

# MCB 5472

Student Projects,  
unix, Perl

*J. Peter Gogarten*

Office: *BPB 404*

phone: *860 486-4061,*

Email: *[gogarten@uconn.edu](mailto:gogarten@uconn.edu)*

# Student Projects

- Should be related to your interests !!!
- Examples for possible projects:

# Example: Evolution of a gene family

- When in the evolution of the interferon (or whatever you are interested in) gene family did gene duplications occur?
- Which of the resulting subfamilies (if any) have acquired a new function?
- What is the phylogenetic distribution of this subfamily? (Would you expect members of this subfamily to be present in insects, fish, chicken, fungi, archaea?)
- Can you detect episodes of positive selection or of relaxed purifying selection?
- Is there anything that would suggest gene conversion events?

The “to-do-list” would include:

- gather data (note for some of the questions mentioned above you’ll need aa **and** nucleotide sequences),
- align sequences
- build phylogenies
- analyze sequences
- assess reliability of branches
- INTERPRET WHAT YOU GOT!

# Example: Can one detect a distinct second divergence peak in the divergence of putatively chimeric genomes?

Genome fusions are the latest rage in evolutionary biology:

For example:

- Koonin EV, Mushegian AR, Galperin MY, Walker DR. *Comparison of archaeal and bacterial genomes: computer analysis of protein sequences predicts novel functions and suggests a chimeric origin for the archaea.* Mol Microbiol. 1997 Aug;25(4):619-37.
- The Eukaryotes are a chimera of at least an archaeal like host cell and a bacterium that evolved into a mitochondrion (+ in some cases a cyanobacterium that evolved into a plastid)
- The Haloarchaea contain many bacterial genes
- The Thermotogales contain many archaeal genes
- Most plants and many fungi (likely including bakers yeast) are aneupolyploids

In most of these instances it is not clear that the transfer (or duplication) really occurred in a single massive event, or if the transfers (duplications) occurred on a gene by gene basis.

(in yeast the type of genes that were duplicated suggest distinct selection pressures, see Benner et al [here](#))

Example: Chimera? continued

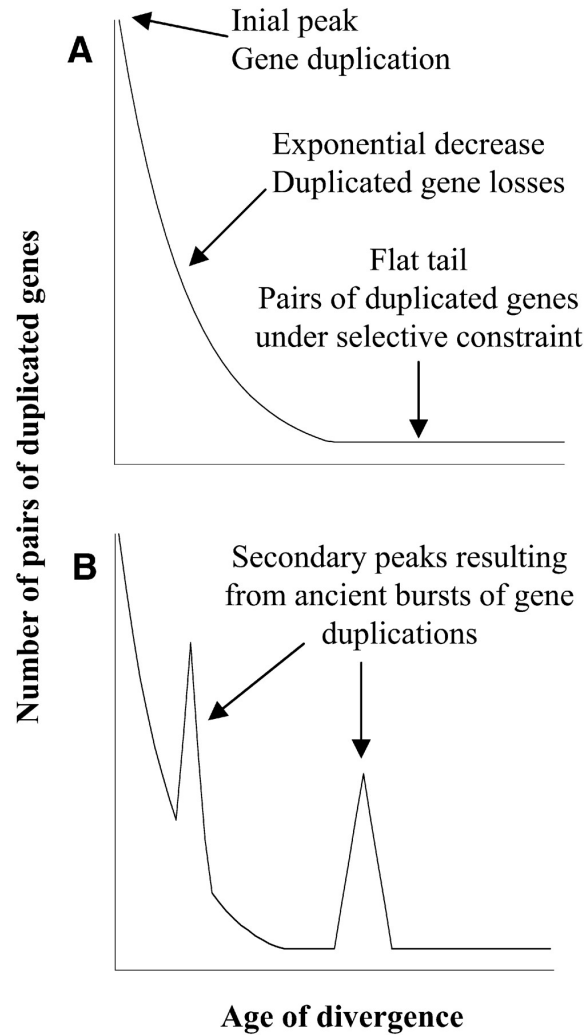
In case of a chimera formed in a single historic event one would expect

A) Two distinct types of phylogenetic affinity.

E.g.: Genes in *Thermotoga maritima* should either group with the sistergroup of the bacterial partner, or with the sistergroup of the archaeal donor

