# Is protein sequence evolution constant over time?

Carolin Kosiol &
Nick Goldman

```
goldman@ebi.ac.uk
http:/www.ebi.ac.uk/goldman
```

# Are Markov process models appropriate for protein sequence evolution?

# NO

# ?

# Evidence of non-Markov evolution of amino acid sequences

**Amino Acid Substitution Matrices From Protein Blocks**

*S. Henikoff and J.G. Henikoff*
*Proceedings of the National Academy of Sciences*
*of the United States of America* 89:10915–10919.  1992

**Tree-based Maximal Likelihood Substitution Matrices and Hidden Markov Models**

*G. Mitchison and R. Durbin*
*Journal of Molecular Evolution* 41:1139–1151.  1995

**Amino Acid Substitution During Functionally Constrained Divergent Evolution of Protein Sequences**

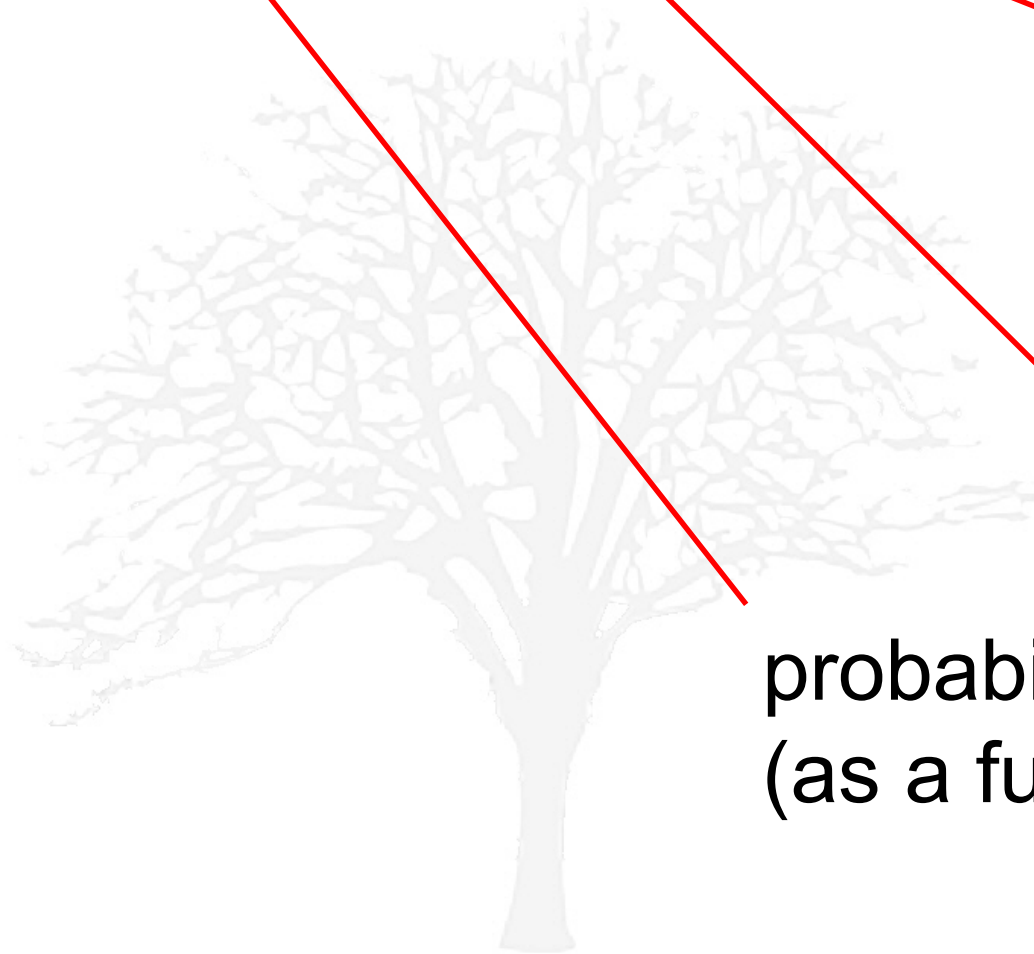*S.A. Benner, M.A. Cohen and G.H. Gonnet*
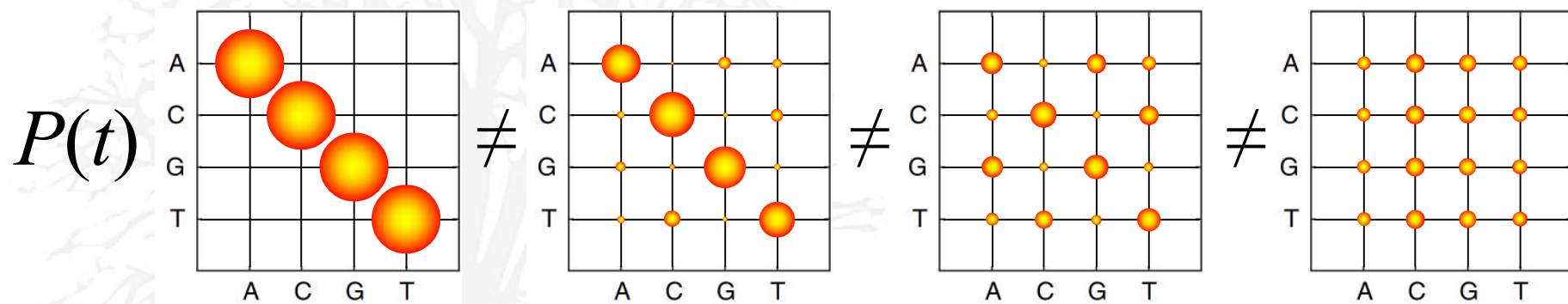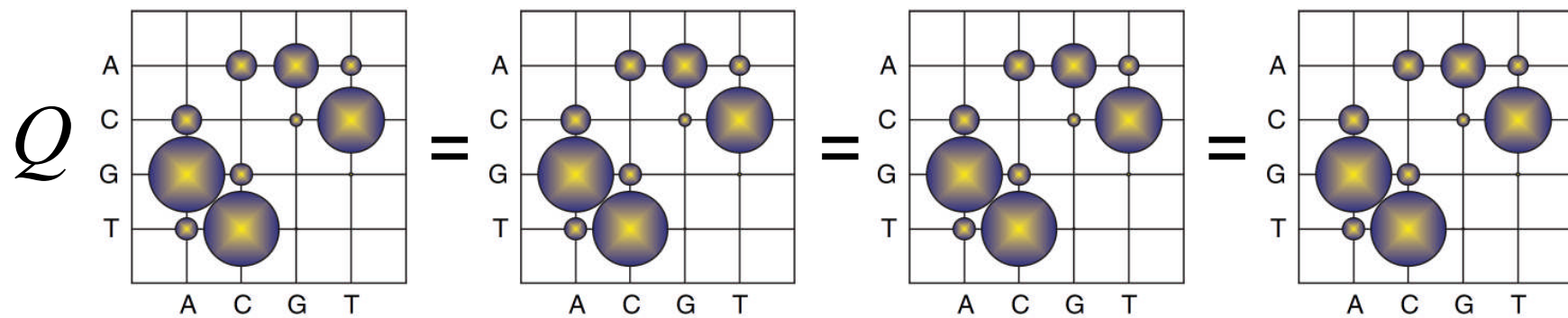*Protein Engineering* 7:1323–1332.  1994

