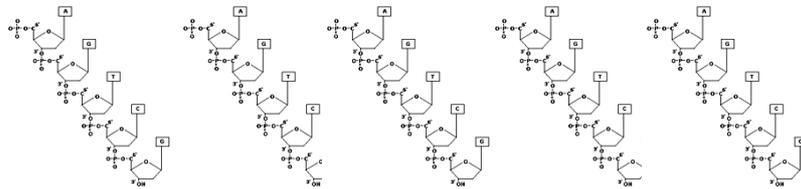


Molecular Phylogeny and Evolution

10 – 14 February 2020

Juan I. MONTOYA BURGOS



Lab of Molecular Phylogeny and Evolution in Vertebrates

Title of the course: **Molecular Phylogeny and Molecular Evolution**

Evolutionary relationships among organisms = tree topology

First phylogenetic methods did not make use of models of molecular evolution (UPGMA, Maximum Parsimony)

A better understanding of molecular evolution improves:

- topology and branch length reconstruction (=phylogenetic tree)

Better phylogenetic trees improve:

- the understanding of evolutionary processes

=> Models of molecular evolution

Lab of Molecular Phylogeny and Evolution in Vertebrates

Why should we care about phylogenies?

Are you using phylogenetics ?

Current phylogenetic methods allow:

- reconstruction of evolutionary relationships

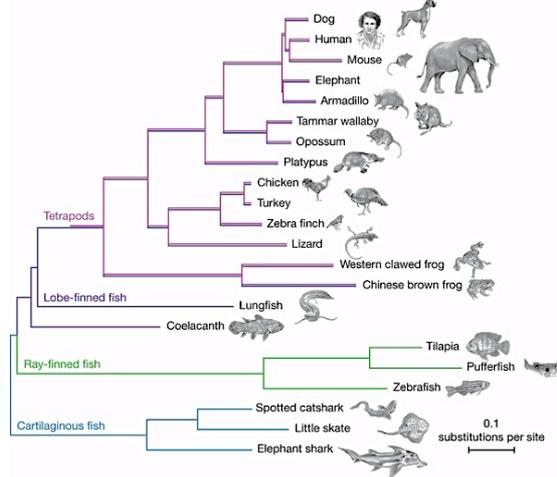
But also the analysis of:

Gene/genome duplication Recombination
Evolutionary rates Divergence time among lineages
Selective pressure Demography
Genetic variability Conservation Biodiversity
Biogeography Spread of contagious disease
Adaptive evolution Discovery of biomedical compounds
Biomonitoring Protein-protein co-evolution

And more

Understanding evolutionary relationships

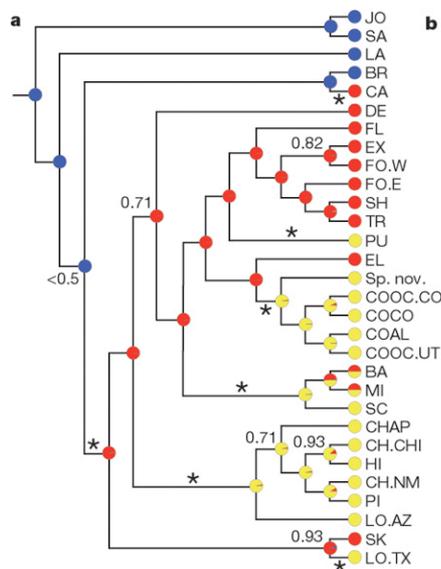
Determining the closest extant relative to tetrapods



Amemiya et al., 2013. The African coelacanth genome provides insights into tetrapod evolution. *Nature* 496.

Lab of Molecular Phylogeny and Evolution in Vertebrates

Understanding morphological evolution



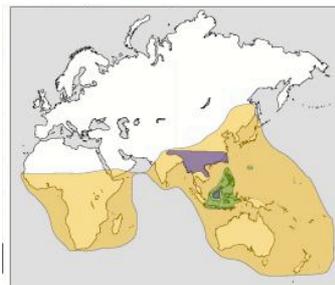
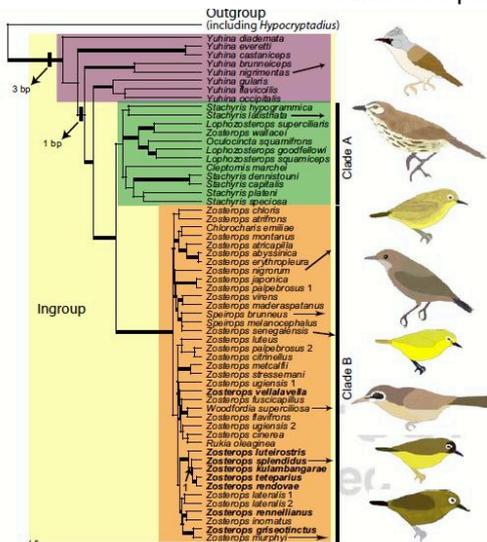
Using a species-level phylogeny of *Aquilegia*, Whittall and Hodges showed a significant evolutionary trend for **longer nectar spur** during directional shifts to **pollinators with longer tongues** (moths).

Whittall and Hodges, 2007. Pollinator shifts drive increasingly long nectar spurs in columbine flowers. *Nature* 447, 706-709

Lab of Molecular Phylogeny and Evolution in Vertebrates

Understanding bio-diversification

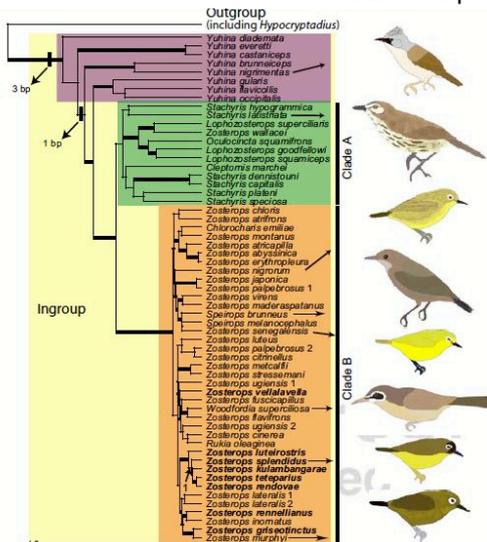
Calibrated phylogeny of *Zosterops* and relatives



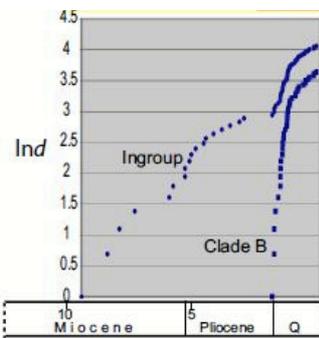
Moyle et al. 2009. Explosive Pleistocene diversification and hemispheric expansion of a "great speciator." PNAS

Understanding bio-diversification

Calibrated phylogeny of *Zosterops* and relatives



Lineage-through-time plots (LTT)

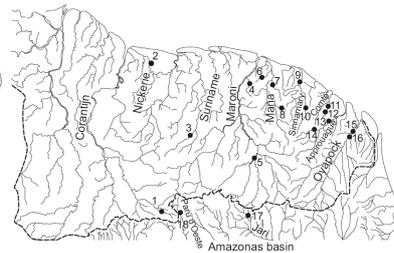


Moyle et al. 2009. Explosive Pleistocene diversification and hemispheric expansion of a "great speciator." PNAS

New biodiversity

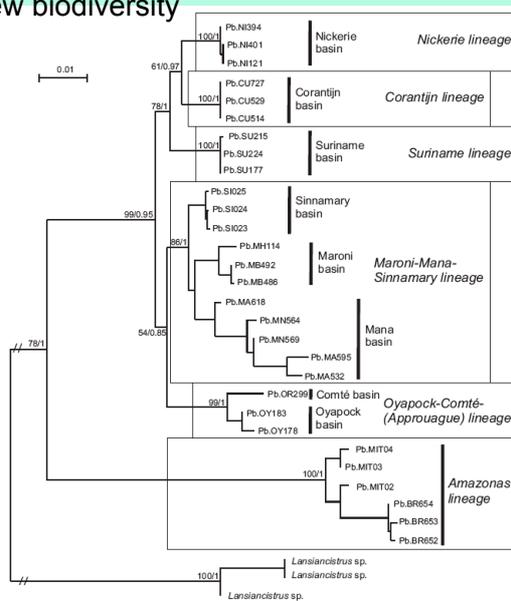
Unexpected diversity in the catfish *Pseudancistrus brevispinis* reveals dispersal routes in a Neotropical center of endemism: the Guyanas Region

YAMILA P. CARDOSO¹ and JUAN I. MONTOYA-BURGOS^{1*}



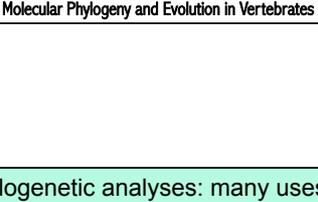
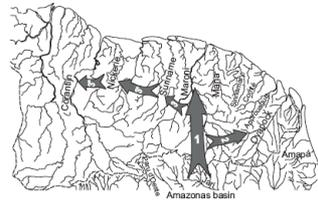
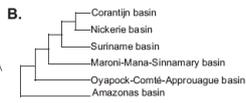
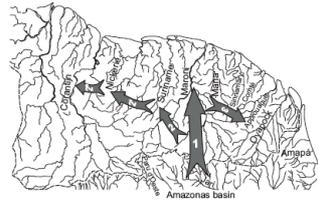
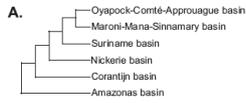
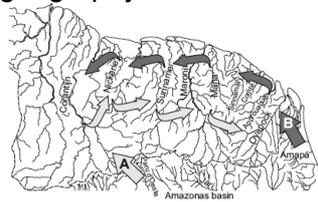
New biodiversity

1 species → 6 species



Morphology		
Plates between adipose and caudal fins	Head profile	Body shape
Three small, one large	Ascending from eyes to dorsal fin	Intermediate
Four small	Rounded from eyes to dorsal fin	Massif
Three small, one large	Horizontal and straight from eyes to dorsal fin	Widest
Two small, two large	Shortest and pointed snout. Smallest interorbital distance	Highest
Three small, one large	Ascending from eyes to dorsal fin	Longest and thinnest
Three small, two large	Descending straight from eyes to dorsal fin	Flattest

Phylogenetic analyses: many uses
Biogeography

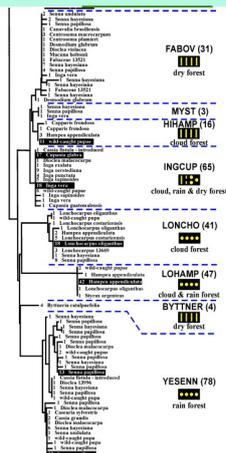
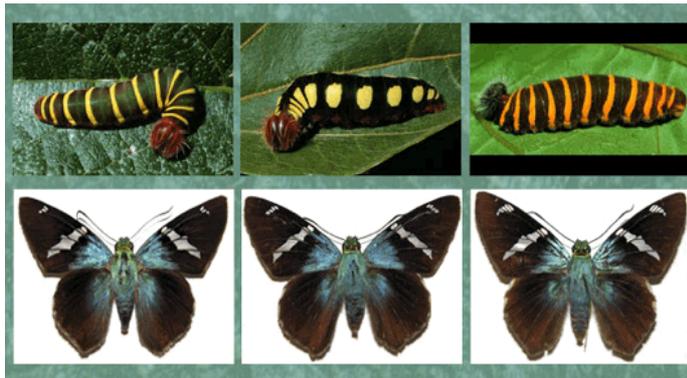