**Related: Ancient genome duplication events are revealed by peaks in the divergence of paralogs.**

# Theoretical Age Distributions of Pairs of Duplicated Genes in a Genome

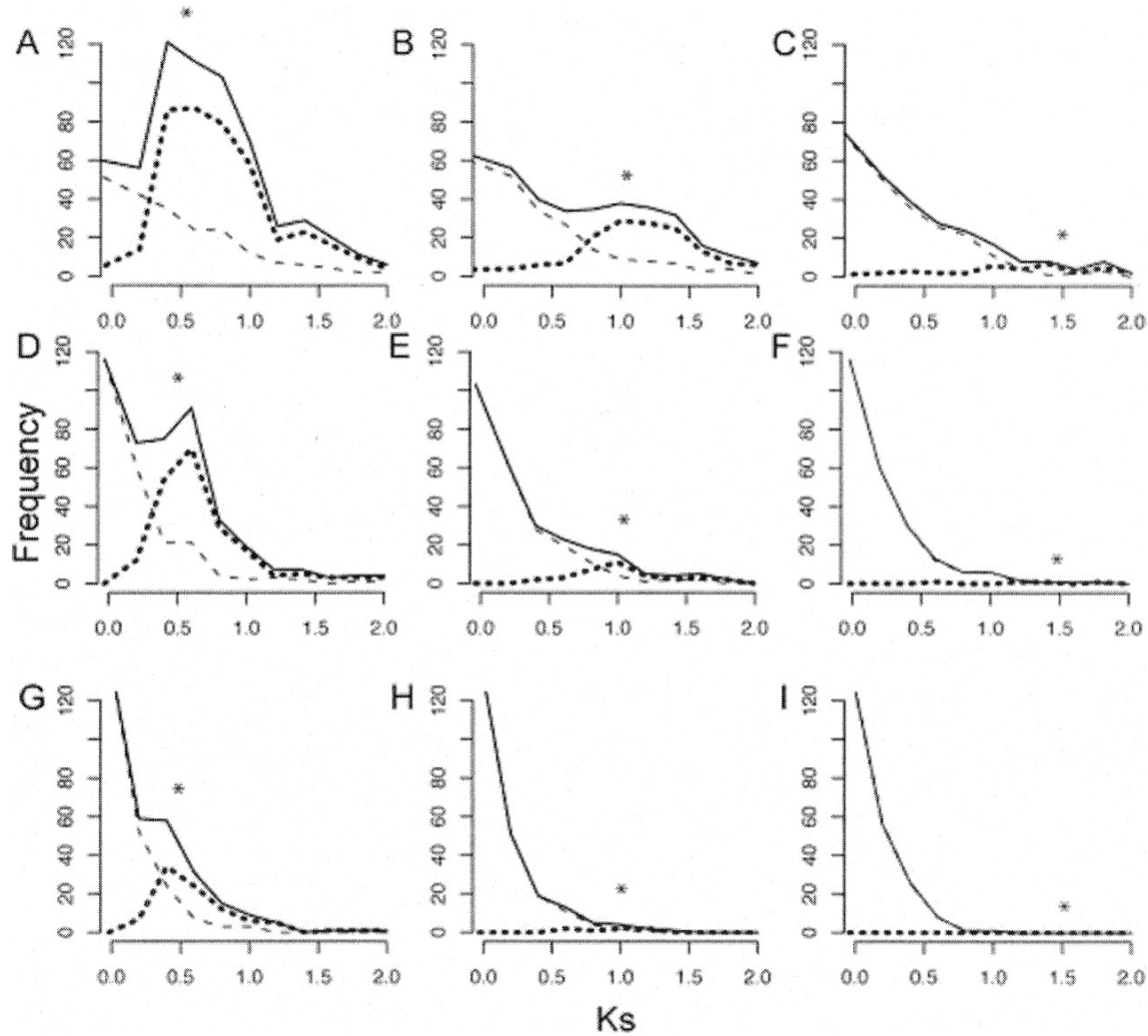


Blanc, G., et al. *Plant Cell* 2004;16:1667-1678



Effect of gene death rate and time of genome duplication on the Ks distribution for paralogs.