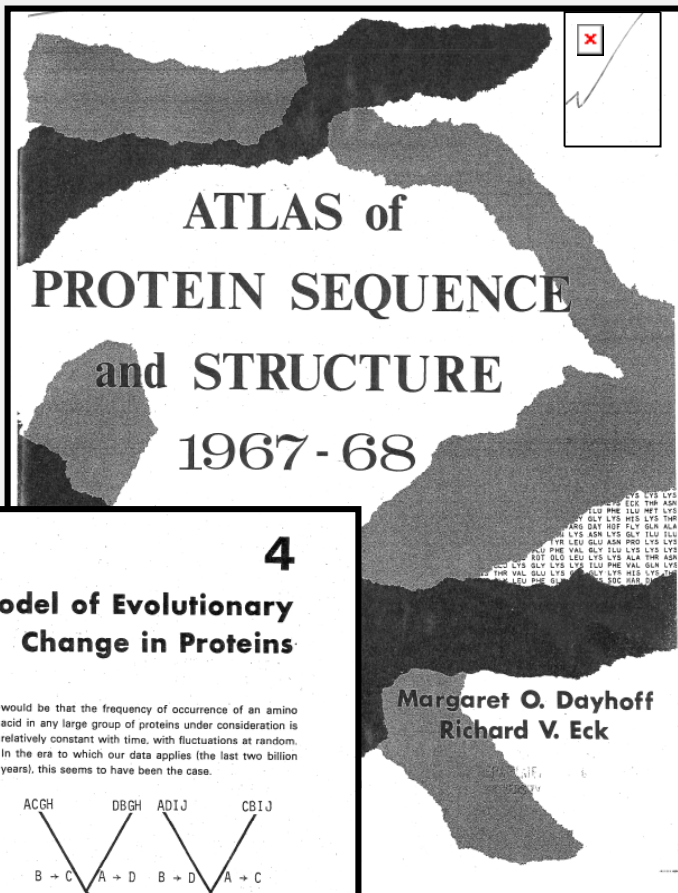$$P(t) = \exp(tQ) = I + tQ + \frac{(tQ)^2}{2!} + \frac{(tQ)^3}{3!} + \cdots$$

instantaneous
rate matrix

time

probability of change
(as a function of time)

