

MCB5472

Computer methods in  
molecular evolution

Lecture 4/14/2014

Lecture 4/14/2014

MCrobot demo?

<http://hydrodictyon.eeb.uconn.edu/people/plewis/software.php>

<http://hydrodictyon.eeb.uconn.edu/people/plewis/software.php>

# OldAssignment

• 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 recognise 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|2292440723|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](#).

• 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 `phym1`)

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](#).

```
#/usr/bin/perl
# namerewrite.pl
use strict; use warnings;

die "usage: namerewrite.pl <in> <out>" unless $ARGV == 1;
my $filename=<ARGV[0]>;
open(IN, "< $filename");
open(OUT, "> namerewrite.out");
my $line=1;
my $species="";
my $proto="";
while(defined(my $line=<IN>)){
    chomp($line);
    if ($line =~ /\w/)
        $line = s"/\w//g;
    $line = s"/\./g;
    $line = s"/\./g;
    $line = s"/\./g;
    my $spc="";
    my $proto="";
    $line = s"/[[:punct:]]//g;
    my $rearrange = "s/$spc/_/$proto";
    print "$rearrange\n";
}
else {
    $line = tr/mcga/ATCG/;
    $line = s"/\./g;
}
print "$line\n";
}
close(IN);
close(OUT);
```

```
# /usr/bin/perl
use strict; use warnings;

die "usage: %s name.range.pl <int> %s" unless $ARGV =~ 1;
my $filenumber=$ARGV[0];
open(CN, "<") or die "cannot open $filenumber";
open(OUT, ">") or die "cannot open $filenumber";
my $line;
my $specname;
while($line=<CN){
    chomp($line);
    if ($line =~ />/) {
        $line =~ s/\>/./g;
        $line =~ s/\</./g;
        $line =~ s/\(log(C)/log(C)/g;
        my $col=$line;
        $line =~ s/((C|))//g;
        my $range = "1..$col";
        print "line:range: $line", $range, "\n";
    }
    else {
        $line =~ s/((log|ATG|C|))//g;
        $line =~ s/\>/./g;
        print "line: $line";
    }
}
close(CN);
close(OUT);
```

[illegible][illegible]

See P17 for info on \$&

```
open(FILE, "gicasy1Transferases.fasta");

while ($line=FILEh>) {
  if ($line =~ /\>/) {
    $line= s/>\/\/g;
    @split = split ('\\', $line);
    @split = split ('\\', $split[1]);
    @split = split ('\\', $line);
    $split[0] =~ s/ / /g;
    print ">$split[1]$split[0] $line";
  } else {
    print $line;
  }
}
```

```
open (FILE, "glcosylTransferases.fasta");

while ($!=$?) {
    if ($! =~ /\>/) {
        $! =~ s/\>/;
        @split = split ('\\', $!);
        @splat = split ('\\', $!);
        $splat[0] =~ s/ /;
        print ">$splat[1]$splat[0] $!";
    } else {
        print $!;
    }
}
```

# HGT as a force a creative force

New biochemical pathways

Oxygen producing PS  
Acetolactate Methanogenesis (new) (cause of the Permian extinctions? new)

New substrates, new weapons, new resistance genes, breaks up linkage in case of selective sweeps.

Discussion:

Selfish genes versus altruism. (Evolutionary stable strategies).

- Group selection? (plasmid sharing in Agrobacteria after plant transformation)
- Under which conditions is it useful for an organism to sacrifice itself (e.g. GTAs), so that other members of the population reap a benefit? -> social parasites

Evolution of the holobiont? (sushi wrapper digesting intestinal symbionts)

- Is selection really acting on the holobiont?

## New biochemical pathways

Oxygen producing PS  
Acetoclastic Methanogenesis ([here](#)) (cause of the Permian extinctions? [here](#))

New substrates, new weapons, new resistance genes, breaks up linkage in case of selective sweeps.

Discussion:

- Selfish genes *versus* altruism. (Evolutionary stable strategies).
  - Group selection? (plasmid sharing in *Agrobacteria* after plant transformation)
  - Under which conditions is it useful for an organism to sacrifice itself (e.g. GTAs), so that other members of the population reap a benefit? → social parasites
- Evolution of the holobiont? (sushi wrapper digesting intestinal symbionts)
  - Is selection really acting on the holobiont?