Biogeography:
four possible
entrance scenarios

Topological tests:
-SH test
-AU test



Phylogenetic analyses: many uses
DNA taxonomy of Animals: species identification



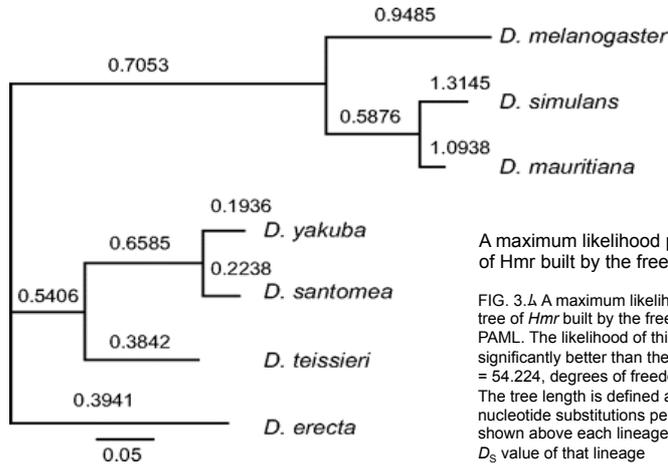
Ten species in one: **DNA barcoding** reveals cryptic species in the neotropical skipper butterfly *Astrapttes fulgerator*

Paul D. N. Hebert^{*,†}, Erin H. Penton^{*}, John M. Burns[‡], Daniel H. Janzen[§], and Winnie Hallwachs[¶]

Hebert et al. PNAS 101:14812 (2004)

Evolutionary rates and selective pressure acting on genes

The evolutionary histories of genes bear the marks of the functional demands to which they have been subjected.



A maximum likelihood phylogenetic tree of *Hmr* built by the free-ratio model in PAML

FIG. 3. A maximum likelihood phylogenetic tree of *Hmr* built by the free-ratio model in PAML. The likelihood of this model was significantly better than the one-ratio model ($2l = 54.224$, degrees of freedom = 11, $P < 10^{-5}$). The tree length is defined as the number of nucleotide substitutions per codon. The number shown above each lineage is the estimated d_N/d_S value of that lineage

Maheshwari, S. et al. Mol Biol Evol 2008 25:2421-2430; doi:10.1093/molbev/msn190

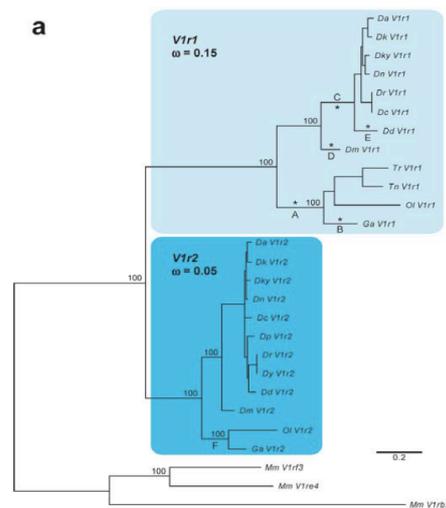
Lab of Molecular Phylogeny and Evolution in Vertebrates



Divergent Evolution among Teleost V1r Receptor Genes

Patrick Pfister^{1,2}, Jerome Randall^{1,2}, Juan I. Montoya-Burgos¹, Ivan Rodriguez^{1,2*}

1 Department of Zoology and Animal Biology, University of Geneva, Geneva, Switzerland. 2 National Center of Competence in Research (NCCR) Frontiers in Genetics, University of Geneva, Geneva, Switzerland



Analysis of:

- duplicated genes
- gene families
- genome duplication / polyploidisation

Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution and selection on proteins, protein interactions and protein networks

REVIEWS

 COMPUTATIONAL TOOLS

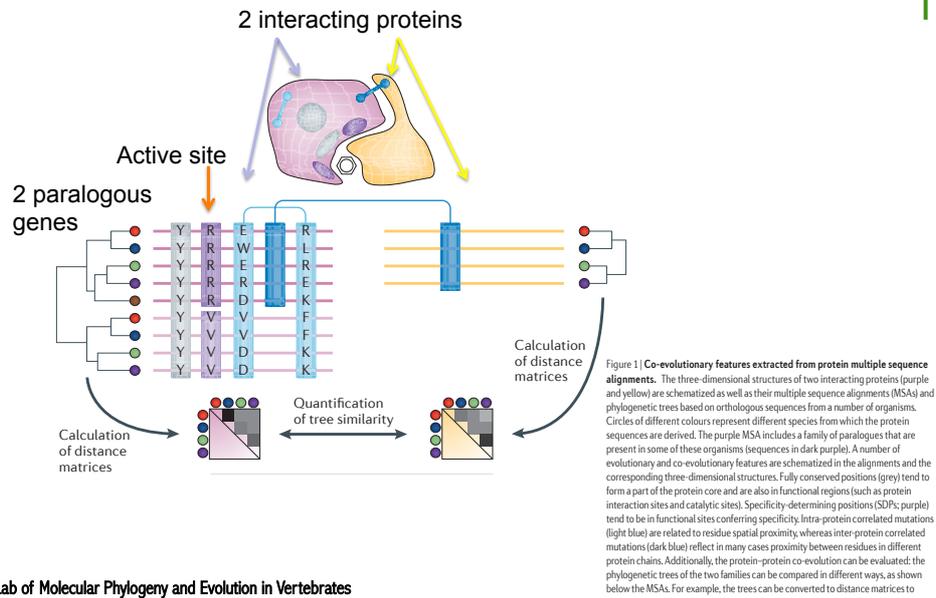
Emerging methods in protein co-evolution

David de Juan¹, Florencio Pazos² and Alfonso Valencia¹

NATURE REVIEWS | GENETICS

VOLUME 14 | APRIL 2013 | 249

Protein evolution



Protein evolution



RESEARCH ARTICLE



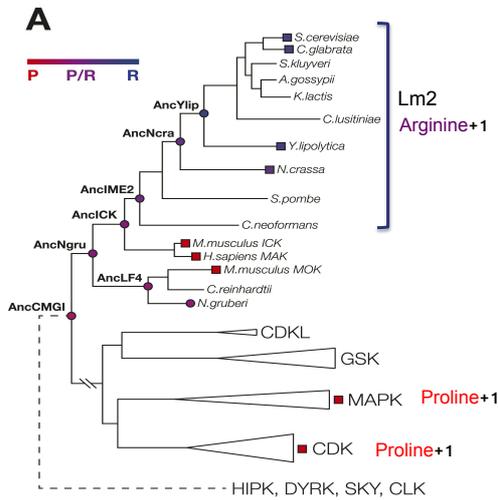
Ancestral resurrection reveals evolutionary mechanisms of kinase plasticity

Conor J Howard^{1†}, Victor Hanson-Smith^{2†}, Kristopher J Kennedy¹, Chad J Miller³, Hua Jane Lou³, Alexander D Johnson², Benjamin E Turk^{3*}, Liam J Holt^{1*}

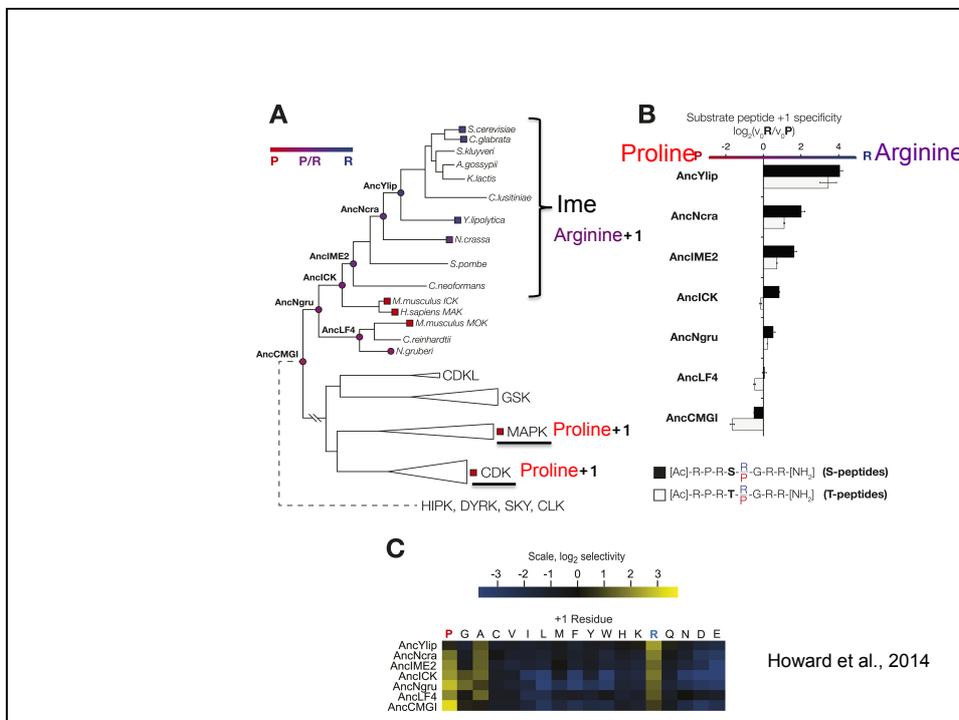
DOI: <http://dx.doi.org/10.7554/eLife.04126>
Cite as eLife 2014;3:e04126

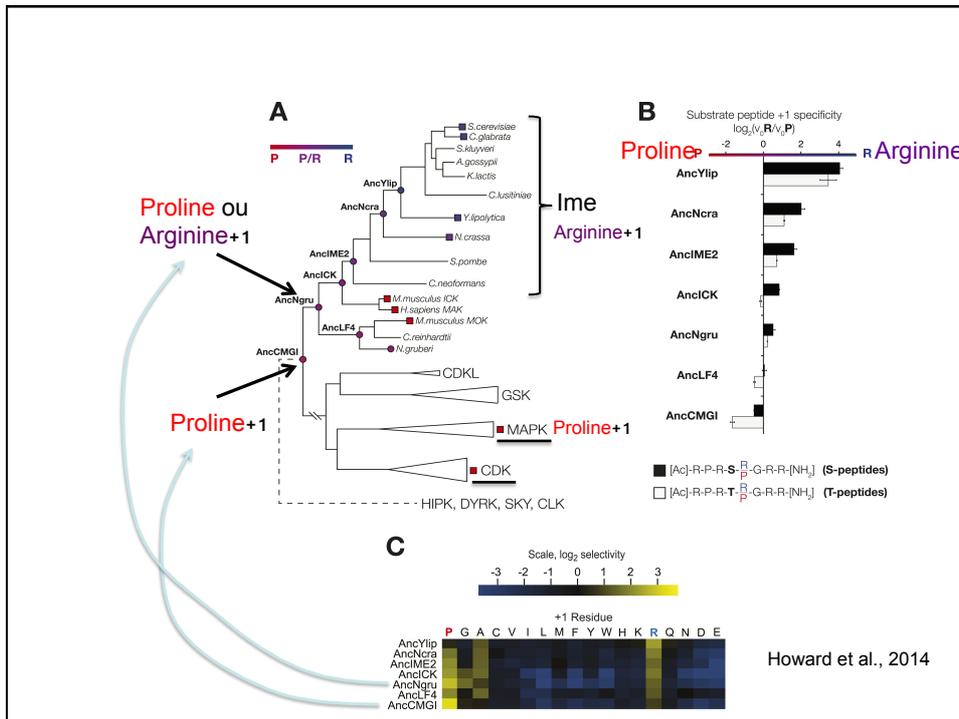
Protein evolution

Ancestral sequence reconstruction using Maximum Likelihood inference.
 Ancestral sequence synthesis and enzymatic activity essays.



