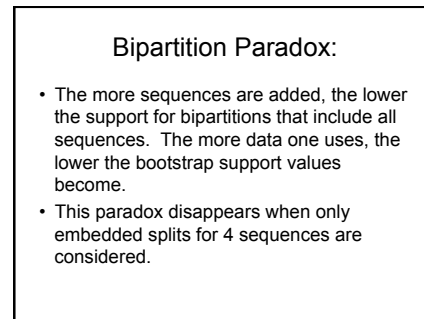
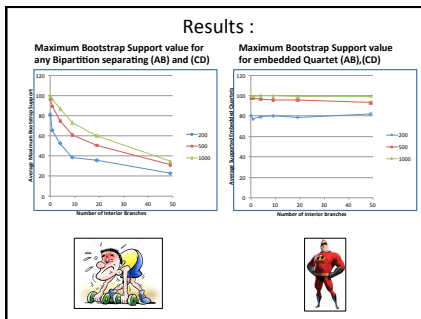
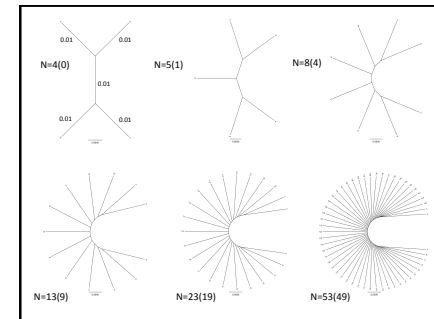
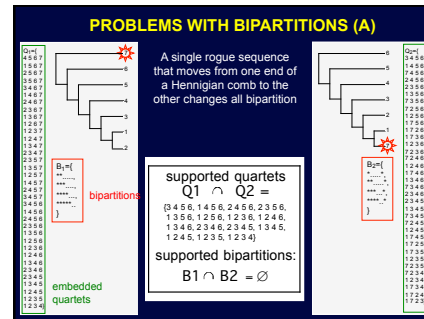
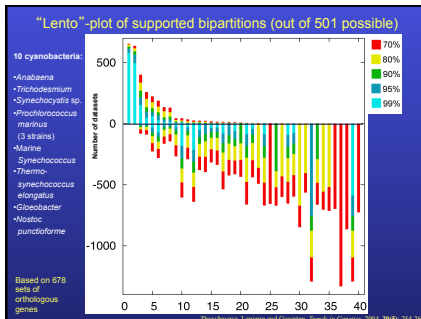
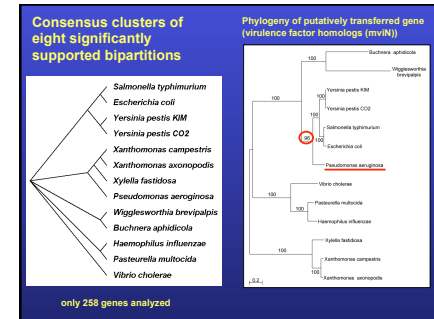
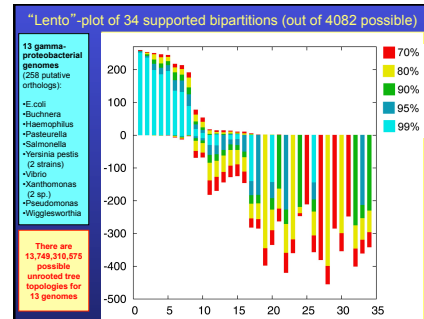
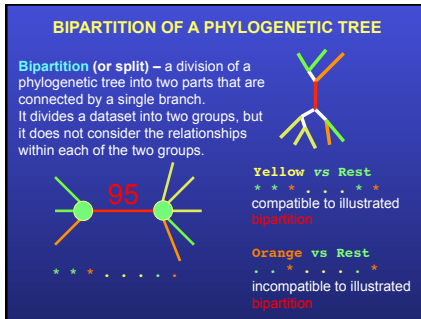
Phylogenetic information present in genomes