# 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.

**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.**

### Genetic Drift

p: 0.5 N: 100 Generations: 250 Start Clear

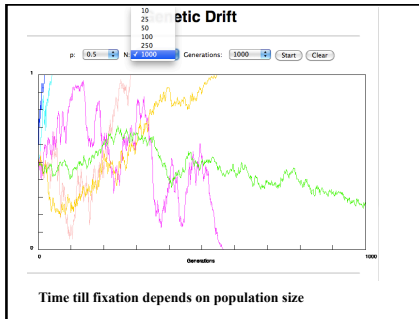
1  
p  
0

0 250  
Generations

If one waits long enough, one of two alleles with equal fitness will be fixed



**If one waits long enough, one of two alleles with equal fitness will be fixed**



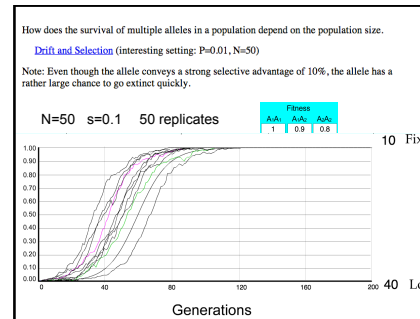
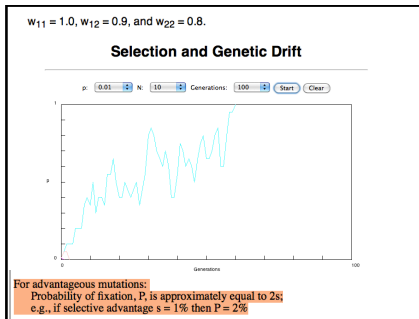
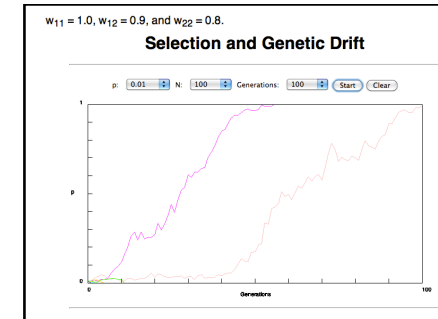
$s=0$

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



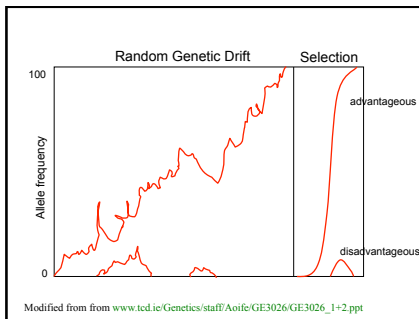
$s > 0$

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

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

$\Rightarrow$  substitution rate of mutation under positive selection is larger than the rate with which neutral mutations are fixed.  
 $\Rightarrow$  This is easily detected in case of diversifying selection.



**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) pepA gene:

Normal sequence:  
 ... CUA TTG ATC TTC CAA TAT GCT ...  
 ... CUA TTG ATC TTC CAA TAT GCT ...

Alanine-resistant mutant:  
 ... CUA TTG ATC TTC CAA TAT GCT ...  
 ... CUA TTG ATC TTC CAA TAT GCT ...

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)

## 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

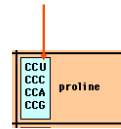
- Replacement of one nucleotide by another
- **Synonymous (Doesn't change amino acid)**
  - Rate sometimes indicated by  $K_s$
  - Rate sometimes indicated by  $d_s$
- **Non-Synonymous (Changes Amino Acid)**
  - Rate sometimes indicated by  $K_a$
  - Rate sometimes indicated by  $d_a$

(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	UCU serine	UAU tyrosine	UGU cysteine
UUC alanine	UCC serine	UAC tyrosine	UGC cysteine
UUA leucine	UCA serine	UAA stop	UGA stop
UUG leucine	UCG serine	UAG stop	UGG tryptophan
CUU leucine	CCU proline	CAU histidine	CGU arginine
CUC leucine	CCC proline	CAC histidine	CGC arginine
CUA leucine	CCA proline	CAA glutamine	CGA arginine
CUG leucine	CCG proline	CAG glutamine	CGG arginine
AUU isoleucine	ACU threonine	AAU asparagine	AGU serine
AUC isoleucine	ACC threonine	AAC asparagine	AGC serine
AUA isoleucine	ACA threonine	AAA lysine	AGA arginine
AUG methionine	ACG threonine	AAG lysine	AGG arginine
GUU valine	GCU alanine	GAU aspartic acid	GGU glycine
GUC valine	GCC alanine	GAC aspartic acid	GGC glycine
GUA valine	GCA alanine	GAA glutamic acid	GGA glycine
GUG valine	GCG alanine	GAG glutamic acid	GGG glycine