The cyclin dependent kinases (CDKs) and mitogen activated protein kinases (MAPKs) require proline (P) at the +1 position of their substrates, while Lme2 prefers arginine (R).

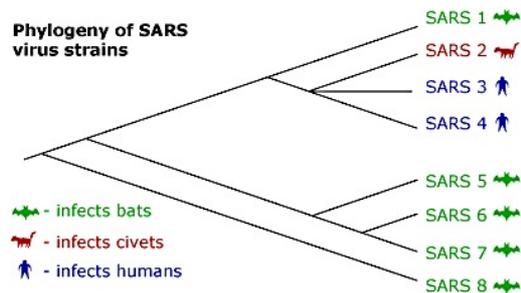




Phylogenetic analyses: many uses

Origin of viruses, disease phylogeny

In 2002-2003, the previously unknown SARS virus generated widespread panic causing 774 deaths and more than 8000 cases of illness. Where did this virus come from?



Source of Severe Acute Respiratory Syndrome coronavirus (SARS) was unknown... phylogenetic analyses traced it to civets and bats.

Phylogenetics as a tool in various fields

Phylogenetic is used as a tool:

- in historical linguistics,
- in security applications for networks and computers such as **artificial immune systems** for computers,
- in forensic studies, for instance in infectious forensics.

Infectious forensics

Phylogenetics offers a way to establish relationships between microbes infecting several individuals and can be used as corroborating evidence when someone is suspected of infecting others with a disease.

- 1 Pathogen genomes can mutate quickly, creating diverse microbial populations in those infected.
- 2 By sequencing highly variable regions of pathogen genomes scientists can build a phylogenetic tree that suggests how the microbes are related.
- 3 The relatedness of the viral populations can support or rule out hypotheses of who infected whom.
- 4 Pathogen diversity can also be used to corroborate time of infection using the mutation rate as a molecular clock.

<http://www.nature.com/news/506424a-i2-jpg-7.15735?article=1.14775>

Lab of Molecular Phylogeny and Evolution in Vertebrates

Today's theory refreshment

Molecular evolution

Molecular evolution

- Molecular basis of evolution.

- Evolution of protein sequences
- Evolution of DNA sequences

Molecular evolution

Molecular basis of evolution

- 1) Evolutionary processes
- 2) Genome structure
- 3) Gene structure and functions
- 4) Mutational changes in DNA
- 5) Codon Usage in coding sequences

1) Evolutionary processes

Causal factor of evolution:

mutational changes in DNA, particularly in coding sequences or regulatory regions

Main mutational changes:

- nucleotide substitution
- insertion/deletion of nucleotides (indels) / mobile genetic elements
- non-homologous recombination
- gene conversion
- duplications of genes / chromosomal region / genomes

Mutation Fixation in a species:

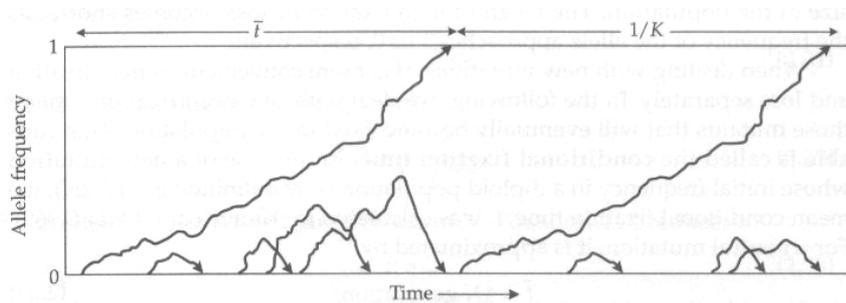
A new mutation can spread in a population over generations and can eventually be fixed (or lost) in the species by genetic drift and/or by natural selection.

- Natural selection: ... - positive Darwinian selection promotes fixation
- negative selection (purifying selection) leads to mutation loss

1) Evolutionary processes

Fixation (or loss) of a mutation by genetic drift or natural selection.

A mutation is **fixed** when it is present in **all members of a species**.



1) Evolutionary processes

Fixation of a mutation

Neutral mutations

A fraction μ of all copies of a gene mutate (μ = mutation rate).

Of these $1/2N$ (equal to the initial frequency of the mutant) succeed in drifting to fixation for the mutant (=probability of fixation).

There are in all $2N$ copies of the gene available to mutate.

The resulting rate of substitution is

$$r = \mu \times 2N \times 1/2N = \mu$$

So the rate of substitution of neutral mutations is equal to the mutation rate.

1) Evolutionary processes

Fixation of a mutation

Selected mutations

A fraction μ of copies of the gene mutate. There are in all $2N$ copies available.

A fraction $2s$ succeed in fixing (=probability of fixation for a selected mutation).

The resulting rate of substitution is:

$$r = \mu \times 2N \times 2s = 4Ns\mu$$

s = selective coefficient.

Note that this is $4Ns$ times as high as for neutral mutants, if the mutation rate in both categories were equal.

1) Evolutionary processes

Fixation of a mutation

Selected mutations

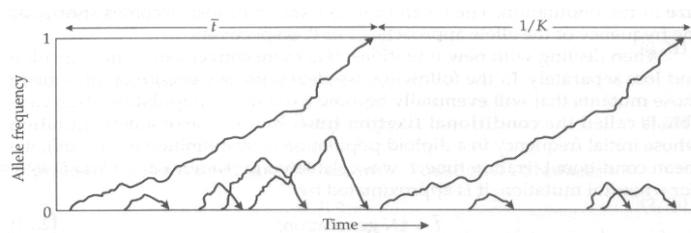
Fixation rate (r) depends on:

- population size (N),
- mutation fitness (s)
- mutation rate (μ).

$$r = 4Ns\mu$$

1) Evolutionary processes

The fitness (s) of a mutation influences the fixation rate



fitness (s)

time to fixation

Neutral mutation:

t_{neu}

Positively selected mutation:

$t_{pos} < t_{neu}$

Negatively selected mutation:

$t_{neg} > t_{neu}$, or $t_{neg} = \infty$

Molecular evolution

Molecular basis of evolution

- 1) Evolutionary processes
- 2) Genome structure
- 3) Gene structure and fonctions
- 4) Mutational changes in DNA
- 5) Codon Usage in coding sequences

2) Genome Structure

Coding and non-coding regions

Non-coding regions (3):

- intergenic DNA or “junk” DNA
- regulatory regions (promoters, enhancers, ...)
- non-coding repetitive elements (Alu, SINE, telomeric repeats)

Coding regions (4): - protein-coding genes with messenger RNAs (mRNAs)

- genes coding for structural RNAs:
 - ribosomal RNAs (rRNA),
 - transfer RNAs (tRNA),
 - small nuclear RNAs (snRNA)
- genes coding for small regulatory RNAs: short interfering RNAs, microRNAs,...
- genes coding for long non-coding RNAs (> 200 bp) : function?

Molecular evolution

Molecular basis of evolution

- 1) Evolutionary processes
- 2) Genome structure
- 3) Gene structure and functions
- 4) Mutational changes in DNA
- 5) Codon Usage in coding sequences

2) Gene structure and function

Structure of protein-coding genes

Transcribed regions: succession of exons and introns
(except: mono-exonic genes)

- Untranslated 5' region (= 5' UTR) in the first (or few first) exon(s).
- Start codon (AUG) usually in the first exon.
- Stop codon usually in the last exon.
- Untranslated 3' region (= 3' UTR) the last (or few last) exon(s).

The pre-mRNA still contains introns.

The **splicing** process discards introns, which leads to a messenger RNA (mRNA)

mRNA still contains **5' and 3' UTRs** in addition to the **coding sequence (CDS)**

The CDS starts at the **Start codon** (AUG) and ends at the codon just before the **Stop codon**.

2) Gene structure and function

The genetic code

Every triplet of the CDS encodes for an amino acid = genetic code

The genetic code is **universal** (or almost universal).
It is essentially the same in Prokaryotes, Eukaryotes, chloroplasts and in mitochondria.

The genetic code is ... **degenerated**.
There are 64 possible codons (4^3), yet only 20 amino acids. Some codons encode for the same amino acid, they are **synonymous**.

Start and Stop codons:

The standard Start codon is AUG which encodes a methionin (Met).
When several methionins are candidates, the "Kozack methionin" is the Start codon, also referred to as **the Translation Initiation Site (TIS)**.
Usually, there are three termination codons (STOP codons): UAA, UAG et UGA.

2) Gene structure and function

universal genetic code

		2nd base in codon				
		U	C	A	G	
1st base in codon	U	Phe Phe Leu Leu	Ser Ser Ser Ser	Tyr Tyr STOP STOP	Cys Cys STOP Trp	U C A G
	C	Leu Leu Leu Leu	Pro Pro Pro Pro	His His Gln Gln	Arg Arg Arg Arg	U C A G
	A	Ile Ile Ile Met	Thr Thr Thr Thr	Asn Asn Lys Lys	Ser Ser Arg Arg	U C A G
	G	Val Val Val Val	Ala Ala Ala Ala	Asp Asp Glu Glu	Gly Gly Gly Gly	U C A G
						3rd base in codon

2) Gene structure and function

mitochondrial genetic codes

Vertebrate mtDNA.

Codon	New translation
AGA	Stop
AGG	Stop
ATA	Met
TGA	Trp

Invertebrate mtDNA.

Codon	New translation
AGA	Ser
AGG	Ser
ATA	Met
TGA	Trp

Yeast mtDNA.

Codon	New translation
ATA	Met
CTA	Thr
CTC	Thr
CTG	Thr
CTT	Thr
TGA	Trp

2) Gene structure and function

mitochondrial genetic codes

Mold, Protozoan and Coelenterate mtDNA.

Codon	New translation
TGA	Trp

Echinoderm mtDNA.

Codon	New translation
AAA	Asn
AGA	Ser
AGG	Ser
TGA	Trp

Ascidian mtDNA.

Codon	New translation
AGA	Gly
AGG	Gly
AGG	Met
TGA	Trp

Flatworm mtDNA.

Codon	New translation
AAA	Asn
AGA	Ser
AGG	Ser
TAA	Tyr
TGA	Trp

2) Gene structure and function

Different genetic codes in the **nucleus**

Ciliate Nuclear Code.

Codon	New translation
TAA	Gln
TAG	Gln

Alternative Yeast Nuclear.

Codon	New translation
CTG	Ser

Molecular evolution

Molecular basis of evolution

- 1) Evolutionary processes
- 2) Genome structure
- 3) Gene structure and functions
- 4) Mutational changes in DNA
- 5) Codon Usage in coding sequences

3) Mutational changes in DNA



Morphological and physiological heritable characters are encoded in the DNA.

➔ Every heritable character change is caused by a change in the DNA.

Four types of changes at the nucleotide level:

1. Substitution of a nucleotide by another nucleotide
2. Insertion
3. Deletion
4. Inversion

Insertions, deletions and inversions may involve more than a single nucleotide.

3) Mutational changes in DNA

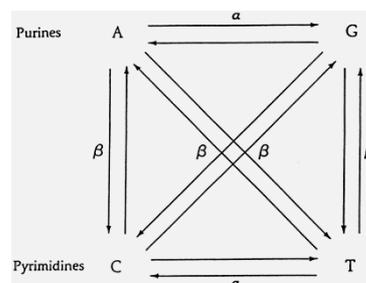
Two types of substitutions: **transitions** and **transversions**

Transitions: substitution between nucleotides belonging to the same family:

a change from a purin (A, G) to a purin
or
from a pyrimidin (C, T) to a pyrimidin (C, T)

Transversions: substitution between nucleotides belonging to different families.