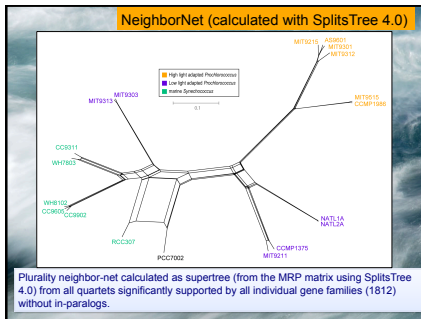
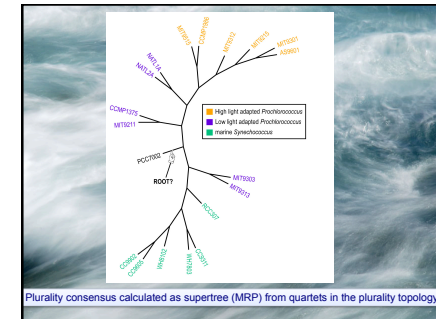
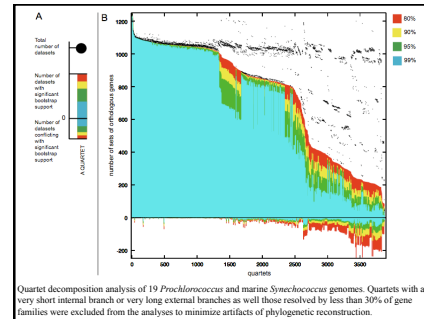
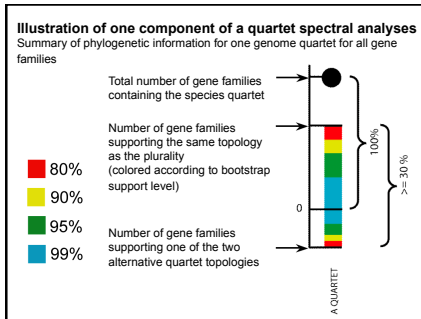
Break information into small quanta of information (bipartitions or embedded quartets)

Analyze spectra to detect transferred genes and plurality consensus.

## TOOLS TO ANALYZE PHYLOGENETIC INFORMATION FROM MULTIPLE GENES IN GENOMES:

**Bipartition Spectra (Lento Plots)**





**Neutral theory:**

The vast majority of observed sequence differences between members of a population are neutral (or close to neutral). These differences can be fixed in the population through random genetic drift. Some mutations are strongly counter selected (this is why there are patterns of conserved residues). Only very seldom is a mutation under positive selection.

The neutral theory **does not** say that all evolution is neutral and everything is only due to genetic drift.

**Nearly Neutral theory:**

Even synonymous mutations do not lead to random composition but to codon bias. Small negative selection might be sufficient to produce the observed codon usage bias.

**How do you define evolution?**

**Richard Goldschmidt 1940**  
hopeful monsters  
Mutationism **HGT/WGD**  
Punctuated Equilibrium  
Few genes / large effect  
Vilified by Mayr; celebrated 1977 Gould & Evo-devo

**Ernst Mayr 1942**  
NeoDarwinian Synthesis  
Natural Selection  
Gradualism  
Many genes/small effect  
Dario - "Fisher right"

**Motoo Kimura 1968**  
Neutral Theory  
Genetic Drift is main force for changing allele frequencies

Slide from Chris Pires

**Duplications and Evolution**

**Susumu Ohno 1970**  
Evolution by gene duplication  
1R and 2R hypothesis  
"Junk DNA" 1972

Ohno postulated that gene duplication plays a major role in evolution

Small scale duplications (SSD)  
Whole genome duplications (WGD)

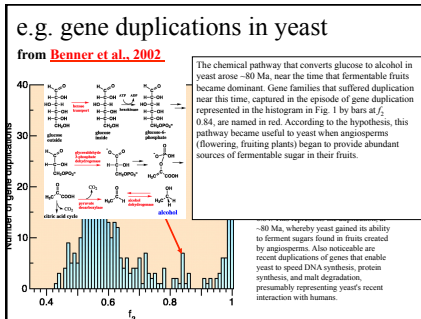
- **Polyploid:** nucleus contains three or more copies of each chromosome
- **Autopolyploid:** formed within a single species  
Diploids **AA** and **A'A'**  $\Rightarrow$  Polyploid **AAA'A'**
- **Allopolyploid:** formed from more than one species  
Diploids **AA** and **BB**  $\Rightarrow$  Polyploid **AA'BB**

Slide from Chris Pires

**What is it good for?**

**Gene duplication** events can provide an outgroup that allows rooting a molecular phylogeny. Most famously this principle was applied in case of the tree of life – the only outgroup available in this case are ancient paralogs (see [http://gogarten.uconn.edu/cvs/Publ\\_Pres.htm](http://gogarten.uconn.edu/cvs/Publ_Pres.htm) for more info). However, the same principle also is applicable to any group of organisms, where a duplication preceded the radiation (**example**). Lineage specific duplications also provide insights into which traits were important during evolution of a lineage.





## the gradualist point of view

Evolution occurs within populations where the fittest organisms have a selective advantage. Over time the advantages genes become fixed in a population and the population gradually changes.