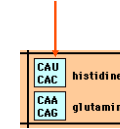
## Genetic Code



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

#s = 2 sites  
#a = 1 site  
#a/#s = 0.5

To assess selection pressures one needs to calculate the rates ( $K_a$ ,  $K_s$ ), 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.ppt](http://mentor.lscf.ucsb.edu/course/spring/eemb102/lecture/Lecture7.ppt)

### Testing for selection using dN/dS ratio

**dN/dS ratio (aka Ka/Ks or ω (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

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

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 nonsynonymous/synonymous substitution rates. This basic model is fitted by specifying model = 0 NSsites = 0, in the control file codeml.ctl. 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

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

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 nonsynonymous/synonymous substitution rates. This basic model is fitted by specifying model = 0 NSsites = 0, in the control file codeml.ctl. 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](#).

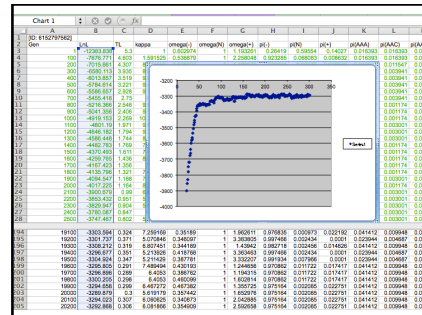
This [article by Yang et al. 2000](#) gives more background on ml approaches to measure omega. The dataset used by Yang et al is here: [flu\\_data.pau.n](#).