It is possible
to infer $P(t)$
from sequence data…

ATLAS of PROTEIN SEQUENCE and STRUCTURE 1967-68

Margaret O. Dayhoff
Richard V. Eck

4

A Model of Evolutionary Change in Proteins

What mutations are most likely to be accepted? Which amino acids are least likely to change? How does the passage of time affect the similarity of related protein sequences?

would be that the frequency of occurrence of an amino acid in any large group of proteins under consideration is relatively constant with time, with fluctuations at random. In the era to which our data applies (the last two billion years), this seems to have been the case.
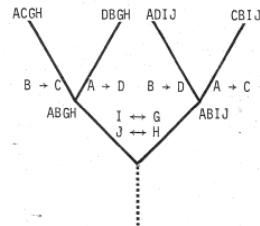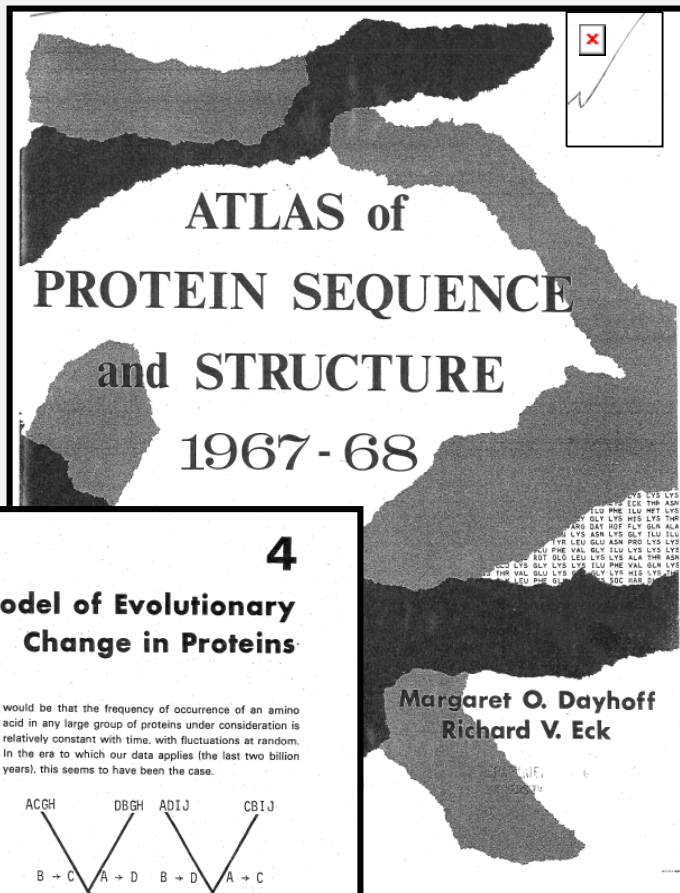
**Accepted Point Mutations**

An *accepted point mutation* is an exchange of one amino acid for another, accepted by natural selection. It is the result of two distinct processes: the first is the occurrence of the mutation in the gene and the second is its acceptance by natural selection as an improvement. To be accepted, the new amino acid side chain usually functions in a similar way to the old one. This plausible conjecture is supported by the chemical and physical similarities between amino acids which are observed to interchange frequently. Some examples are given in Chapter 5.

…and possible to infer *Q* from *P(t)*

**Dayhoff**

**Different Versions of the Dayhoff Rate Matrix**

*Carolin Kosiol and Nick Goldman*
*Mol. Biol. Evol. 22(2):193–199. 2005*

Many phylogenetic inference methods are based on Markov models of sequence evolution. These are usually expressed in terms of a matrix ($Q$) of instantaneous rates of change but some models of amino acid replacement, most notably the PAM model of Dayhoff and colleagues, were originally published only in terms of time-dependent probability matrices ($P(t)$). Previously published methods for deriving $Q$ have used eigen-decomposition of an approximation to $P(t)$. We