Substitution scheme



3) Mutational changes in DNA

Consequences at the level of coding regions:

Due to the degenerated genetic code, the impact of a substitution will vary:

A will not change the encoded amino acid.

A will change the encoded amino acid.

A will generate a STOP codon.

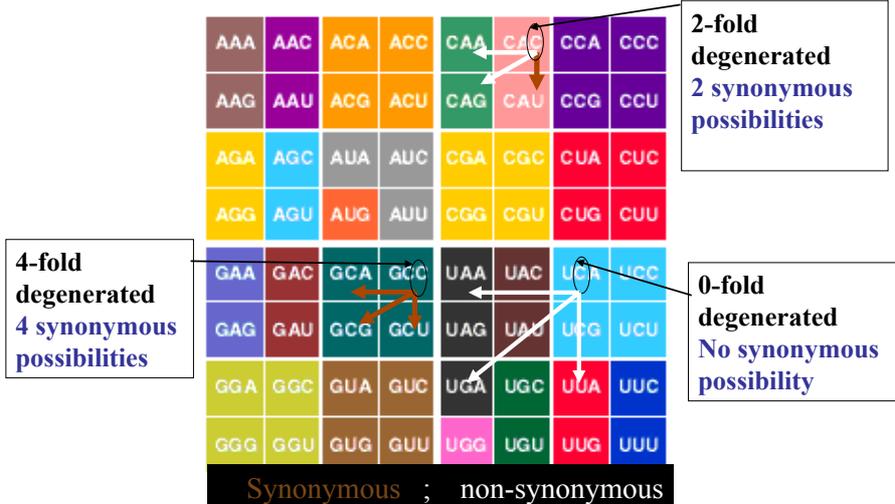
Due to the degenerated genetic code:

- most of the **synonymous** mutations take place at the **3rd** codon position,
- a few **synonymous** mutations take place at the **1st** codon position.

Mutations occurring at the 2nd codon position will be either **non-synonymous** or **nonsense** mutations.

3) Mutational changes in DNA

Sites are classified into 3 categories according to their degeneration level



3) Mutational changes in DNA

4-fold degenerated: 32 possibilities in 3rd codon position over 61 codons

AAA	AAC	ACA	ACC	CAA	CAC	CCA	CCC
AAG	AAU	ACG	ACU	CAG	CAU	CCG	CCU
AGA	AGC	AUA	AUC	CCA	CCU	CUA	CUC
AGG	AGU	AUG	AUU	CCG	CCU	CUG	CUU
GAA	GAC	GCA	GCC	UAA	UAC	UCA	UCC
GAG	GAU	GCG	GCU	UAG	UAU	UCG	UCU
GGA	GGC	GUA	GUC	UGA	UGC	UUA	UUC
GGG	GGU	GUG	GUU	UGG	UGU	UUG	UUU

Lab of Molecular Phylogeny and Evolution in Vertebrates

3) Mutational changes in DNA

2-fold degenerated: 25 possibilities in 3rd codon position
8 possibilities in 1st codon position.

AAA	AAC	ACA	ACC	CAA	CAC	CCA	CCC
AAG	AAU	ACG	ACU	CAG	CAU	CCG	CCU
AGA	AGC	AUA	AUC	CCA	CCU	CUA	CUC
AGG	AGU	AUG	AUU	CCG	CCU	CUG	CUU
GAA	GAC	GCA	GCC	UAA	UAC	UCA	UCC
GAG	GAU	GCG	GCU	UAG	UAU	UCG	UCU
GGA	GGC	GUA	GUC	UGA	UGC	UUA	UUC
GGG	GGU	GUG	GUU	UGG	UGU	UUG	UUU

Lab of Molecular Phylogeny and Evolution in Vertebrates

3) Mutational changes in DNA

Consequences at the level of coding regions:

Theoretical proportions of the three substitution types.

If we assume that:

- codons are used with the same frequency
- nucleotides can be substituted by any other nucleotide with the same frequency

then:

non-synonymous substitutions =	71%
synonymous substitutions =	25%
non-sense substitutions =	4%

Is this found in real data ?

These theoretical proportions are not found in real data !
(the assumptions are not valid)

3) Mutational changes in DNA

Insertions / deletions : main causes

Small indels (involving few nucleotide sites) are generally due to **DNA replication errors** during meiosis (as for most nucleotide substitutions).

Large indels may result from:

- **unequal crossing overs**
- **transposition** via transposons or other transposable elements (retrotransposons, LINEs, SINEs, Alu, ...).

Some insertions may be due to **horizontal gene transfers**,
= integration into the genome of a species of a piece a DNA coming from a different species (via viruses, plasmids, transposons...?).

Molecular evolution

Molecular basis of evolution

- 1) Evolutionary processes
- 2) Genome structure
- 3) Gene structure and functions
- 4) Mutational changes in DNA
- 5) Codon Usage in coding sequences

4) Codon Usage in Coding Sequences

Are **synonymous codons** used with **equal frequency**?

$UCU = UCC = UCA = UCG$??? (Serine)

If we assume:

- absence of selective pressure
- no bias in nucleotide substitution frequency

then **synonymous** codons should be present at the same frequency across all protein-coding genes in a genome.

Real data:

Synonymous codons are not equally frequent (assumptions above are not realistic).
The bias in synonymous codon usage is highly variable across organisms and genes.

Example: *E.coli* RNA polymerase genes (in parenthesis: *RSCU*)
Relative Synonymous Codon Usage

Val	GUU: 55 (1.53)	Pro	CCU: 9 (0.48)
	GUC: 21 (0.58)		CCC: 0 (0.00)
	GUA: 34 (0.94)		CCA: 11 (0.59)
	GUG: 34 (0.94)		CCG: 55 (2.93)

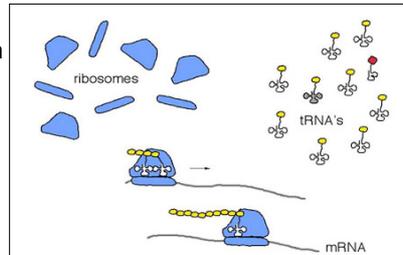
UCU	} Ser
UCC	
UCA	
UCG	
CCU	} Pro
CCG	
CCA	
ACU	} Thr
ACC	
ACA	
ACG	

4) Codon Usage in Coding Sequences

What can explain the bias in synonymous codons usage?

Observations in highly transcribed genes:

- the bias is high
- the usage frequency is high for codons with abundant tRNA in the sites of protein synthesis



Why is it so?

Explanation 1: For genes that need to be **highly transcribed**, alleles that are composed of codons that possess abundant tRNAs can be rapidly transcribed and are favored by natural selection. In other words, **purifying selection** eliminates alleles bearing codons for which tRNAs are rare.

4) Codon Usage in Coding Sequences



HOME | ABC

Search

New Results

Biased gene conversion drives codon usage in human and precludes selection on translation efficiency

Fanny Pouyet, Dominique Mouchiroud, Laurent Duret, Marie Semon

doi: <https://doi.org/10.1101/086447>

This article is a preprint and has not been peer-reviewed [what does this mean?]

Abstract | Info/History | Metrics | Supplementary material | Preview PDF

Abstract

In humans, as in other mammals, synonymous codon usage (SCU) varies widely among genes. In particular, genes involved in cell differentiation or in proliferation display a distinct codon usage, suggesting that SCU is adaptively constrained to optimize translation efficiency in distinct cellular states. However, in mammals, SCU is known to correlate with large-scale fluctuations of GC-content along chromosomes, caused by meiotic recombination, via the non-adaptive process of GC-biased gene conversion (gBGC). To disentangle and to quantify the different factors driving SCU in humans, we analyzed the relationships between functional categories, base composition, recombination, and gene expression. We first demonstrate that SCU is predominantly driven by large-scale variation in GC-content and is not linked to constraints on tRNA abundance, which excludes an effect of translational selection. In agreement with the gBGC model, we show that differences in SCU among functional categories are explained by variation in intragenic recombination rate, which, in turn, is strongly negatively correlated to gene expression levels during meiosis. Our results indicate that variation in SCU among

expression levels during meiosis. Our results indicate that variation in SCU among functional categories (including variation associated to differentiation or proliferation) result from differences in levels of meiotic transcription, which interferes with the formation of crossovers and thereby affects gBGC intensity within genes. Overall, the gBGC model explains 81.3% of the variance in SCU among genes. We argue that the strong heterogeneity of SCU induced by gBGC in mammalian genomes precludes any optimization of the tRNA pool to the demand in codon usage.

4) Codon Usage in Coding Sequences

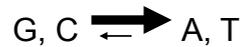
What can explain the bias in synonymous codons usage?

Explanation 2:

Bias in the substitution rate among nucleotides.

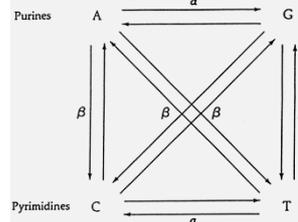
Examples:

1. In some bacteria (*Mycoplasma* sp.), the substitution rate is biased in favor of A or T (Muto at Osawa, 1987)



Consequence on synonymous codon usage?

All 3rd codon positions are A or T, leading to a bias in the use of synonymous codons



4) Codon Usage in Coding Sequences

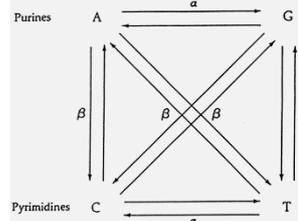
What can explain the bias in synonymous codons usage?

Explanation 2:

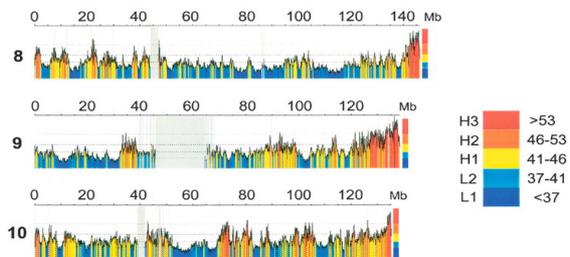
Bias in the substitution rate among nucleotides.

Examples:

2. Isochores in mammals: alternating chromosomal regions with high G+C content (60%) and low G+C content (30%).



The substitution bias can thus vary within a genome!



An isochore map of human chromosomes, Constantini et al., 2006

4) Codon Usage in Coding Sequences

Statistical measure of the synonymous codon usage bias.

The **RSCU** (Relative Synonymous Codon Usage) **characterizes a codon** and gives its usage level as compared to the other synonymous codons in a given gene (or set of genes).

$$RSCU_i = \frac{Obs_i}{Exp_i}$$

where:

Obs_i is the observed number of codons i in the gene

Exp_i is the expected number of codons i in the gene.

Exp_i is obtained as follows:

$$Exp_i = \frac{\sum aa_i}{\sum sym_i} \quad \begin{array}{l} \text{occurrence of the aa encoded by the codon } i \text{ in the gene} \\ \text{number of } \mathbf{synonymous\ codons} \text{ encoding that aa} \end{array}$$

4) Codon Usage in Coding Sequences

Statistical measure of the synonymous codon usage bias.

The **CAI** (Codon Adaptation Index) **characterizes a gene** [Sharp and Li \(1987\)](#).

The CAI uses a set of highly expressed genes as a reference for the usage of synonymous codons. The CAI varies from 0 to 1.