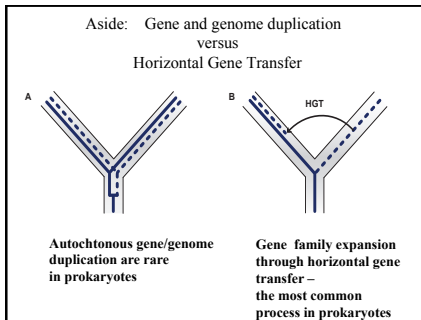
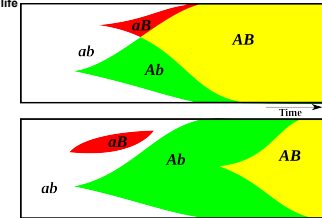
Note: this is not in contradiction to the theory of neutral evolution. (which says what?)

Processes that MIGHT go beyond inheritance with variation and selection?

- Horizontal gene transfer and recombination
- Polyploidization (botany, vertebrate evolution) see [here](#) or [here](#)
- Fusion and cooperation of organisms (Kefir, lichen, also the eukaryotic cell)
- Targeted mutations (?), genetic memory (?) (see [Foster's](#) and [Hall's](#) reviews on directed/adaptive mutations; see [here](#) for a counterpoint)
- Random genetic drift
- [Gratuitous complexity](#)
- Selfish genes (who/what is the subject of evolution??)
- Parasitism, altruism, [Morons](#).
- [Evolutionary capacitors](#)
- [Hopeless monsters](#) (in analogy to Goldschmidt's [hopeful monsters](#))

## Gene Transfer, Sex, and Recombination:

- Inventions do not need to be made sequentially
- Gene transfer, followed by homologous or non-homologous recombination, allows inventions to be shared across the tree of life



## Horizontal Gene Transfer (HGT) and the Acquisition of New Capabilities

- Most important process to adapt microorganisms to new environments. E.g.: Antibiotic and heavy metal resistance, pathways that allow acquisition and breakdown of new substrates.

- Creation of new metabolic pathways.

- HGT not autochthonous gene duplication is the main process of gene family expansion in prokaryotes.

- Also important in the recent evolution of multicellular eukaryotes (HGT between fish species and between grasses).

## Selection acts on the Holobiont (= Host + Symbionts)

- To adapt to new conditions, new symbionts can be acquired, or existing symbionts can acquire new genes through HGT.



## Gene Transfer in Eukaryotes

Published online 7 April 2010 | Nature | doi:10.1038/nature09037

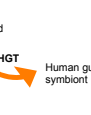
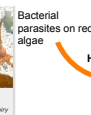
**A genetic gift for sushi eaters**

Seaweed-rich diet leaves its mark on gut microbes.

[Nobu Lohoff](#)

**Transfer of carbohydrate-active enzymes from marine bacteria to Japanese gut microbiota**

Jan-Hendrik Hehemann<sup>1,2,3</sup>, Gaëlle Correia<sup>1,2</sup>, Tristan Barbeyron<sup>1,2</sup>, William Helbert<sup>1,2</sup>, Mirjam Czjzek<sup>1,2</sup> & Gervan Michel<sup>1,2</sup>



Porphyria – also known as nori.

M.D. Guiry

HGT

Human gut symbiont

## Gene Transfer in Eukaryotes – Example 2

### Highlights

- Key genes for  $C_4$  photosynthesis were transmitted between distantly related grasses
- These genes contributed to the adaptation of the primary metabolism
- Their transmission was independent from most of the genome

Curr Biol. 2012 Mar 6;22(5):445-9. Epub 2012 Feb 16.

### Adaptive Evolution of $C(4)$ Photosynthesis through Recurrent Lateral Gene Transfer.

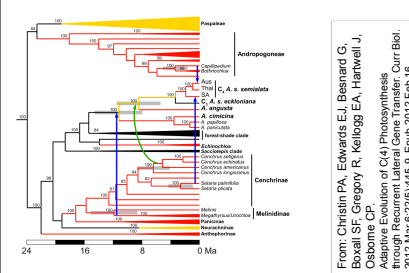
Christin PA, Edwards EJ, Besnard G, Boxall SF, Gregory R, Kellogg EA, Hartwell J, Osborne CP.

### $C_4$ Photosynthesis: Need a Gene? Borrow One!

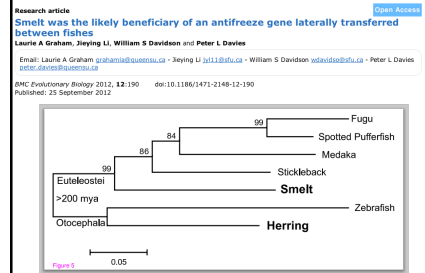
Eric H. Roalson Current Biology Vol 22 No 5 R162

Horizontal gene transfer has been increasingly documented between eukaryotes, but a new study suggests a much larger role for horizontal gene transfer in physiological adaptation through the transfer of photosynthetic pathway genes.