**"Matrix space"**

non-stochastic matrices
stochastic matrices

Dayhoff

*not* constant, according to Henikoff x 2 (BLOSUM) and Mitchison & Durbin

$\delta = 1.0$

$\delta^*$

$P_D(\delta)$

PAM1 $(\delta = 0.01)$

$\delta = 0.02$    $\delta = 0.05$

$(\text{PAM1})^n$

$P(t) = e^{tQ}$

$I$

$0$

time $t$

Benner *et al.* found rate matrix elements varied with observed divergence

They argued that the genetic code influences the matrix strongly at early stages of divergence, while physicochemical properties are dominant at later stages

Mitchison & Durbin found the accumulation of amino acid replacements that could be generated by a single nucleotide change was inconsistent with a simple Markov process

Time / Evolutionary distance

Low divergence

Medium divergence

High divergence

Time / Evolutionary distance

Time / Evolutionary distance

Medium divergence  ence

So, how **will** we explain the
evidence of non-Markov behaviour? —
the <u>a</u>ggregated <u>M</u>arkov <u>p</u>rocess (AMP):

$$\ldots \quad \rightarrow \quad X(t_k) = CTT \quad \rightarrow \quad X(t_{k+1}) = CCT \quad \rightarrow \quad \ldots$$

Markov process
(codon evolution)

$\downarrow f$ $\qquad\qquad\qquad$ $\downarrow f$

Deterministic
function on states
(genetic code)

$$\ldots \qquad Y(t_k) = L \qquad\qquad Y(t_{k+1}) = P \qquad \ldots$$

Non-Markov process
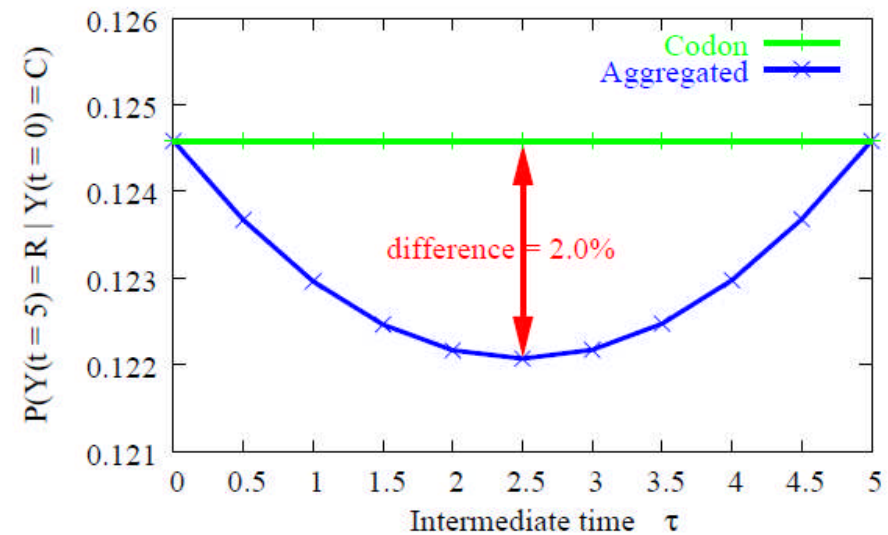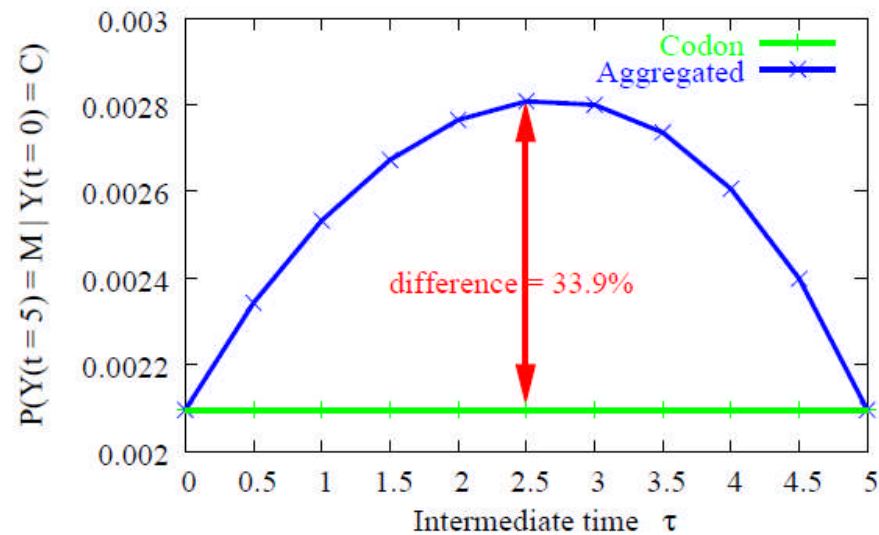(protein evolution)