High CAI (0.6 to 1) indicates that a gene “behaves” like a highly expressed gene (i.e. has a similar bias in synonymous codon usage).

This Index is used for:

- estimating the expression level of a gene
- comparing the synonymous codon usage bias among organisms

$$\ln(\text{CAI}) = \sum_{i=1}^{61} f_i \ln W_i$$

f_i is the relative frequency of codon i in the analyzed sequence

W_i is the ratio of the frequency of the codon i over the frequency of the most frequently used synonymous codon in the set of highly expressed genes.

4) Codon Usage in Coding Sequences

CAI (Codon Adaptation Index) [Sharp and Li \(1987\)](#)

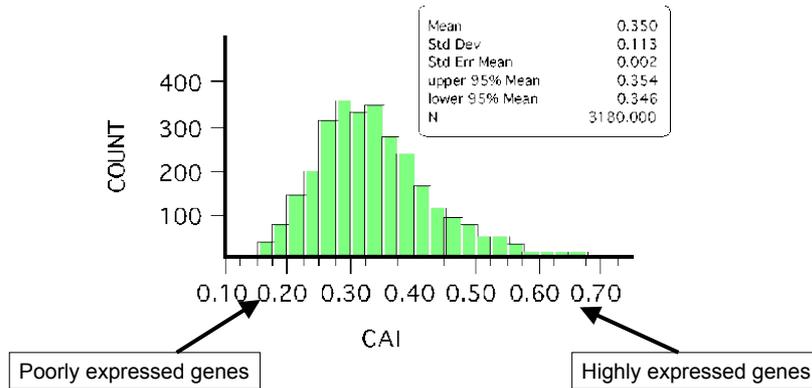


Figure: Distribution of CAI for 3180 *E. coli* genes

Molecular evolution

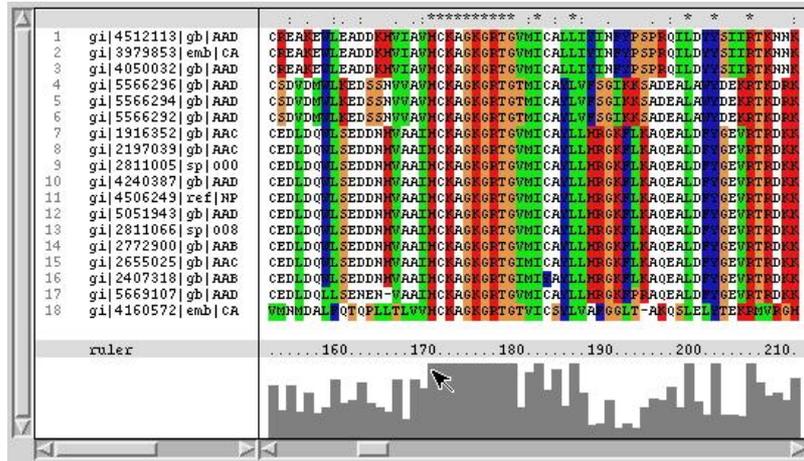
- Molecular basis of evolution.

- Evolution of protein sequences

- Evolution of DNA sequences

Evolution of protein sequences

How to measure the difference between two sequences ?



ClustalX window

Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of protein sequences

How to measure the difference between two sequences ?

A simple measure is to count the number of different amino acids and divide it by the total number of positions compared = **p distance**

	1	5	10	15	20	25	30																						
<i>α-DTX</i>	E	P	R	R	K	L	C	I	L	H	R	N	P	G	R	C	Y	D	K	I	P	A	F	Y	Y	N	K	K	K
<i>DTX-I</i>	Q	P	L	R	K	L	C	I	L	H	R	N	P	G	R	C	Y	Q	K	I	P	A	F	Y	Y	N	K	K	K

$$p = n_d / n$$

$$p = 3/30 = 0.1$$

n_d – number of amino acid differences between two sequences; n – number of aligned amino acids.

The **p distance** is an exact measure of divergence if:

- all the changes that occurred during the independent evolution of the two sequences can be observed today.

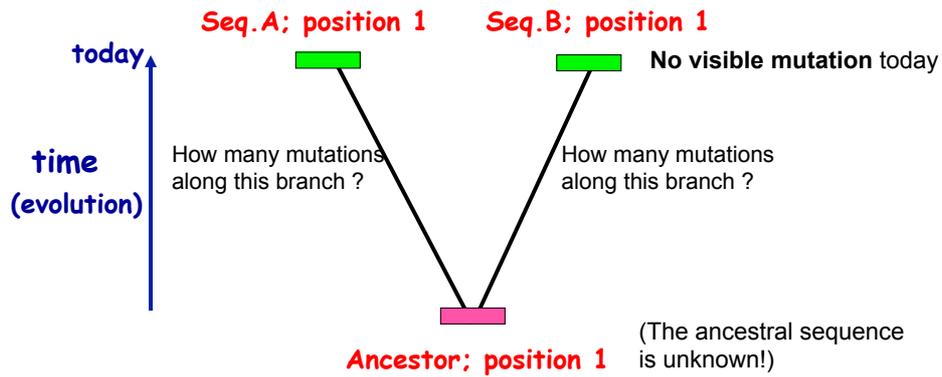
Is this always the case ???

Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of protein sequences

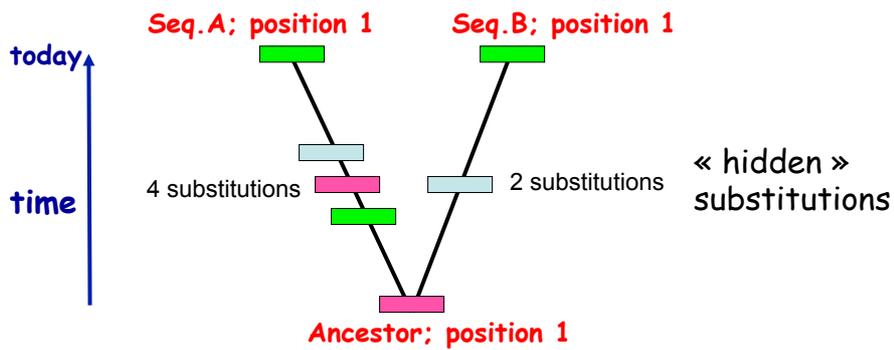
How to measure the difference between two sequences ?

Are all mutations visible today ??



Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of protein sequences



The number of hidden substitutions increases with time.

➔ The p distance is not proportional to time.

Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of protein sequences

The ***p* distance** must be “corrected” to give an accurate measure of divergence.

How to correct the ***p* distance**?

Two ways:

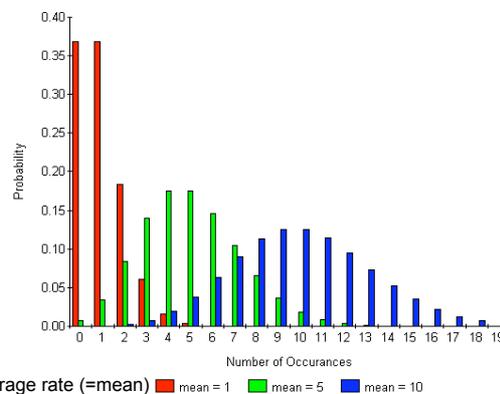
- Theoretical correction: Poisson corrected distance (PC)
- Empirical correction: substitution matrices (PAM, JTT,...)

Evolution of protein sequences

Theoretical correction: Poisson corrected distance (PC)

The Poisson distribution is a discrete probability distribution.

It expresses the probability of a number of events occurring in a fixed period of time if these events occur with a **known average rate** (λ) and independently of the time since the last event.



Evolution of protein sequences

Theoretical correction: Poisson corrected distance (PC)

The probability of having exactly k events ($k = 0, 1, 2, 3, \dots$) is given by:

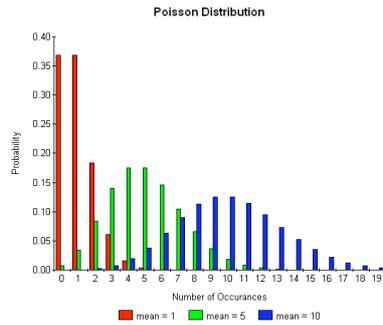
$$f(k; \lambda) = \frac{e^{-\lambda} \lambda^k}{k!},$$

where:

e is the base of the natural logarithm ($e=2.71828\dots$)

$k!$ is the factorial of k

λ is the average number of occurrences that occur during the given time interval

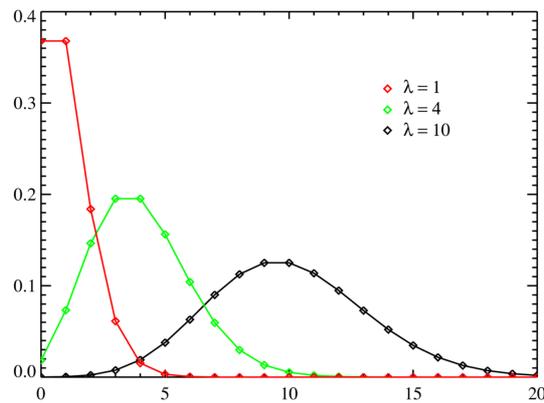


λ is the only variable that determines the distribution.

Evolution of protein sequences

Poisson distribution as a function of λ .

(λ = average number of occurrences that occurs during the given time interval)



Evolution of protein sequences

Theoretical correction: Poisson corrected distance (PC)

For us, what is the value of λ ?

(λ = mean number of occurrences that occur during the given interval)

Let's suppose that r is the substitution rate of any amino acid to any another amino acid during a year, and that this rate is the same for all sites

r = mean number of events per year

The mean substitution rate per site and for a period of time t years is: rt

rt = mean number of events over a period of t years

For us, λ is thus equal to rt

Evolution of protein sequences

Theoretical correction: Poisson corrected distance (PC)

$$f(k; \lambda) = \frac{e^{-\lambda} \lambda^k}{k!},$$

The probability of no change ($k = 0$) along time t for a given site is:

$$P(0; t) = f(0; rt) = e^{-rt} \quad (\text{remember that } n^0 = 1 \text{ and that } 0! = 1)$$



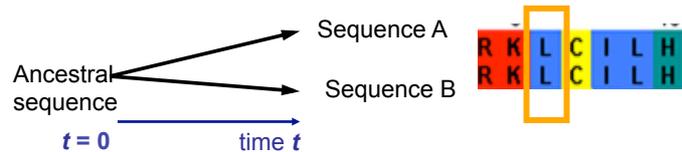
This equation is not useful because the ancestral sequence at $t = 0$ is unknown.

Solution: the proportion of **unchanged** amino acids **between two sequences** along time t is empirically obtained.

Evolution of protein sequences

Theoretical correction: Poisson corrected distance (PC)

The amount of unchanged amino acids is calculated by comparing two sequences that diverged since a time t .