Higher death rate of duplicated genes

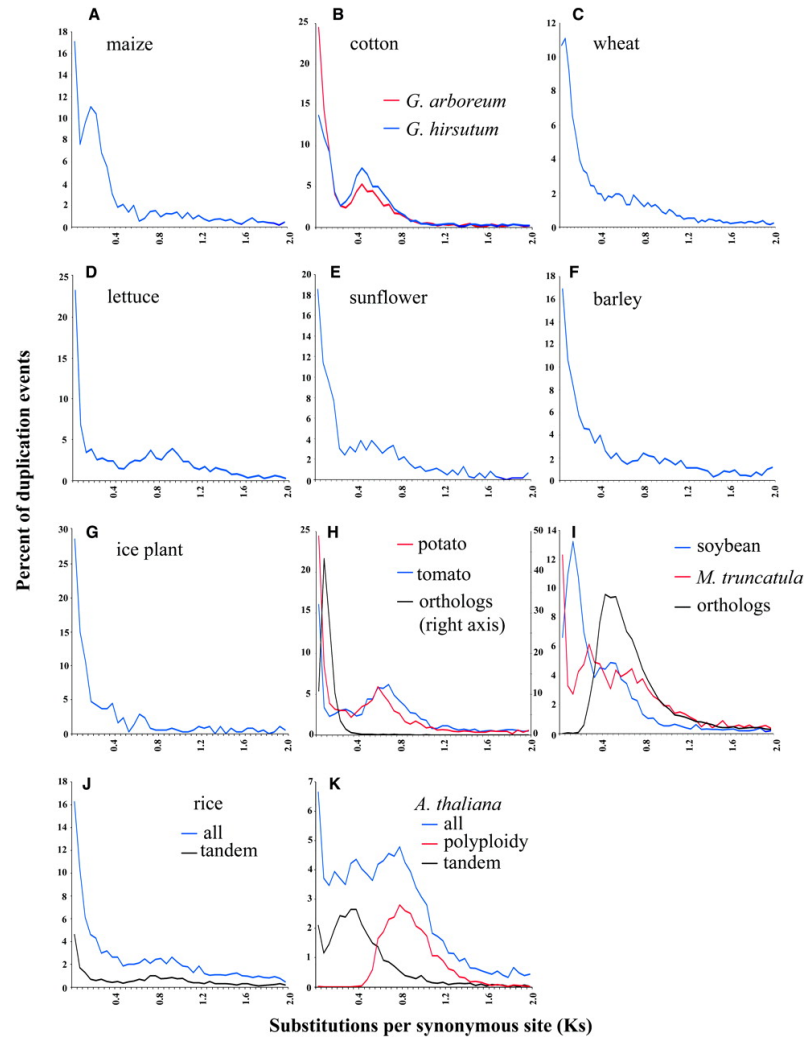


Cui L et al. Genome Res. 2006;16:738-749

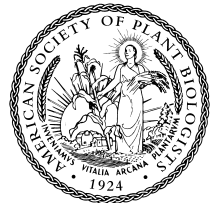
Duplication (\*) more ancient



# Distributions of the Fraction of Duplication Events as a Function of Their Levels of Synonymous Substitution for 14 Model Plant Species



Blanc, G., et al. Plant Cell 2004;16:1667-1678



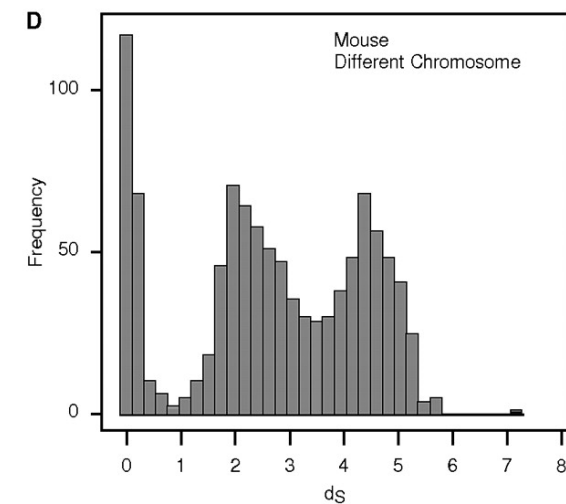
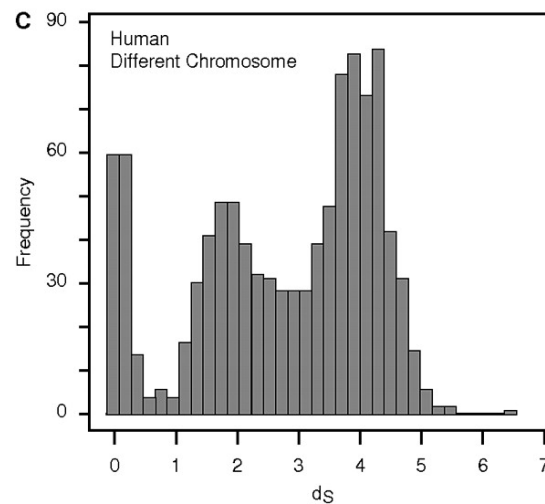
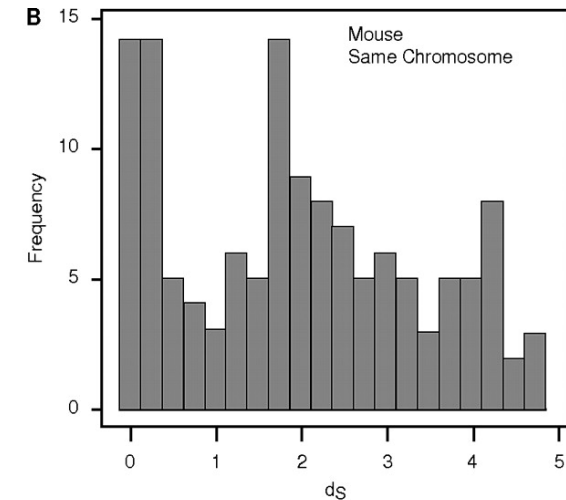
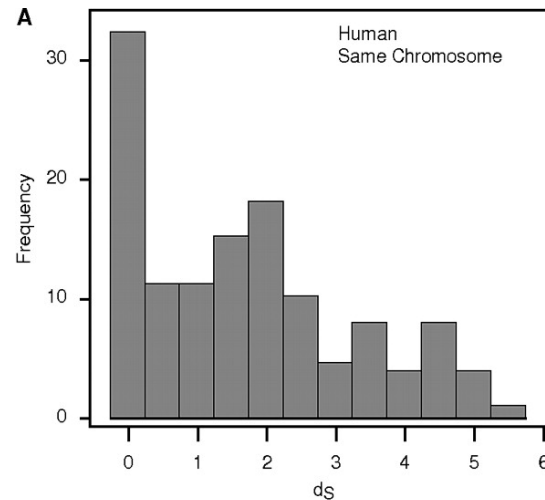