## Gene Transfer in Eukaryotes – Example 2

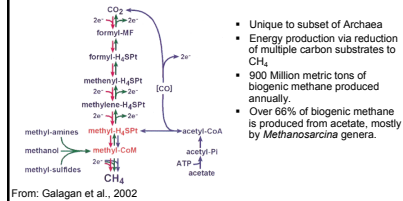


## Gene Transfer in Eukaryotes – Example 3



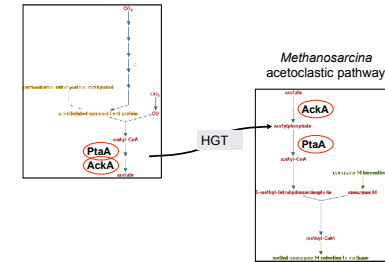
## HGT as a force creating new pathways

### HGT as a force creating new pathways – Example 1 Acetoclastic Methanogenesis



Fourier and Gogarten (2008) Evolution of Acetoclastic Methanogenesis in Methanosarcina via Horizontal Gene Transfer from Cellulolytic Clostridia. *J. Bacteriol.* 190(3): 1124-7

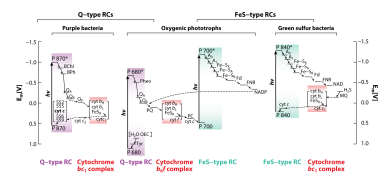
### Clostridia acetogenic pathway



Figures drawn with Metacyc (www.metacyc.org)

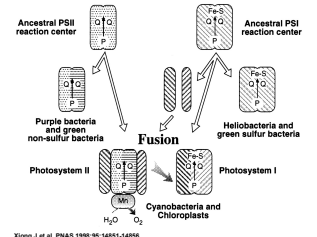
### HGT as a force creating new pathways – Example 2

#### Oxygen producing photosynthesis



Hohmann-Marriott MF, Blankenship RE. 2011. *Annu. Rev. Plant Biol.* 62:515-48

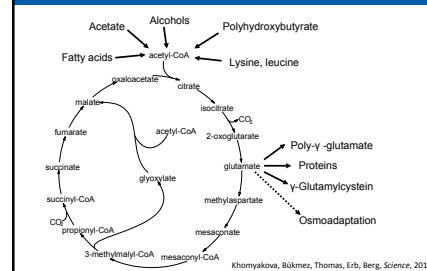
### A heterologous fusion model for the evolution of oxygenic photosynthesis based on phylogenetic analysis.



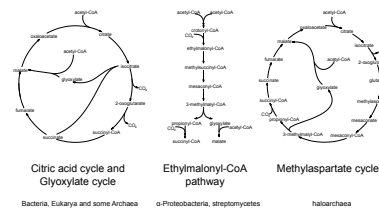
Xiong J et al. PNAS 1998;95:14881-14886

PNAS

### HGT as a force creating new pathways – Example 3 Acetyl-CoA Assimilation: Methylaspartate Cycle



## Comparison of different anaplerotic pathways

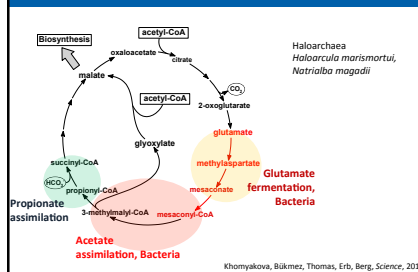


Bacteria, Eukarya and some Archaea

α-Proteobacteria, streptococci

Haloarchaea

### HGT as a force creating new pathways – Example 3 Acetyl-CoA Assimilation: methylaspartate cycle



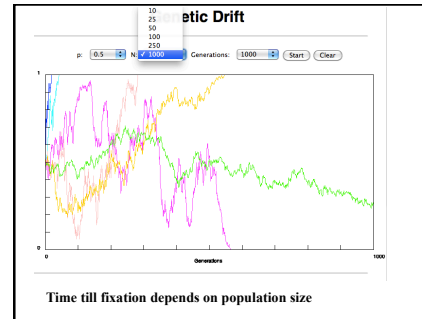
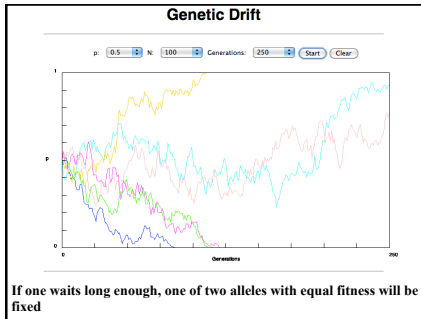
Khomyakova, Bükmez, Thomas, Erb, Berg, Science, 2011

## selection versus drift

The larger the population the longer it takes for an allele to become fixed.

Note: Even though an allele conveys a strong selective advantage of 10%, the allele has a rather large chance to go extinct.