The probability (q) that none of the two aa of an homologous site in sequences A and B has been substituted along time t is given by:

$$q = (e^{-rt}) \times (e^{-rt}) = e^{-2rt}$$

The probability q can be estimated by $q = 1 - p$

$p = p \text{ distance}$ = proportion of observed differences between sequences A and B

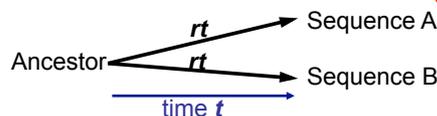
Evolution of protein sequences

Theoretical correction: Poisson corrected distance (PC)

Using equation $q = 1 - p$
 and $q = e^{-2rt}$

we can derive $1 - p = e^{-2rt}$ and by extracting the exponent $2rt = -\ln(1-p)$

Also, the "true" distance separating sequences A and B is: $d = 2rt$
 ($r = \text{substitution rate}$)



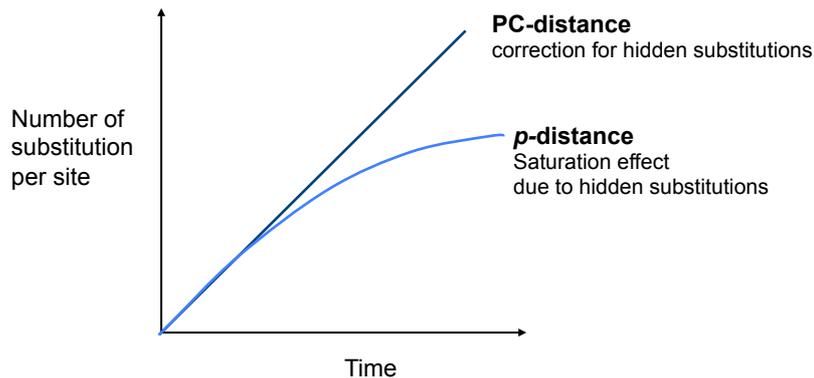
Thus, $d = -\ln(1-p)$ This is the Poisson corrected distance !

The variance is: $V(d) = p / [(1-p)n]$

Evolution of protein sequences

Theoretical correction: Poisson corrected distance (PC)

Relationship between **p-distance** and **Poisson corrected-distance (PC)**



Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of protein sequences

The **Poisson-corrected distance** method assumes that:

-the substitution probability is equal across all sites of the polypeptide sequence.

Is this assumption generally correct ? **NO**

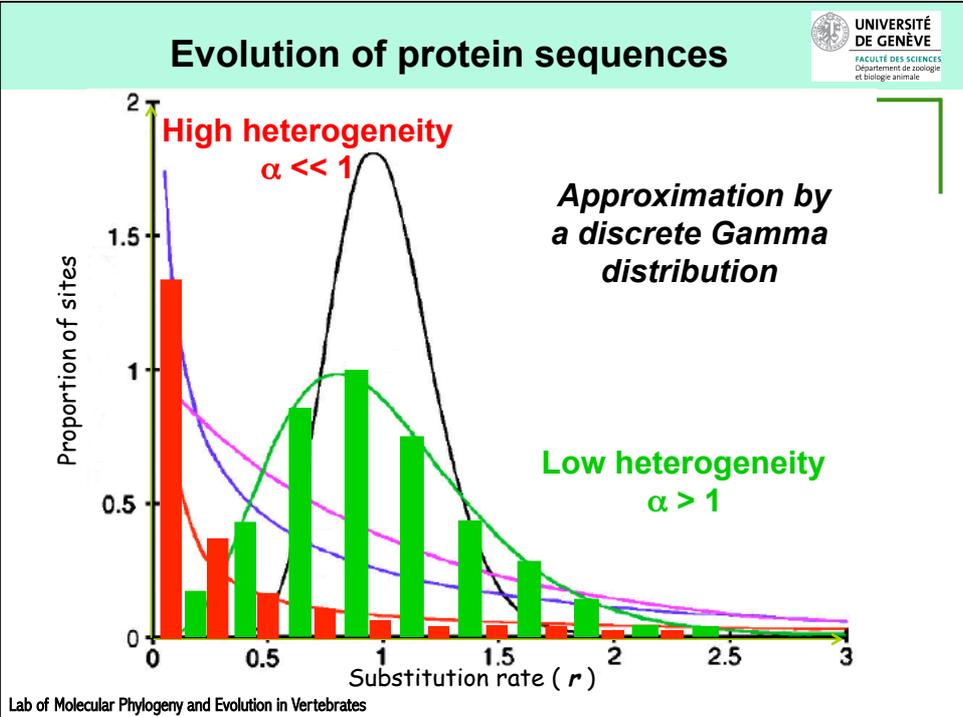
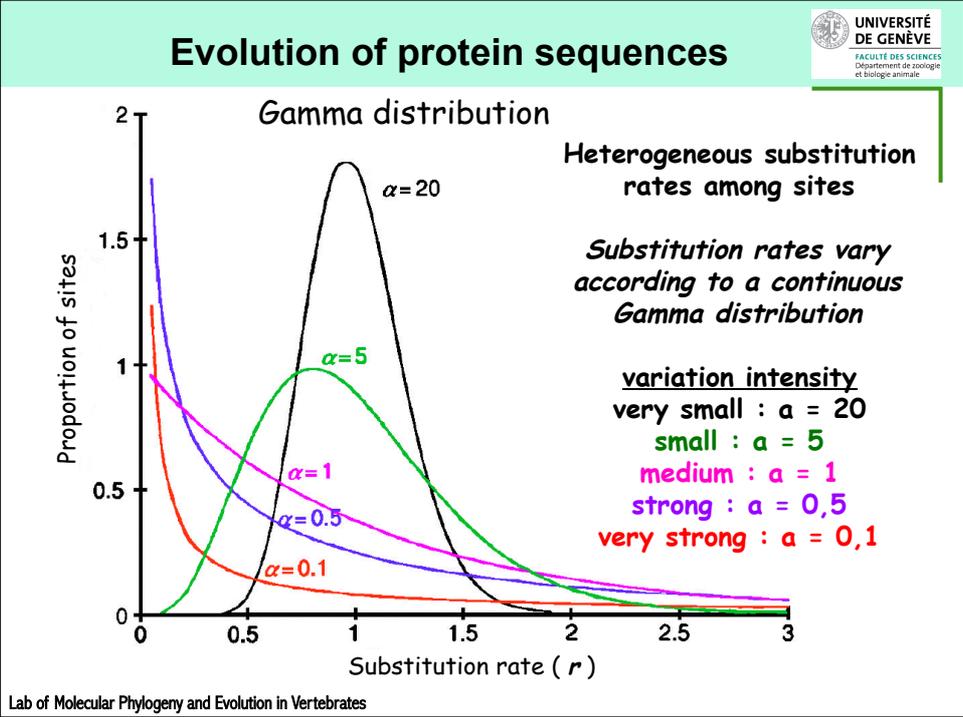
The distribution of the number of substitution per site (**k**) has a variance larger than the variance of the Poisson distribution (its an overdispersed Poisson distribution).

The distribution of **k** corresponds better to a **negative binomial distribution**, which has one more parameter than the Poisson distribution (Uzzell and Corbin, 1971).

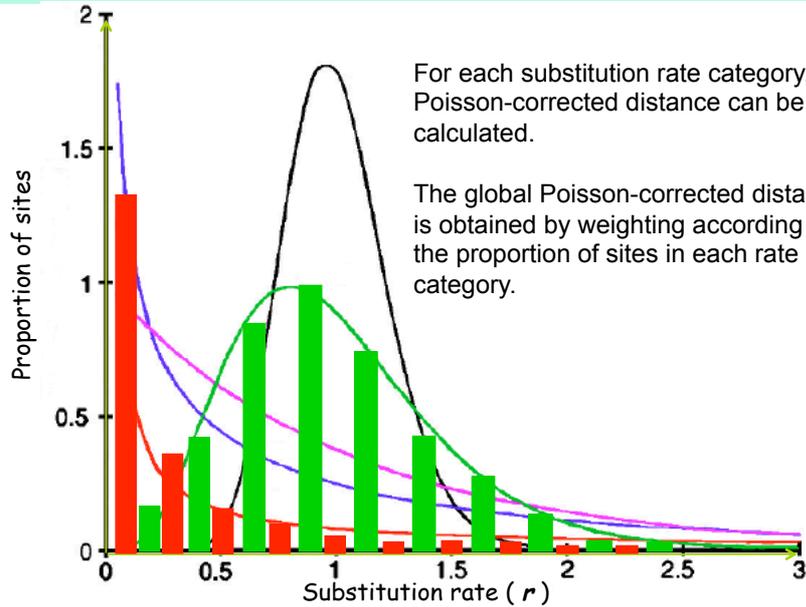
→ This implies that the substitution rate (**r**) varies among sites according to a **gamma distribution**.

Gamma distribution: a single parameter (**α**) defines the shape of the distribution.

Lab of Molecular Phylogeny and Evolution in Vertebrates



Evolution of protein sequences



Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of protein sequences

Examples of Poisson-corrected distance without / with incorporation of substitution rate heterogeneity across sites.

hemoglobin alpha chain	P-distance	PC-distance	PC + Gamma-distance
Human/cow	0.121	0.129	0.134
Human/kangaroo	0.186	0.205	0.216
Human/carp	0.486	0.665	0.789

However:

The **Poisson-corrected + Gamma** model of aa sequence evolution is still too simplistic.

Empirical corrections fit better to the complexity of the evolution of aa sequences.

Lab of Molecular Phylogeny and Evolution in Vertebrates

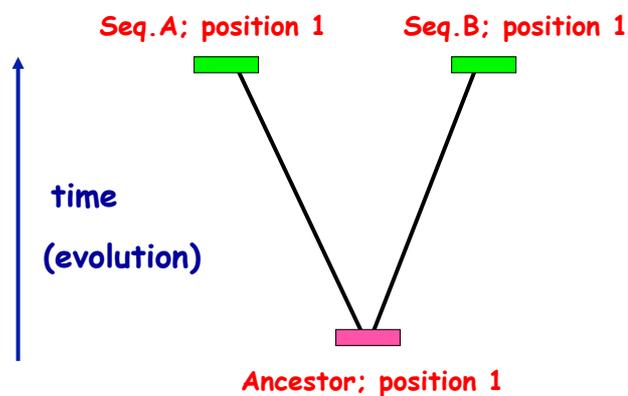
Molecular evolution

- Molecular basis of evolution.
- Evolution of protein sequences
- Evolution of DNA sequences

Evolution of DNA sequences

Observed distance or p -distance:

Number of different nucleotides / Number of compared nucleotides



Evolution of DNA sequences

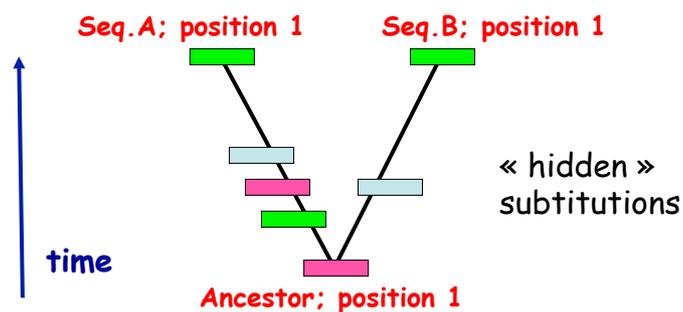
Calculating p -distance

Séq X	C	G	A	T	G	A	T	A	A	C
Séq Y	C	G	A	A	A	A	C	A	G	C
			*	*		*			*	

$$D_{xy} = \frac{\text{Nb of different sites (k)}}{\text{Nb of sites compared (N)}} = 4 / 10 = 0.4$$

Evolution of DNA sequences

p -distance vs evolutionary distance

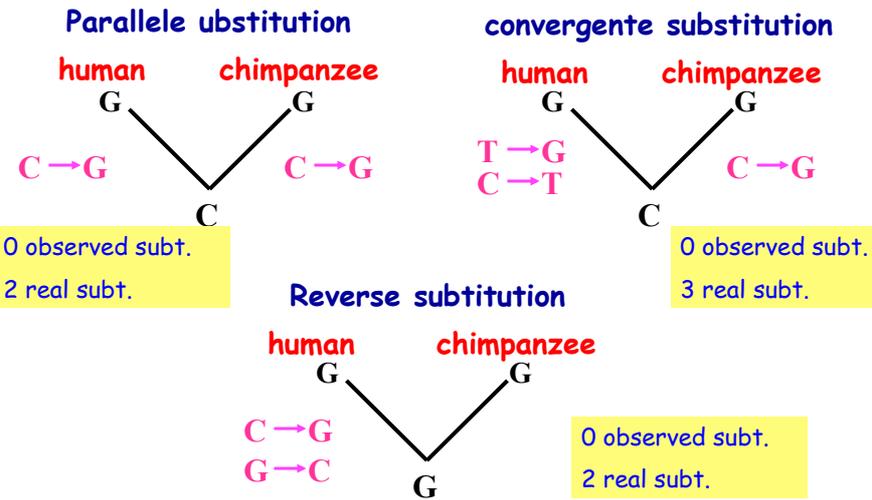


The amount of hidden substitutions increases with evolutionary time.