# Example: Gene versus Genome Duplications

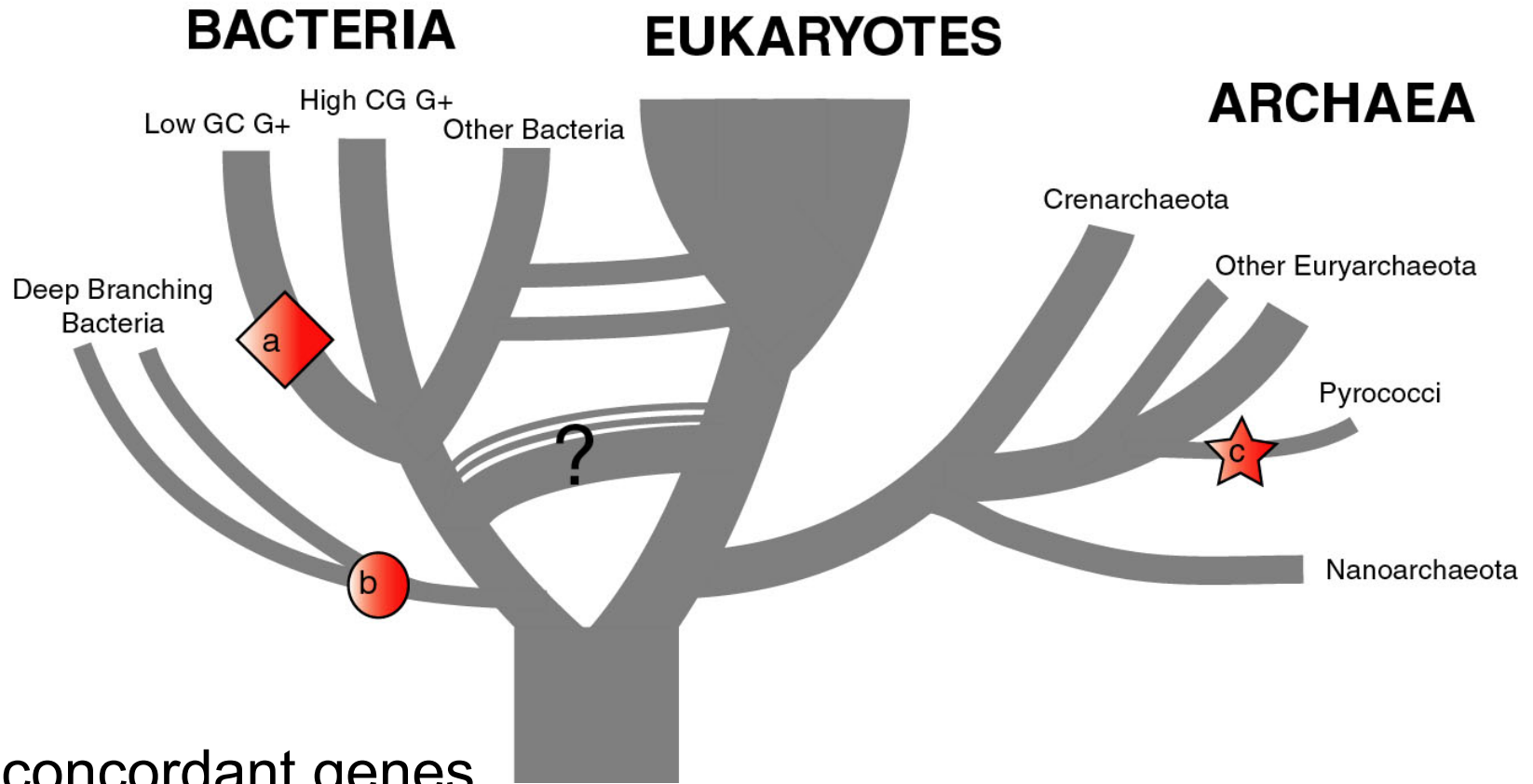
The same approach as suggested for the chimera formation can be applied to the question was the whole genome or a large segment of an organism's genome duplicated, or did the duplications occur in a piecemeal fashion?

Frequency distributions of  $d_S$  in human and mouse between the members of two-member gene families located on the same and different chromosomes

From: Robert Friedman and Austin L. Hughes: *Two Patterns of Genome Organization in Mammals: the Chromosomal Distribution of Duplicate Genes in Human and Mouse*. *Mol. Biol. Evol.* 21(6):1008–1013. 2004



# The Phylogenetic position of *Thermotoga maritima*



- (a) concordant genes,
- (b) according to 16S (and other conserved genes)
- (c) according to phylogenetically discordant genes

Gophna, Doolittle & Charlebois: Weighted genome trees: refinements and applications. *J. Bacteriol.* [here](#)

Gogarten & Townsend: Horizontal gene transfer, genome innovation, and evolution  
Nature Reviews in Microbiology 3(9) 679-687 ([pdf](#))

# Chimera Example, continued

The “to-do-list” would include:

- Formulate the question you want to address
- Download and analyze the required genomes
- Run blastall (this might take a couple of hours)
- Analyze the results in an Excel spreadsheet
- Selected some genes (e.g., the ones that are most archaeal), assemble gene families and reconstruct their phylogenies.
- **INTERPRET YOUR RESULTS!** What does it all mean?