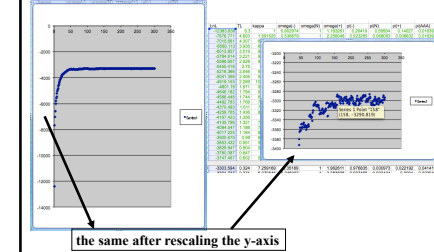
## sites model in MrBayes

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;
```



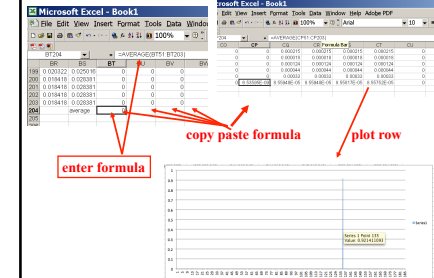
## plot LogL to determine which samples to ignore



F236										
DC 8157/9192										
Site	1	2	3	4	5	6	7	8	9	10
Seq	1	2	3	4	5	6	7	8	9	10
1	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
2	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
3	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
4	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
5	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
6	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
7	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
8	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
9	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
10	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000

Site	1	2	3	4	5	6	7	8	9	10
Seq	1	2	3	4	5	6	7	8	9	10
1	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
2	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
3	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
4	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
5	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
6	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
7	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
8	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
9	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
10	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000

## for each codon calculate the the average probability



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

Site	1	2	3	4	5	6	7	8	9	10
Seq	1	2	3	4	5	6	7	8	9	10
1	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
2	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
3	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
4	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
5	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
6	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
7	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
8	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
9	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
10	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000

Select values for the parameter, sampled after the burning.

Copy paste to a new spreadsheet,

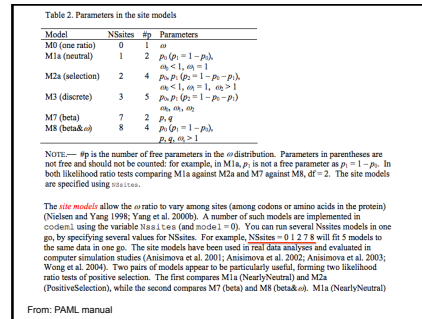
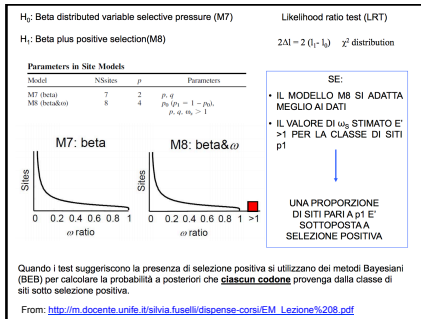
Site	1	2	3	4	5	6	7	8	9	10
Seq	1	2	3	4	5	6	7	8	9	10
1	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
2	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
3	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
4	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
5	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
6	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
7	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
8	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
9	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
10	12180.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000

- Sort values according to size,
- Discard top and bottom 2.5%
- Remainder gives 95% credibility interval.

Slides on codeml are at

<http://abacus.gene.ucl.ac.uk/zhigeng/data/pam1DEMO.pdf>  
[http://m.docente.unife.it/silvia.fuselli/dispense-corsi/EM\\_Lezione%208.pdf](http://m.docente.unife.it/silvia.fuselli/dispense-corsi/EM_Lezione%208.pdf)





Hy-Phy - Hypothesis Testing using Phylogenies.

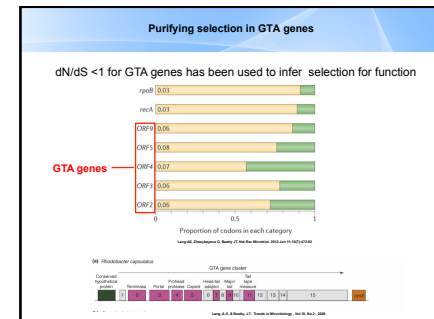