Note#2: Fixation is faster under selection than under drift.



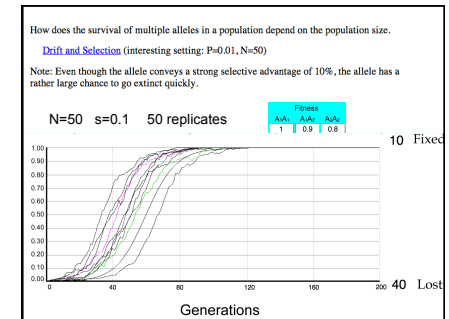
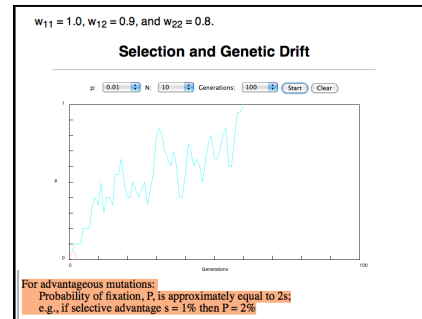
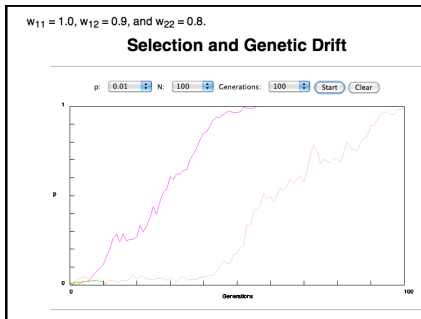
$s=0$

Probability of fixation,  $P$ , is equal to frequency of allele in population.  
 Mutation rate (per gene/per unit of time) =  $u$  ;  
 freq. with which allele is generated in diploid population size  $N = u \cdot 2N$   
 Probability of fixation for each allele =  $1/(2N)$

**Substitution rate** =  
 frequency with which new alleles are generated \* Probability of fixation =  
 $u \cdot 2N \cdot 1/(2N) = u = \text{Mutation rate}$

Therefore:  
 If  $s=0$ , the substitution rate is independent of population size, and equal to the mutation rate !!!! (NOTE: Mutation unequal Substitution! )  
 This is the reason that there is hope that the molecular clock might sometimes work.

Fixation time due to drift alone:  
 $t_{fix} = 4N_e$  generations  
 $(N_e = \text{effective population size; For } n \text{ discrete generations})$   
 $N_e = n / (1/N_1 + 1/N_2 + \dots + 1/N_n)$



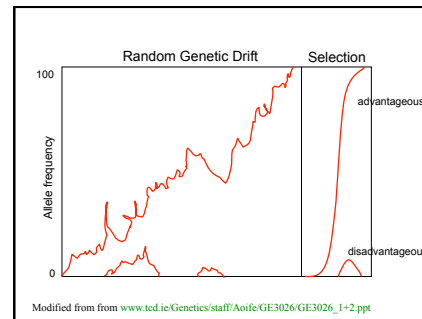
$s>0$

Time till fixation on average:  
 $t_{fix} = (2/s) \ln (2N)$  generations  
 (also true for mutations with negative "s"! discuss among yourselves)

E.g.:  $N=10^6$ ,  
 $s=0$ : average time to fixation:  $4 \cdot 10^6$  generations  
 $s=0.01$ : average time to fixation: 2900 generations

$N=10^4$ ,  
 $s=0$ : average time to fixation: 40,000 generations  
 $s=0.01$ : average time to fixation: 1,900 generations

=> substitution rate of mutation under positive selection is larger than the rate with which neutral mutations are fixed.



**Positive selection ( $s>0$ )**

- A new allele (mutant) confers some increase in the **fitness** of the organism
- Selection acts to favour this allele
- Also called adaptive selection or Darwinian selection.

NOTE: **Fitness** = ability to survive and reproduce

Modified from from [www.tcd.ie/Genetics/staff/Aoife/GE3026/GE3026\\_1+2.ppt](http://www.tcd.ie/Genetics/staff/Aoife/GE3026/GE3026_1+2.ppt)

## Advantageous allele

Herbicide resistance gene in nightshade plant

*Solanum nigrum* (nightshade) petA gene:

Normal sequence:  
 ... C G A T T G A T C C A A T A T G C T A C A C A C T C T ...  
 R L I F Q Y A S P N H S

Atrazine-resistant mutant:  
 ... C G A T T G A T C C A A T A T G C T A C A C A C T C T ...  
 R L I F Q Y A S P N H S

The mutation is a G to A change in the 5th codon, changing the amino acid from Serine to Glutamine.