This problem is greater for DNA sequences (4 possible states per site) than for amino acid sequences (20 possible states per site).

Evolution of DNA sequences

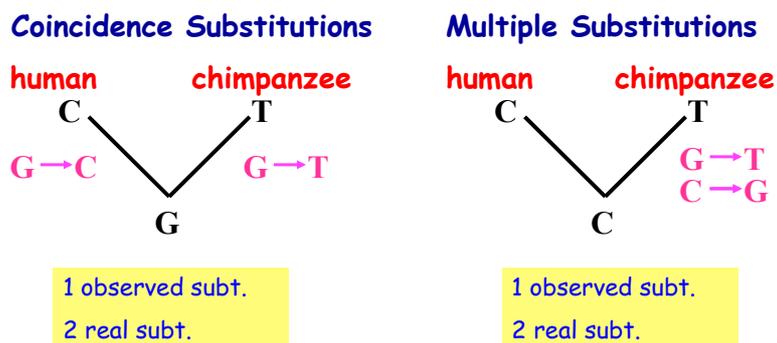
Hidden substitutions



Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of DNA sequences

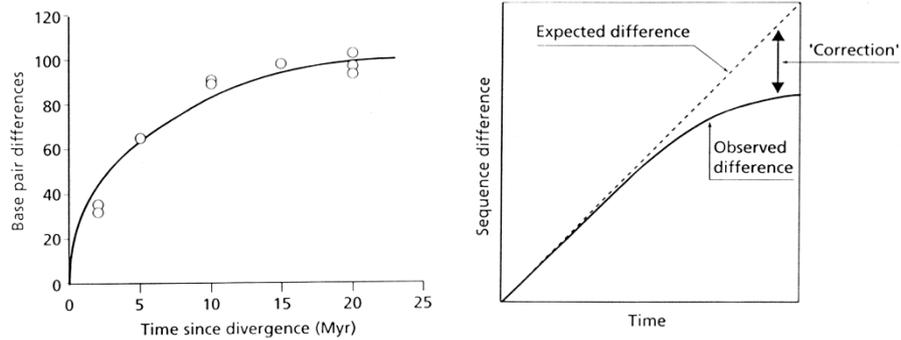
Hidden substitutions



Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of DNA sequences

Correcting p -distance



Page & Holmes, Molecular Evolution, 1998

Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of DNA sequences

To correct the p -distance, we need models of DNA sequence evolution.

First and simplest model of DNA sequence evolution:

the model of Jukes and Cantor, 1969 (JC69)

Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of DNA sequences

Jukes-Cantor (JC69) model

$$\text{Substitution Matrix } P_t = \begin{array}{c|cccc|c} & \text{A} & \text{C} & \text{G} & \text{T} & \\ \hline & \cdot & \alpha & \alpha & \alpha & \text{A} \\ \hline & \alpha & \cdot & \alpha & \alpha & \text{C} \\ \hline & \alpha & \alpha & \cdot & \alpha & \text{G} \\ \hline & \alpha & \alpha & \alpha & \cdot & \text{T} \\ \hline \end{array}$$

a single parameter: α

Parameters: $p_A = p_T = p_G = p_C$ and $\alpha = \beta$ α : transitions
 β : transversions

Formula: corrected distance $d_{xy} = -3/4 \ln(1 - 4/3 D)$

where D is the p -distance

Evolution of DNA sequences

Example

Two sequences of 500 nucleotides differing by 50 substitutions have a p -distance of:

$$D = 50 / 500 = \mathbf{0.1}$$

The evolutionary distance calculated according to the Jukes and Cantor model is:

$$d_{xy} = \mathbf{0.1073}$$

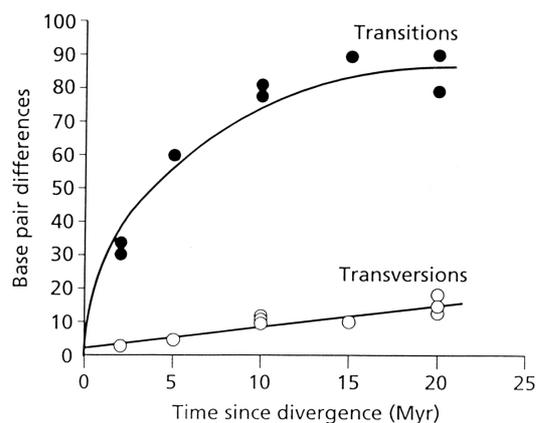
The corrected number of substitutions is 53.66,
that is 3-4 substitutions that were not observed

Evolution of DNA sequences

More realistic and complex models have been developed, with the help of the increase of computation power.

Evolution of DNA sequences

Transitions are generally more frequent than transversions



	Ts/Tv ratio
mt DNA	9.0
12 S rRNA	1.75
α - et β -globines	0.66

Accumulation of transitions and transversions in COI sequences of bovids

Evolution of DNA sequences

Kimura 2 parameters model (K2P)

$$P_t = \begin{array}{c|cccc|c} & A & C & G & T & \\ \hline & \cdot & \beta & \alpha & \beta & A \\ & \beta & \cdot & \beta & \alpha & C \\ & \alpha & \beta & \cdot & \beta & G \\ & \beta & \alpha & \beta & \cdot & T \\ \hline \end{array}$$

2 free parameters:

α = transitions

β = transversions

Parameters: $pA = pT = pG = pC$
• $\alpha \neq \beta$

• Formula: $d_{xy} = -1/2 \ln(1-2P-Q) + 1/4 \ln(1-2Q)$

• where P = proportion of transitions; Q = proportion of transversions

Evolution of DNA sequences

More complex models of sequence evolution

F81 Model (Felsenstein, 1981)

$\pi_A \neq \pi_C \neq \pi_G \neq \pi_T$ $\alpha = \beta$

4 free parameters

HKY85 Model (Hasegawa, Kishino, Yano, 1985)

$\pi_A \neq \pi_C \neq \pi_G \neq \pi_T$ $\alpha \neq \beta$

5 free parameters

Evolution of DNA sequences

The more complex model of sequence evolution

GTR Model
General Time Reversible
(Simon and Tavaré, 1986)

$$P_t = \begin{vmatrix} \cdot & a(pAC) & b(pAG) & c(pAT) \\ a(pCA) & \cdot & d(pCG) & e(pCT) \\ b(pGA) & d(pGC) & \cdot & f(pGT) \\ c(pTA) & e(pTC) & f(pTG) & \cdot \end{vmatrix}$$

Parameters: $\pi_A \neq \pi_C \neq \pi_G \neq \pi_T$

Every substitution has its own probability

$a = A \leftrightarrow C$ $b = A \leftrightarrow G$ $c = A \leftrightarrow T$
 $d = C \leftrightarrow G$ $e = C \leftrightarrow T$ $f = T \leftrightarrow G$

9 free parameters

Evolution of DNA sequences

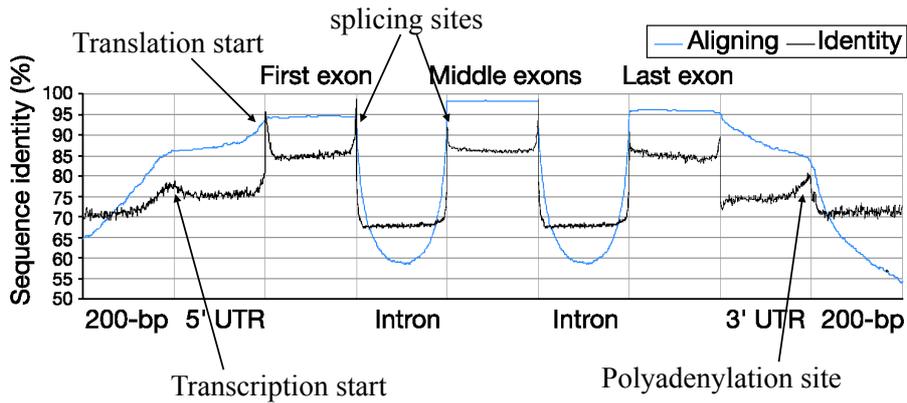
Most phylogenetic methods assume that:

1. Substitutions are independent from one another.
2. The substitution rate is constant in time and among lineages.
3. Base frequency is homogeneous among lineages.
4. The substitution rate is equal among sites and constant in time.

These assumptions are not always satisfied!

Evolution of DNA sequences

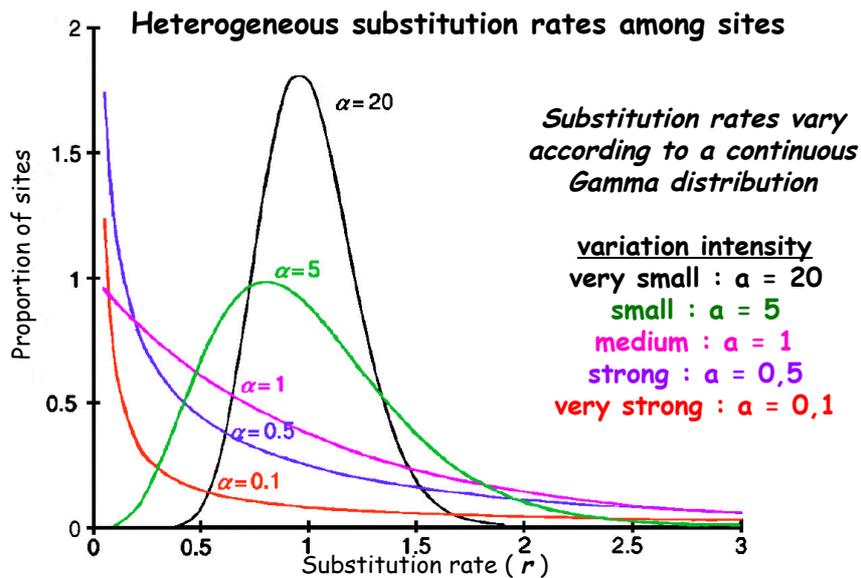
Substitution rate along a 'model' gene
Mean sequence identity of 3,165 human-mouse comparisons