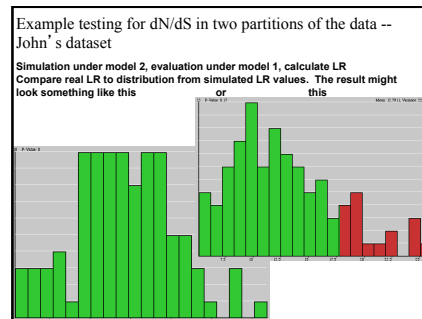
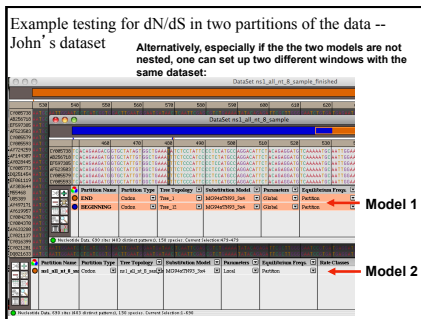
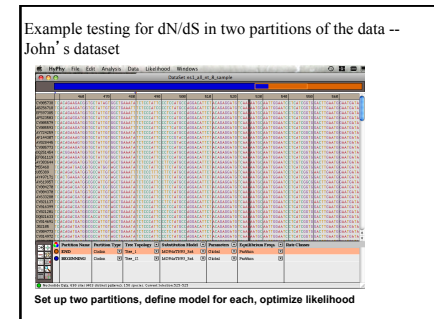
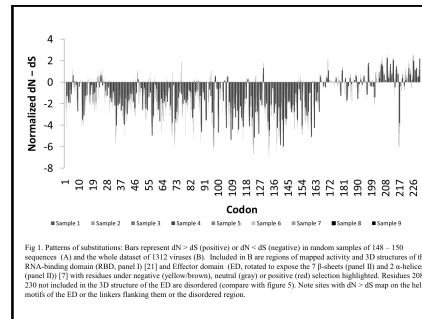
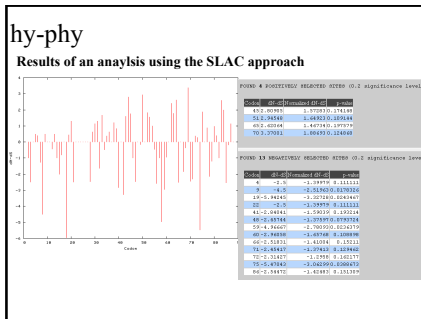
Using Batchfiles or GUI

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

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

**DATAMONKEY**  
 NEW HY-PHY ONLINE ANALYSIS

Supporting hy-phy that you cannot do the following questions: (1) How many sites are under positive selection? (2) Which sites are under positive selection? (3) Which sites are under negative selection? (4) Which sites are under neutral selection? (5) Which sites are under purifying selection? (6) Which sites are under diversifying selection? (7) Which sites are under balancing selection? (8) Which sites are under directional selection? (9) Which sites are under stabilizing selection? (10) Which sites are under disruptive selection? (11) Which sites are under fluctuating selection? (12) Which sites are under episodic selection? (13) Which sites are under relaxed selection? (14) Which sites are under increased selection? (15) Which sites are under decreased selection? (16) Which sites are under no selection? (17) Which sites are under unknown selection? (18) Which sites are under unknown selection? (19) Which sites are under unknown selection? (20) Which sites are under unknown selection?



## Purifying selection in *E. coli* ORFans

$dN-dS < 0$  for some ORFan *E. coli* clusters seems to suggest they are functional genes.

Gene groups	Number	$dN-dS < 0$	$dN-dS < 0$	$dN-dS < 0$
<i>E. coli</i> ORFan clusters	3773	944 (25%)	1953 (52%)	876 (23%)
Clusters of <i>E. coli</i> sequences found in <i>Salmonella</i> sp. <i>Croacacter</i> sp.	610	104 (17%)	423 (69%)	83 (14%)
Clusters of <i>E. coli</i> sequences found in some <i>Enterobacteriaceae</i> only	373	8 (2%)	365 (98%)	0 (0%)

Adapted after Yu, G. and Storz, A. *Genome Biol Evol* (2012) 16: 4 (176-187)

Vincent Daubin and Howard Ochman: Bacterial Genomes as New Gene Homes: The Genealogy of ORFans in *E. coli*.  
*Genome Research* 14:1036-1042, 2004

The ratio of non-synonymous to synonymous substitutions for genes found only in the *E. coli* - *Salmonella* clade is lower than 1, but larger than for more widely distributed genes.

A scatter plot showing the ratio of non-synonymous to synonymous substitutions ( $Ka/Ks$ ) on the y-axis (ranging from 0 to 0.25) against phylogenetic depth on the x-axis (categorized as  $n_2$ ,  $n_3$ , and  $n_4$ ). The legend indicates that open circles represent HOPs and filled circles represent ORFans. A horizontal dashed line is drawn at  $Ka/Ks = 0.05$ . A red arrow at the bottom points from left to right, labeled 'Increasing phylogenetic depth'. A green arrow points from the text on the left to the data point for ORFans at  $n_2$ .

Phylogenetic Depth	HOPs ( $Ka/Ks$ )	ORFans ( $Ka/Ks$ )
$n_2$	~0.07	~0.19
$n_3$	~0.09	~0.09
$n_4$	~0.06	~0.08

Fig. 3 from Vincent Daubin and Howard Ochman, *Genome Research* 14:1036-1042, 2004

**Trunk-of-my-car analogy.** Hardly anything in there is the result of providing a selective advantage. Some items are removed quickly (purifying selection), some are useful under some conditions, but most things do not alter the fitness.

The image shows the open trunk of a dark-colored car. Inside the trunk, there is a large yellow potato on the left, crossed out with a large red 'X'. Next to it is a blue and yellow bag of Doritos, also crossed out with a red 'X'. In the center, there is a white box of tissues and a clear plastic bottle of water. To the right, there is a small white container. On the far right, two red cylindrical objects labeled 'TNT' are shown, each crossed out with a red 'X'. The background of the trunk is lined with a grey material.