Modified from [www.ted.ie/Genetics/staff/Aoife/GE3026/GE3026\\_1+2.ppt](http://www.ted.ie/Genetics/staff/Aoife/GE3026/GE3026_1+2.ppt)

## Negative selection ( $s < 0$ )

- A new allele (mutant) confers some decrease in the fitness of the organism
- Selection acts to remove this allele
- Also called purifying selection

Modified from [www.ted.ie/Genetics/staff/Aoife/GE3026/GE3026\\_1+2.ppt](http://www.ted.ie/Genetics/staff/Aoife/GE3026/GE3026_1+2.ppt)

## Neutral mutations

- Neither advantageous nor disadvantageous
- Invisible to selection (no selection)
- Frequency subject to 'drift' in the population
- Random drift** – random changes in small populations

## Types of Mutation-Substitution

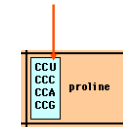
- Replacement of one nucleotide by another
- Synonymous (Doesn't change amino acid)
  - Rate sometimes indicated by  $K_s$
  - Rate sometimes indicated by  $d_s$
- Non-Synonymous (Changes Amino Acid)
  - Rate sometimes indicated by  $K_a$
  - Rate sometimes indicated by  $d_n$

(this and the following 4 slides are from [mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/Lecture7.ppt](http://mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/Lecture7.ppt))

## Genetic Code – Note degeneracy of 1<sup>st</sup> vs 2<sup>nd</sup> vs 3<sup>rd</sup> position sites

UUU phenylalanine	UUC	UUA leucine	UUG	UUU phenylalanine	UUC	UUA leucine	UUG
CUU leucine	CUC	CUA leucine	CUG	CUU leucine	CUC	CUA leucine	CUG
AUU isoleucine	AUC	AUA isoleucine	AUG methionine	AUU isoleucine	AUC	AUA isoleucine	AUG methionine
GUU valine	GUC	GUA valine	GUG	GUU valine	GUC	GUA valine	GUG

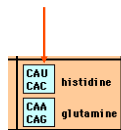
## Genetic Code



Four-fold degenerate site – Any substitution is synonymous

From: [mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/Lecture7.ppt](http://mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/Lecture7.ppt)

## Genetic Code



Two-fold degenerate site – Some substitutions synonymous, some non-synonymous

From: [mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/Lecture7.ppt](http://mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/Lecture7.ppt)

## Genetic Code

Degeneracy of 1<sup>st</sup> vs 2<sup>nd</sup> vs 3<sup>rd</sup> position sites results in 25.5% synonymous changes and 74.5% non synonymous changes (Yang & Nielsen, 1998).

First Position	U	C	A	G
Second Position	UUU phe UUC phe UUA leu UUG leu CUU leu CUC leu CUA leu CUG leu AUU ile AUC ile AUA ile AUG met GUU val GUC val GUA val GUG val	UCU phe UCC phe UCA ser UCG ser CCU pro CCC pro CCA pro CCG pro ACU thr ACC thr ACA thr ACG thr GCU ala GCC ala GCA ala GCG ala	UAU tyr UAC tyr UAA stop UAG stop CAU his CAC his CAA his CAG his AAU asp AAC asp AAA glu AAG glu	UGU cys UGC cys UGA stop UGG trp CGU arg CGC arg CGA arg CGG arg AGU ser AGC ser AGA arg AGG arg GGU gly GGC gly GGA gly GGG gly

## Measuring Selection on Genes

- Null hypothesis = neutral evolution
- Under neutral evolution, synonymous changes should accumulate at a rate equal to mutation rate
- Under neutral evolution, amino acid substitutions should also accumulate at a rate equal to the mutation rate

From: [mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/Lecture7.ppt](http://mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/Lecture7.ppt)

## Counting #s/#a

Species1	Ser	Ser	Ser	Ser	Ser
	TGA	TGC	TGT	TGT	TGT
Species2	Ser	Ser	Ser	Ser	Ala
	TGT	TGT	TGT	TGT	GGT

#s = 2 sites  
#a = 1 site

#a/#s=0.5

To assess selection pressures one needs to calculate the rates (Ka, Ks), i.e. the occurring substitutions as a fraction of the possible syn. and nonsyn. substitutions.

Things get more complicated, if one wants to take transition transversion ratios and codon bias into account. See chapter 4 in Nei and Kumar, Molecular Evolution and Phylogenetics.