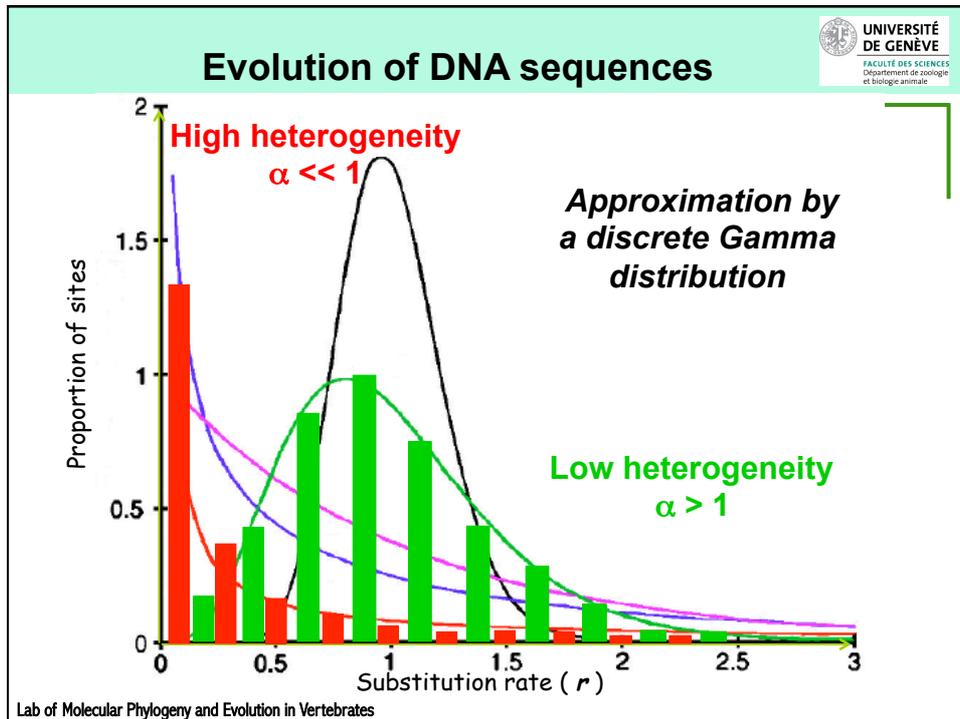
Lab of Molecular Phylogeny and Evolution in Vertebrates

MGSC Nature (2002)

Evolution of DNA sequences



Lab of Molecular Phylogeny and Evolution in Vertebrates



Evolution of DNA sequences



UNIVERSITÉ DE GENÈVE
FACULTÉ DES SCIENCES
Département de zoologie et biologie animale

Estimation of the parameter α (Yang, 1996):

• albumin	1.05
• insulin	0.40
• prolactin	1.37
• 16S rRNA (stem)	0.29
• 16S rRNA (loop)	0.58
• 12S rRNA mt	0.16
• D-loop	0.17

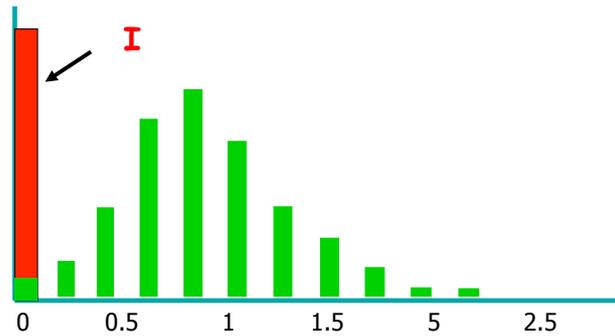
Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of DNA sequences

Invariant sites

Sometimes, the category of invariant sites does not fit to the distribution.

The category of invariant sites can then be treated as an additional parameter: **I**



The Gamma distribution is then used only for the remaining categories.

Evolution of DNA sequences

ModelTest: 56 models

Posada and Crandall, 1998

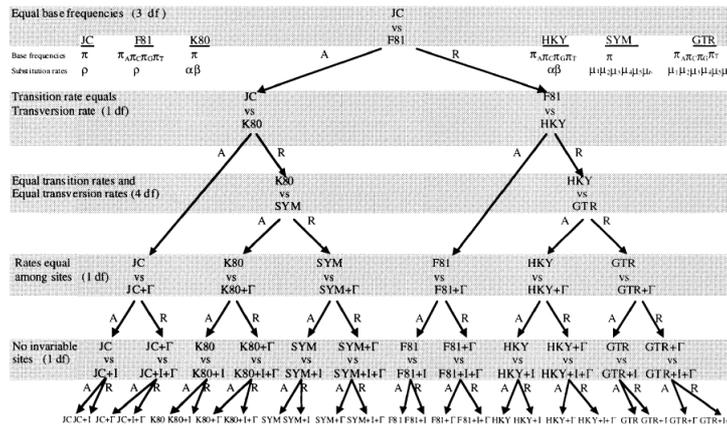


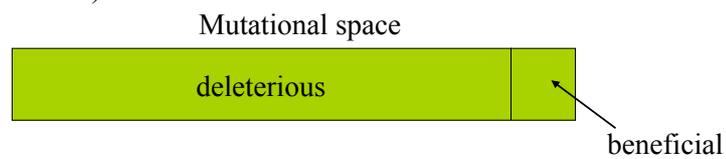
Fig. 1. Hierarchical hypothesis testing in MODELTEST. At each level the null hypothesis (upper model) is either accepted (A) or rejected (R). The models of DNA substitution are: JC (Jukes and Cantor, 1969), K80 (Kimura, 1980), SYM (Zharkikh, 1994), F81 (Felsenstein, 1981), HKY (Hasegawa *et al.*, 1985), and GTR (Rodriguez *et al.*, 1990). Γ : shape parameter of the gamma distribution; I: proportion of invariable sites. df: degrees of freedom. 1: equal base frequencies (0.25). π_A : frequency of adenine, π_C : frequency of cytosine, π_G : frequency of guanine, π_T : frequency of thymine. ρ : equal substitution rate, α : transition rate, β : transversion rate, μ_1 : A=C rate, μ_2 : A=G rate, μ_3 : A=T rate, μ_4 : C=G rate, μ_5 : C=T rate, μ_6 : G=T rate.

Evolution of DNA sequences

Selective pressure acting on DNA sequences

Evolution of DNA sequences

Functional view of mutations before the Neutral Theory of Evolution
(Kimura 1968)



- Most mutations were thought to be deleterious, reducing the fitness of the organism
 - the fate of such mutations is to disappear via **purifying selection**.
- Some mutations might be beneficial, increasing the fitness of the organism
 - such mutations are kept via **positive Darwinian selection (adaptive)**

The fitness is a function of survival and fecundity (number of offspring)

Evolution of DNA sequences

Motoo Kimura

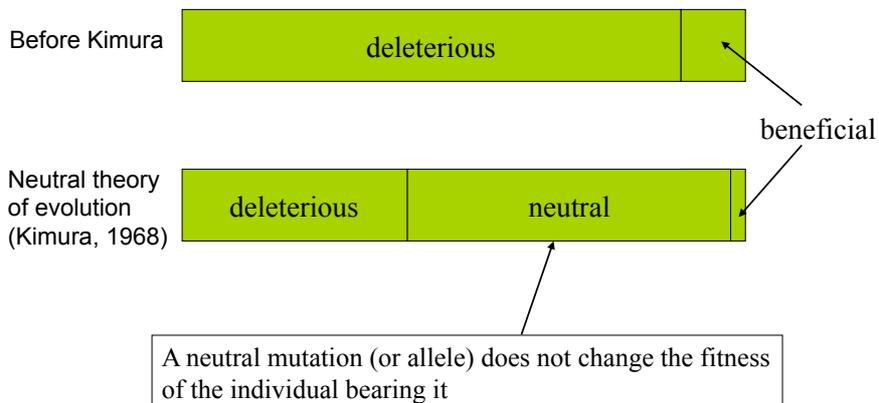


In 1968 Kimura analyzed the patterns of substitution in several genes (hemoglobin, cytochrom c and triosephosphate dehydrogenase)

Kimura (1968) *Nature* 217 624-626

Evolution of DNA sequences

Kimura proposed that most of the substitutions are neutral (no change in fitness) and that very few have been kept due to positive Darwinian selection.



Evolution of DNA sequences

Synonymous mutations can be neutral mutations

Predictions:

- If most changes in gene sequence were due to positive selection then we would expect that most changes would occur in 1st or 2nd codon positions (as they generally change the amino acid)
- If changes in gene sequence include neutral mutations, then 3rd codon positions would change at a higher rate because they are mostly synonymous mutations without any effect on fitness (=neutral mutations)

Which prediction is correct ?

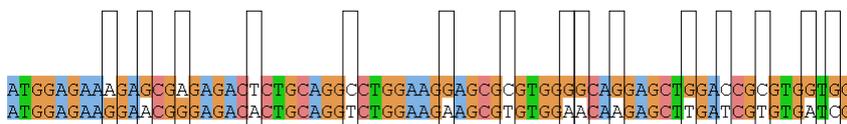
King, J. L., and Jukes, T. H. 1969. Non-Darwinian evolution, Science 164, 788-798.

Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of DNA sequences

Synonymous and non-synonymous substitution rates

- K_S = number of Synonymous substitutions per synonymous site
- K_A = number of non-synonymous substitutions (Altering) per non-synonymous site



It is possible to estimate K_a and K_s empirically, by examining the alignments.

However, appropriate methods have been developed to estimate K_a and K_s .

Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of DNA sequences

Synonymous substitution rates are generally higher than non-synonymous substitution rates

TABLE 8.1 Numbers of nucleotide substitutions per 100 sites between species^a

Species pair	Synonymous sites		Nonsynonymous sites	
	K_S	L^b	K_A	L^b
Mouse–rat	18.0 ± 0.7	4,229	1.8 ± 0.1	15,217
Mouse–hamster	30.3 ± 1.0	4,229	2.9 ± 0.1	15,217
Rat–hamster	31.3 ± 1.0	4,229	2.7 ± 0.1	15,217
Mouse–human	53.4 ± 1.5	4,229	5.2 ± 0.2	15,217
Rat–human	51.6 ± 1.5	4,229	5.0 ± 0.2	15,217
Hamster–human	52.3 ± 1.5	4,229	5.1 ± 0.1	15,217

From O'hUigin and Li (1992).

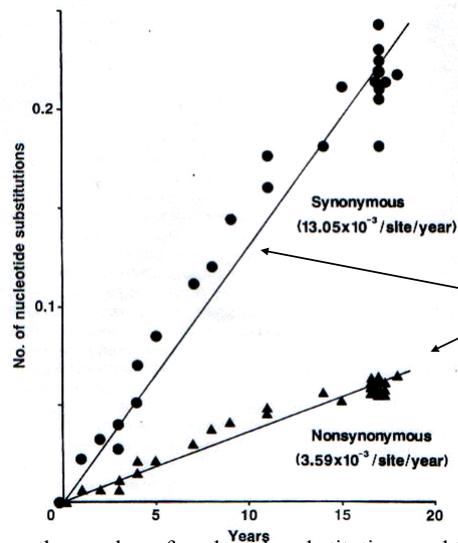
^aComputed by Li et al.'s (1985b) method.

^bNumber of sites compared.

From Li, W-H Molecular Evolution who took it from:
Oh' Uigin and Li. JME 1992 35: 377-384

Lab of Molecular Phylogeny and Evolution in Vertebrates

The Molecular Clock of Viral Evolution



Relationship between the number of nucleotide substitutions and the difference in the year of isolation for the H3 hemagglutinin gene of human influenza A viruses. All sequence comparisons were made with the strain isolated in 1968.

Gojobori et al. 1990 PNAS 87 10015-10018

Lab of Molecular Phylogeny and Evolution in Vertebrates

Evolution of DNA sequences

The ratio K_A/K_S is used to assess the selective pressure acting on coding regions.

$K_A/K_S \approx 1$ -> neutral evolution, no selection.

$K_A/K_S > 1$ -> positive selection

$K_A/K_S \ll 1$ -> negative selection (purifying selection)