*Could some of the inferred purifying selection be due to the acquisition of novel detrimental characteristics (e.g., protein toxicity; HOPELESS MONSTERS)?*

## Other ways to detect positive selection

Selective sweeps -> fewer alleles present in population  
(see contributions from archaic Humans for example)

Repeated episodes of positive selection -> high  $dN$

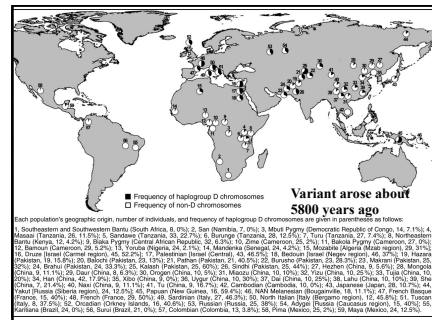


Fig. 3. Global frequencies of Microcephalin haplogroup D chromosomes (defined as having the derived C allele at the G3795C diagnostic SNP) in a panel of 1184 individuals

P. D. Evans et al., Science 309, 1717–1720 (2005)

The age of haplogroup D was found to be ~37,000 years

Published by AAS

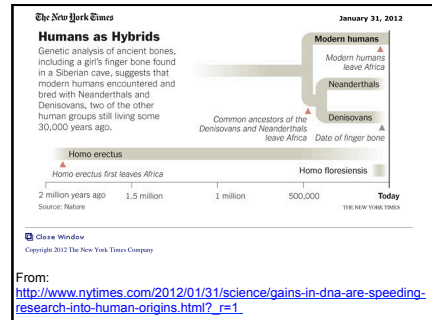
**Figure 1. Evidence that the adaptive allele of the brain size gene ENTPKG1 has been under positive selection in the archaic *Homo* lineage.**

**A**

Phylogenetic tree showing the divergence of hominids. The tree is rooted at the bottom left. The branches are labeled with time in millions of years (mya): 1.7 mya, 1.0 mya, 0.5 mya, 0.2 mya, and 0.1 mya. The tree shows the divergence of the archaic *Homo* lineage (red triangle) and the modern *Homo* lineage (blue triangle). The archaic *Homo* lineage is shown as a red triangle, and the modern *Homo* lineage is shown as a blue triangle. The tree is rooted at the bottom left.

**B**

Plot of nucleotide diversity ( $\pi$ ) for ENTPKG1 across the genome. The x-axis represents the genome in Mb. The y-axis represents  $\pi$  (0.000 to 0.002). The plot shows a significant peak in diversity in the archaic *Homo* lineage (red triangle) and a smaller peak in the modern *Homo* lineage (blue triangle). The archaic *Homo* lineage is shown as a red triangle, and the modern *Homo* lineage is shown as a blue triangle.



### The multiregional hypothesis

The diagram illustrates the multiregional hypothesis of human evolution. It features a single, wide-based tree structure. At the base is a common ancestor labeled 'H. erectus'. This ancestor branches into four distinct lineages, each representing a modern population. The lineages are labeled as follows: 'Modern Africans' (top left), 'Modern Asians' (top center-left), 'Modern Europeans' (top center-right), and 'Modern Australians' (top right). Each modern population is connected to its respective 'H. erectus' ancestor by a vertical line. The tree structure suggests that all modern human populations share a common ancestor and have evolved in parallel from that ancestor, with no significant interbreeding or gene flow between the populations over time.

Modern Africans

Modern Asians

Modern Europeans

Modern Australians

Archaic Africans

Archaic Asians

Archaic Europeans

Archaic Australians

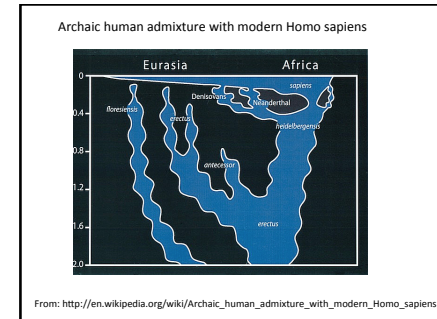
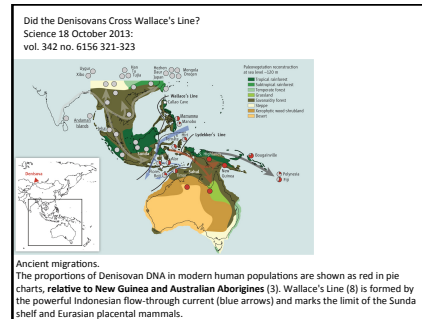
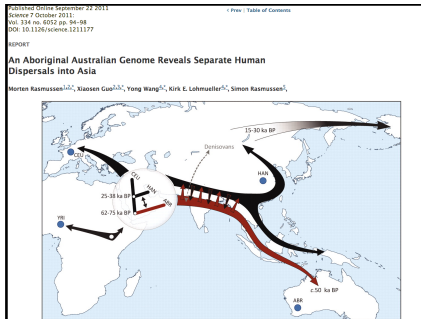
Africans  
H. erectus

Asians  
H. erectus

Europeans  
H. erectus

Australians  
H. erectus

From [http://en.wikipedia.org/wiki/Multiregional\\_Evolution](http://en.wikipedia.org/wiki/Multiregional_Evolution)



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>