Modified from: [mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/lecture7.pdf](http://mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/lecture7.pdf)

## Testing for selection using dN/dS ratio

dN/dS ratio (aka Ka/Ks or  $\omega$  (omega) ratio) where

dN = number of non-synonymous substitutions / number of possible non-synonymous substitutions

dS = number of synonymous substitutions / number of possible non-synonymous substitutions

dN/dS > 1 positive, Darwinian selection

dN/dS = 1 neutral evolution

dN/dS < 1 negative, purifying selection

## PAML (codeml) the basic model

$$q_{ij} = \begin{cases} 0, & \text{if the two codons differ at more than one position,} \\ \pi_{ij}, & \text{for synonymous transversion,} \\ \kappa\pi_{ij}, & \text{for synonymous transition,} \\ \omega\pi_{ij}, & \text{for non-synonymous transversion,} \\ \omega\kappa\pi_{ij}, & \text{for non-synonymous transition,} \end{cases}$$

The equilibrium frequency of codon  $j$  ( $\pi_j$ ) can be considered a free parameter, but can also be calculated from the nucleotide frequencies at the three codon positions (control variable CodonFreq). Under this model, the relationship holds that  $\omega = d_N/d_S$ , the ratio of non-synonymous/synonymous substitution rates. This basic model is fitted by specifying model = 0 NSites = 0, in the control file codeml.c. It forms the basis for more sophisticated models implemented in codeml.

## dambe

Three programs worked well for me to align nucleotide sequences based on the amino acid alignment,

One is **DAMBE** (works well for windows). This is a handy program for a lot of things, including reading a lot of different formats, calculating phylogenies, it even runs codeml (from PAML) for you.

The procedure is not straight forward, but is well described on the help pages. After installing DAMBE go to HELP -> general HELP -> sequences -> align nucleotide sequences based on ...->

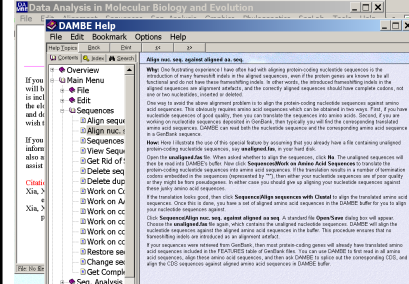
If you follow the instructions to the letter, it works fine.

DAMBE also calculates Ka and Ks distances from codon based aligned sequences.

Alternatives are

- **tranalign** from the **EMBOSS** package, and
- **Seaview** (see below)

## dambe (cont)

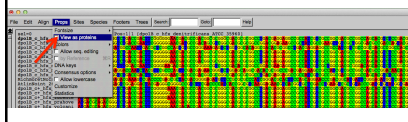


## Codon based alignments in Seaview

Load nucleotide sequences (no gaps in sequences, sequence starts with nucleotide corresponding to 1<sup>st</sup> codon position)



Select view as proteins



## Codon based alignments in Seaview

With the protein sequences displayed, align sequences



Select view as nucleotides



## PAML (codeml) the basic model

$$q_{ij} = \begin{cases} 0, & \text{if the two codons differ at more than one position,} \\ \pi_{ij}, & \text{for synonymous transversion,} \\ \kappa\pi_{ij}, & \text{for synonymous transition,} \\ \omega\pi_{ij}, & \text{for non-synonymous transversion,} \\ \omega\kappa\pi_{ij}, & \text{for non-synonymous transition,} \end{cases}$$

The equilibrium frequency of codon  $j$  ( $\pi_j$ ) can be considered a free parameter, but can also be calculated from the nucleotide frequencies at the three codon positions (control variable CodonFreq). Under this model, the relationship holds that  $\omega = d_N/d_S$ , the ratio of non-synonymous/synonymous substitution rates. This basic model is fitted by specifying model = 0 NSites = 0, in the control file codeml.c. It forms the basis for more sophisticated models implemented in codeml.

## sites versus branches

You can determine omega for the whole dataset; however, usually not all sites in a sequence are under selection all the time.

PAML (and other programs) allow to either determine omega for each site over the whole tree, **Branch Models**, or determine omega for each branch for the whole sequence, **Site Models**.

It would be great to do both, i.e., conclude codon 176 in the vacuolar ATPases was under positive selection during the evolution of modern humans – alas, a single site does not provide much statistics ...

## Sites model(s)

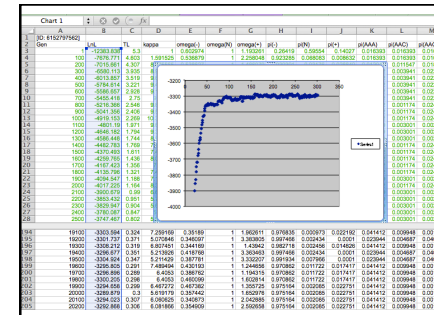
work great have been shown to work great in few instances.  
The most celebrated case is the influenza virus HA gene.

A talk by Walter Fitch (slides and sound) on the evolution of this molecule is [here](#).  
This [article by Yang et al. 2000](#) gives more background on ml approaches to measure omega. The dataset used by Yang et al is here: [flu\\_data.pau.n](#).