# Background for group selection Example:

## Selection acts on

- **genes** (as in the selfish gene theory, the genes are the replicators that build the body of the organism). According to this all genes are selfish, most are cooperating with one another, a few are not. To distinguish the latter from the former, I call them parasitic genes (or molecular parasites).
- **individuals** in a population (the survival of the fittest).
- **groups** of organisms (group selection). The group that has properties that allows it to adapt better, or to evolve faster, or to make better use of resources will be selected. In this case the group (community, *not necessarily all belonging to the same species*) is the unit of selection. (see group selection entry at [wikipedia](#))

**Note: in general this is controversial. To what extent is group selection reflecting kin-selection: the organism acting to guarantee survival of genes that are related to its own genes (bees in a beehive are all closely related).**

# Examples for “group selection” in microbes: (a) *Agrobacteria*

*Agrobacteria* that carry a Ti plasmid can transform plant cells with a T DNA. As result of a successful transformation the plant cell has integrated the T DNA into its genome and expresses the encoded genes. This results in the transformed cells forming a tumor, and, in addition, the transformed plant cells also produce a strange amino acid that cannot be utilized by the plant cells, but that serves as a carbon and nitrogen source for the *Agrobacteria*. The genes responsible for transferring the Ti plasmid between different *Agrobacteria* (*tra* genes) are under the control of quorum sensing. The effect is that if one *Agrobacterium* strain has successfully transformed a plant, and now lives from the plant produced strange amino acid, other *Agrobacteria* can receive the Ti plasmid, which contains the T DNA transferred into the plant and in addition encodes enzymes that allow the metabolism of the strange amino acids. The *Agrobacteria*, which receive the Ti-plasmid thus participate in the utilization of the plant produced carbon and nitrogen source. This observation **might be described as group selection**: the population of *Agrobacteria* avoids a selective sweep and carries larger genetic diversity into the population living on the transformed plant. The increased diversity will facilitate future adaptations to a changing environment, and will avoid the fixation of slightly deleterious mutations that might have been carried by the *Agrobacterium* that transformed the plant cell. On the other hand, **one can consider this process the outcome of the "selfishness" of the *tra*-genes and of the Ti plasmid**. These genes manage to move themselves into the growing part of the population, and they will benefit from a more diverse group of host organisms.

Examples for “group selection” in microbes (b):  
*Metal resistance genes in microbial communities  
inside rocks in the dry valleys of Antarctica*

These rocks have high concentrations of toxic heavy metals. The endolithic microbial community readily shares heavy metal resistant genes with microbes that might be able to become part of the community. At the community level the outcome is a higher diversity, and a richer network of metabolic reactions. Presumably the more diverse communities are more stable towards perturbations, and provided the community can propagate as a whole, this would provide a **selective advantage to the community**. However from the **selfish gene point of view**, the resistance gene increases its chances of long term survival by invading as many additional species as possible.

## Examples for “group selection” in microbes ( c): Gene Transfer Agents (GTA) in alpha proteobacteria

GTA are prophages that do not specifically pack their own DNA, but that unselectively pack host DNA into the phage head (see [here](#)).

- Are these just defective prophages that lost their sequence specificity in DNA packaging?
- Is this an illustration that HGT is beneficial and under group selection?

(Aside: In general, HGT might reflect uptake of DNA for food, recombination might be a negligible side effect (Rosi Redfield, e.g. [here](#)), or HGT might reflect the selfishness of the transferred DNA.

# Are the genes involved in gene transfer under purifying selection

Possible hypotheses:

- GTAs are defective prophages that lost their sequence specificity in DNA packaging?
- GTAs evolved from phages but now benefit the group and are under group selection?

Under #2:

- The GTA should be more related to one another than to functioning phage
- Their molecular phylogeny should reflect the phylogeny of the organism (as measured by rRNA and ribosomal proteins)
- The genes encoding the GTA should be under strong purifying selection (under #1 they should be pseudogenes).



A manuscript reporting GTA genes under purifying selection is in preparation.

Does a  $dN/dS$  ratio  $<1$  reflect purifying selection for function, or could it reflect selection against becoming detrimental.

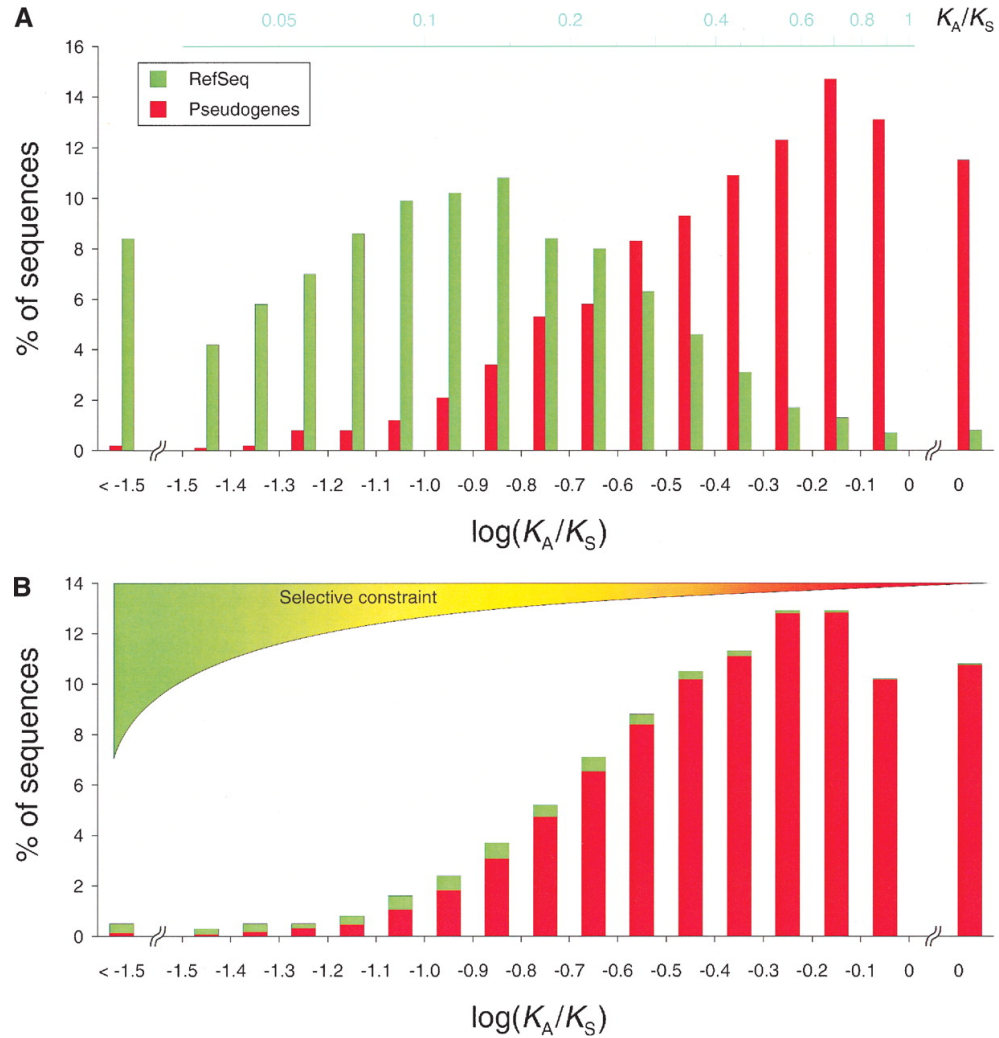
How many sites are inferred to be under purifying selection?

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.



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

# KA/KS distributions of benchmark and candidate sets (Human genome).



Torrents D et al. Genome Res. 2003;13:2559-2567



# 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.

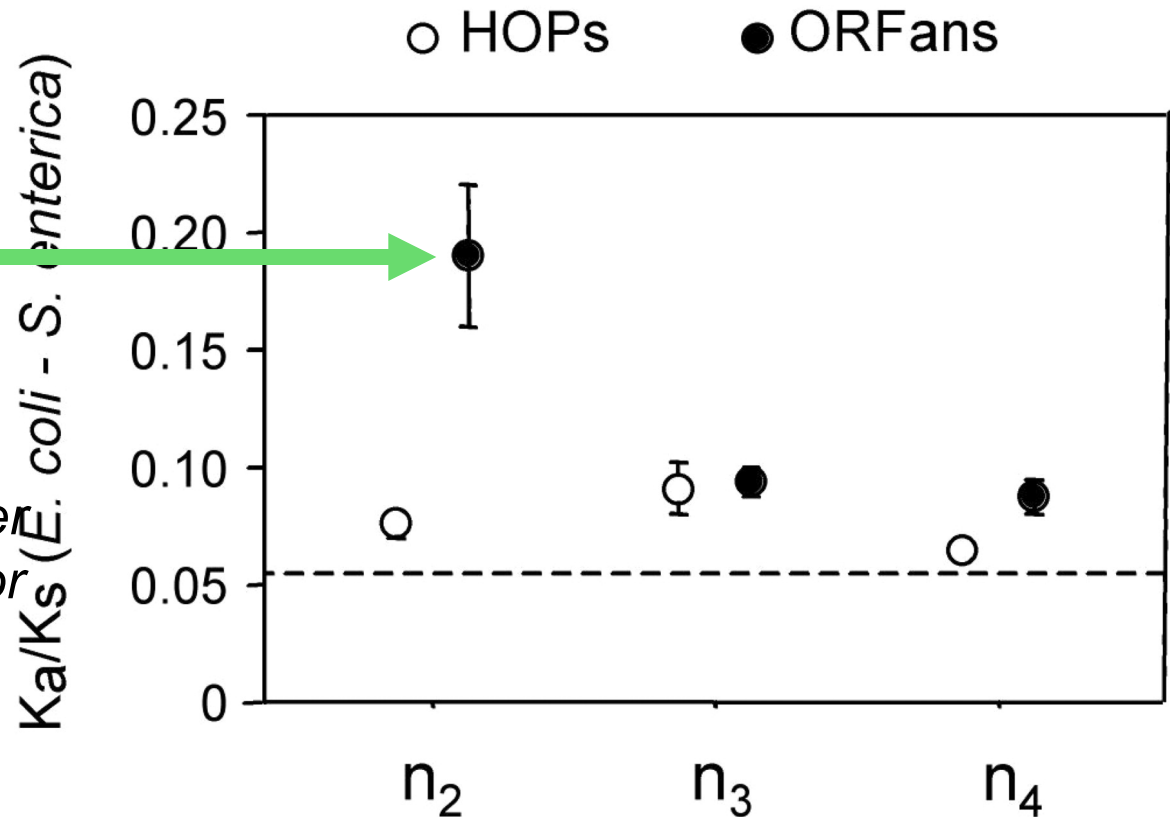


Fig. 3 from Vincent Daubin and Howard Ochman, *Genome Research* 14:1036-1042, 2004

A manuscript reporting GTA genes under purifying selection is in preparation.

Does a  $dN/dS$  ratio  $<1$  reflect purifying selection for function, or could it reflect selection against becoming detrimental.

How many sites are inferred to be under purifying selection?

Could one use prophage or transposases that have survived in a lineage for a long time as a test case?