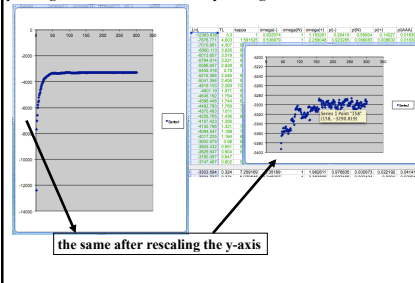
## sites model in MrBayes

The MrBayes block in a nexus file might look something like this:

```
begin mrbayes;
set autoclose=yes;
lset nst=2 rates=gamma nucmodel=codon omegaivar=Ny98;
mcmc samplefreq=500 printfreq=500;
mcmc ngen=500000;
sump burnin=50;
sumt burnin=50;
end;
```



## plot LogL to determine which samples to ignore



f236										
	A	B	C	D	E	F	G	H	I	K
2	1	2	3	4	5	6	7	8	9	10
3	1	2	3	4	5	6	7	8	9	10
4	1	2	3	4	5	6	7	8	9	10
5	1	2	3	4	5	6	7	8	9	10
6	1	2	3	4	5	6	7	8	9	10
7	1	2	3	4	5	6	7	8	9	10
8	1	2	3	4	5	6	7	8	9	10
9	1	2	3	4	5	6	7	8	9	10
10	1	2	3	4	5	6	7	8	9	10

BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU

for each codon calculate the average probability										
	A	B	C	D	E	F	G	H	I	K
2	1	2	3	4	5	6	7	8	9	10
3	1	2	3	4	5	6	7	8	9	10
4	1	2	3	4	5	6	7	8	9	10
5	1	2	3	4	5	6	7	8	9	10
6	1	2	3	4	5	6	7	8	9	10
7	1	2	3	4	5	6	7	8	9	10
8	1	2	3	4	5	6	7	8	9	10
9	1	2	3	4	5	6	7	8	9	10
10	1	2	3	4	5	6	7	8	9	10

## To determine credibility interval for a parameter (here omega<1):

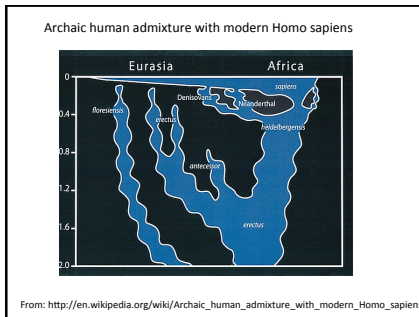
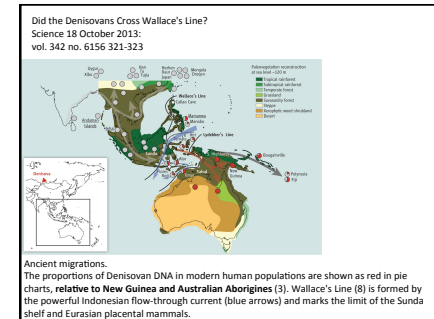
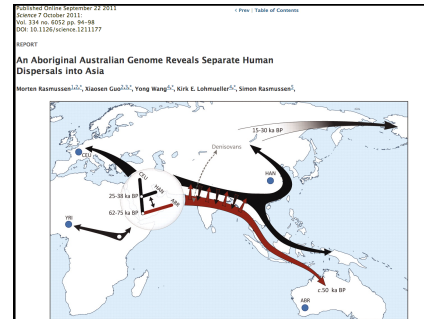
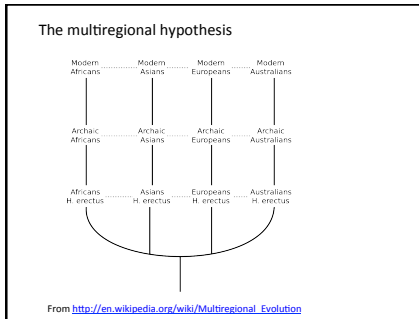
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU

BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU
BT	BU	BU	BU	BU	BU	BU	BU	BU	BU	BU









For more discussion on archaic and early humans see:  
[http://en.wikipedia.org/wiki/Denisova\\_hominin](http://en.wikipedia.org/wiki/Denisova_hominin)  
<http://www.nytimes.com/2012/01/31/science/gains-in-dna-are-speeding-research-into-human-origins.html>  
<http://www.sciencedirect.com/science/article/pii/S0002929711003958>  
<http://www.abc.net.au/science/articles/2012/08/31/3580500.htm>  
<http://www.sciencemag.org/content/334/6052/94.full>  
<http://www.sciencemag.org/content/334/6052/94/F2.expansion.html>  
<http://haplogroup-a.com/Ancient-Root-AJHG2013.pdf>