# other ideas:

- Write a script that uses the 100+ known intein alleles each as a seed in PSI BLAST, and stores the profiles. Write a second script that uses these profiles to detect putative inteins in completely sequenced genomes, in metagenomes, and in whole genome sequences (in case of ngs, you want to use the 6 frame translation of the genome).  
(The current [InBase](#) is out of date, and rumors have it that NEB is no longer supporting research in inteins)
- Same as above but use transposases, integrases, homing endonucleases, or a molecular parasite of your choice as a seed.
- Form families for all genes from Thermotogales, add the fifteen most similar sequences from reference genomes, calculate phylogenies, screen for polyphyly of Thermotogales, screen for conflict with consensus.
- *Thermus thermophilus* is a naturally competent bacterium (with an archaeal type ATPase).  
Screen *Thermus* genomes (there are several from closely related organisms) for (a) genes recently acquired (using darkhorse – blast based; and composition based approaches), (b) gene replacements (ancient and recent). We are especially interested in genes that were replaced by a distant homolog, and that might be candidates for experimental reversion of the replacement.

# UNIX

## Basic UNIX commands

ls, cd, chmod, cp, rm, mkdir, more (or) less, vi, ps, kill -9, man  
A brief listing is [here](#)

**chmod** is a particular pain in the ... .

Under unix every file has an owner and the owner, his group and everyone else have permissions to read, write and/or execute the file (or they don't). If you want to see which permissions are currently assigned to your files, type `ls -l` at the command prompt.

`chmod a+x *.pl` gives everyone execute permission for all files that end with `.pl` the `*` is a wildcard. (warning don't ever use `rm` in conjunction with `*`)

For more on `chmod` type `"man chmod"` or see [here](#).

(In the OSX GUI you can control click at a file, and change permissions in the info box). Most ssh clients (FUGU and SSH) allow you to use a GUI to change file permissions (in FUGU ctrl click).

# Unix - command line interface

**If you tried to execute a command, and you made a mistake, for example, you mistyped a file name, you can recall the last command using the up arrow (down arrow for more recent).**

**If you are tired typing long filenames, you can use the tab key to complete the line, provided there is only one way to complete the line. E.g: `cd /Desktop` could be replaced by `cd /D<tab>`**

**If there are two or more choices you hear a boing, if you hit `<tab>` again, you get a list of choices.**



# writing Perl scripts

Use unix/ linux /OsX if possible.

A) open a terminal window ; type "which perl <return>"

B) SSH to a unix machine (cluster OsX), log in, type "which perl <return>"

C) to check the version type perl -v <return>The response of the system should tell you, where Perl is installed on your machine (you need to know this for the first line of your perl program, which tells the operating system how to interpret what follows. On most installations this is `#!/usr/bin/perl` ).

**WINDOWS:** If you use a windows machine, you can use an ssh program to connect to the biotech cluster. A good ssh client is available at <ftp://ftp.ssh.com/pub/ssh/>- highly recommended. I am sure that there are editors available that are more useful than notepad, but I don't know of them. :(

**MAC OsX:** If you use a Mac under OS X, and you do not want to (only) use the PERL locally, you want to install both jellyfish (ssh terminal) and fugu (a secure file transfer program) or FileZilla . Both are available at <ftp://ftp.uconn.edu/pub/packages/ssh/mac/> or through the people who wrote the software - GOOGLE)

Also, the bbcxsr1 is available as a server using ssh or afp. You can connect to it from the finder menu (-> GO -> Connect to Server) pasting the following into the menu box `afp://bbcxsr1.biotech.uconn.edu` (select your account).

**LINUX:** Most editors on linux systems recognize Perl programs and provide context dependent coloring. Ssh and Konquerer work well for file transfer.

# characters at the end of lines

File transfers from Windows to UNIX and return:

End of Line characters are a problem. Under Windows DO NOT use notepad, it does not understand UNIX newline symbols '\n'.

**Best** write your programs under UNIX using vi or vim (or any other editor you are comfortable with)

**2nd** best is to use a text editor like [textwrangler](#) (very nice and free program for UNIX). Like vi and vim it provides context dependent coloring.

**3rd** best is to remove end of line symbols in a UNIX editor or use sed (Stream Editor) after you transferred the file:

```
sed s/.$// name_of_WINDOWS_infile > name_of_UNIX_outfile
```

(This replaces the last non letter character before the eol (\$) with nothing)

Some versions of office allow to change files as UNIX textfiles, but ...

A related problem is encountered by Mac users. Most text editors will use MAC carriage returns at the end of the line. Most unix programs will not be able to handle these. In a terminal window you could use the following command to convert your file:

```
tr '\r' '\n' < name_of_the_Mac_file > name_of_the_unix_file
```

If you are working in a GUI environment, you also could use the convertNewLines.app program (install it in your application folder, drag the file you want to convert into the icon). The program is available [here](#). This is very inconvenient, but there really is no easy solution, tough luck; and you better know about this incase something goes wrong.

# vi

A vi tutorial is at <http://www.eng.hawaii.edu/Tutor/vi.html> -- however, if you run into problems google usually helps (e.g. google: vi replace unix gives you many pages of info on how to replace one string with another under vi)

```
vi myprogram.pl #starts the editor and loads the file myprogram.pl into the editor
```

The following should get you started:

The arrow keys move the cursor in the text

(if you have a really dumb terminal you can use the letter hjkl to move the cursor)

`x` deletes the character under the cursor  
`esc` (i.e. the escape key) leaves the edit mode  
`i` enters the edit mode and inserts before the cursor  
`a` enters the edit mode and appends

`esc` : opens a command line (here you can start searches, and replacements)

`:w` #saves the file

`:w new_name_of_file` #writes the file into a new file.

`:wq` #saves the file and exits vi

`:q!` #exits vi without saving

# customizing vi

One of the beauties of vi is that usually it provides context dependent coloring.

You need to tell vi which terminal you use.

One way to do so is to add a file called `.vimrc` to your home directory.

The following works under both, MAS OSX and using ssh via the secure shell program under windows:

```
vi .vimrc #opens vi to edit .vimrc (Files that start with a dot are not listed if you list a directory. List with ls -a)
```

```
set term=xterm-color #tells the editor that you use a terminal that conforms to some standard
```

```
syn on # tells the editor program that you want to use syntax dependent coloring.
```

```
esc:wq
```

This might seem a little inconvenient, but it really comes in handy to trouble shoot the program in the same environment where you want to run it.

(comment on textwrangler alternative, ssh is included inside the program)

# PERL conventions and rules

Basic Perl Punctuation:

line ends with “;”

empty lines in program are ignored

comments start with #

first line points to path to interpreter:

```
#! /usr/bin/perl
```

# “#!” is known as “shebang”;

keep one command per line for readability

For shell scripts the first line would refer to the type of shell to be used, e.g.:

```
#!/bin/bash
```

use indentation do show program blocks.

Variables start with **\$**scalars, **@**rarrays, or **%**ashes

**Scalars:** floating point numbers, integers,  
non decimal integers, strings

Scalar variables are placeholders that can be assigned a scalar value (either number or string).

## Scalar variables begin with \$

```
$n=3; #assigns the numerical value 3 to the variable $n.  
#Variables are interpolated, for example if you print text
```

```
$b = 4 + ($a = 3); # assign 3 to $a, then add 4 to that  
# resulting in $b getting 7  
$d = ($c = 5); # copy 5 into $c, and then also into $d  
$d = $c = 5; # the same thing without parentheses
```

```
$a = $a + 5; # without the binary assignment operator  
$a += 5; # with the binary assignment operator
```

```
$str = $str . " "; # append a space to $str  
$str .= " "; # same thing with assignment operator
```

```
"hello" . "world" # same as "helloworld"  
'hello world' . "\n" # same as "hello world\n"  
"fred" . " " . "barney" # same as "fred barney"  
"fred" x 3 # is "fredfredfred"  
"barney" x (4+1) # is "barney" x 5, or # "barneybarney....."  
(3+2) x 4 # is 5 x 4, or really "5" x 4, which is "5555"
```

Note: these are not mathematical equations but assignments!

# Numbers can be manipulated using the typical symbols:

$2 + 3$  # 2 plus 3, or 5

$5.1 - 2.4$  # 5.1 minus 2.4, or approximately 2.7;

$3 * 12$  # 3 times 12 = 36;

$2^{**}3$  # 2 taken to the third power =  $2*2*2 = 8$

$14 / 2$  # 14 divided by 2, or 7;

$10.2 / 0.3$  # 10.2 divided by 0.3, or approximately 34;

$10 / 3$  # always floating point divide, so approximately 3.3333333...

## Special characters:

`\n #newline`

`\t #tab`



Double quoted strings are interpolated by the Perl interpreter:

```
"hello world\n" # hello world, and a newline  
"coke\tsprite" # a coke, a tab, and a sprite
```

The backslash can precede many different characters to mean different things (typically called a backslash escape).

# Variable interpolation - single quoted strings are not interpolated:

```
'hello' # five characters: h, e, l, l, o
'don\'t' # five characters: d, o, n, single-quote, t
'' # the null string (no characters)
'silly\\me' # silly, followed by backslash, followed by me
'hello\n' # hello followed by backslash followed by n
'hello
there' # hello, newline, there (11 characters total)
```

# Assignments for next week:

Think about a topic for your student project!  
Please, don't hesitate to send me an email in case you have a question.

Let me know what you are interested in (email).  
What we will do in this course will in part depend on your interests.

# Assignment for next Monday

- 1) On the computer that you plan to use for your project set up a connection (or connections) to `bbcxsrv1` and `bbcsrcv3` that allows you
  - (a) ssh to the server using a command line interface
  - (b) allows you to drop and drag files from your computer to the server.
- 2) check that your vi editor on `bbcxsrv1` is set up to have context dependent coloring (do this, even if you don't plan to use vi on the server!).
- 3) if you do not want to use vi, install an editor on your computer that provides context dependent coloring.

[4] Read through U1-U26 of the Unix and Perl Primer for Biologists (available at the Korfflab at [http://korfflab.ucdavis.edu/Unix\\_and\\_Perl/](http://korfflab.ucdavis.edu/Unix_and_Perl/)]

Read through pages 53-61 of the Unix and Perl Primer for Biologists

5) Read U35 and <http://kb.iu.edu/data/abdb.html> on the `chmod` command in unix

6) Create first Perl Program- "Hello, world!" [make file executable using `chmod`]

```
#!/usr/bin/perl -w  
print ("Hello, world! \n");
```

What happens if you leave out the new line character?

You can run the program by typing `./program_name.pl`, if the file containing the program is made executable (using `chmod u+x *.pl`).