time $t$

$$q_{ij,i \neq j} = \begin{cases} 0 & \text{if } i \text{ or } j \text{ is a stop codon or requires} > 1 \text{ nucleotide substitution} \\ \\ \pi_j & \text{if } i \rightarrow j \text{ synonymous transversion} \\ \\ \pi_j \kappa & \text{if } i \rightarrow j \text{ synonymous transition} \\ \\ \pi_j \omega & \text{if } i \rightarrow j \text{ nonsynonymous transversion} \\ \\ \pi_j \kappa \omega & \text{if } i \rightarrow j \text{ nonsynonymous transition} \end{cases}$$
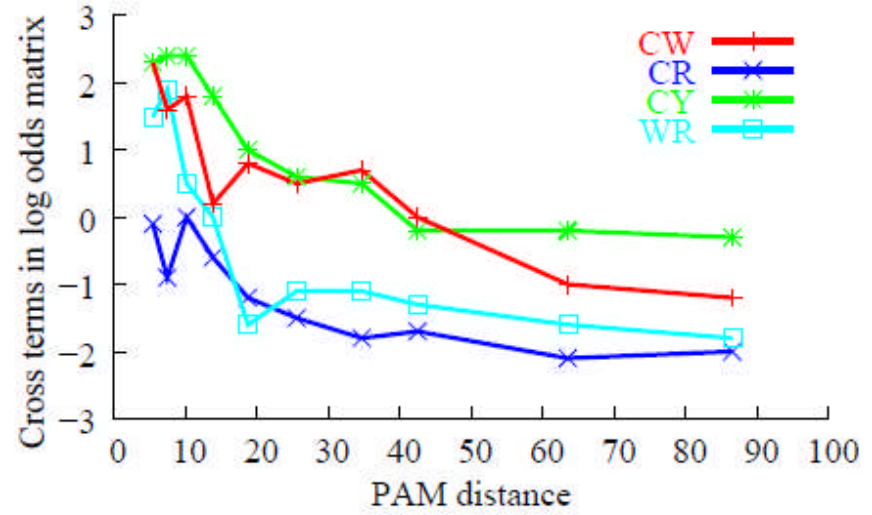
$$\kappa = 2.5 \qquad \omega = 0.2$$
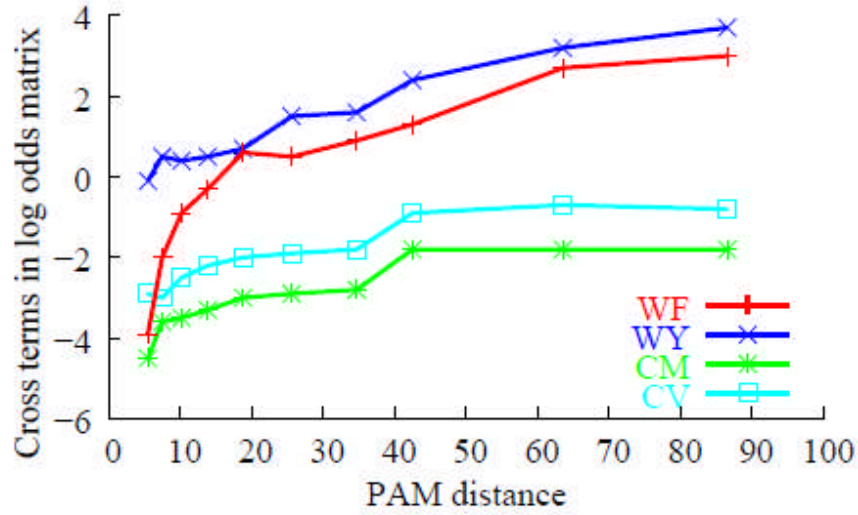
$$r_1 = 0.00001, \quad r_2 = 0.0001, \quad r_3 = 0.0001, \quad r_4 = 0.001,$$
$$r_5 = 0.01, \qquad r_6 = 0.1 \qquad r_7 = 0.15, \qquad r_8 = 0.2,$$
$$r_9 = 0.3, \qquad r_{10} = 0.5, \qquad r_{11} = 2.0, \qquad r_{12} = 8.73889$$

# Aggregated Markov processes are not Markov:

$$P(Y(t_1) = \mathrm{M}|Y(t_0) = \mathrm{C}) = \sum_{i=1}^{20} P(Y(t_1) = \mathrm{M}|Y(\tau) = \mathrm{A}_i) \times P(Y(\tau) = \mathrm{A}_i|Y(t_0) = \mathrm{C})$$
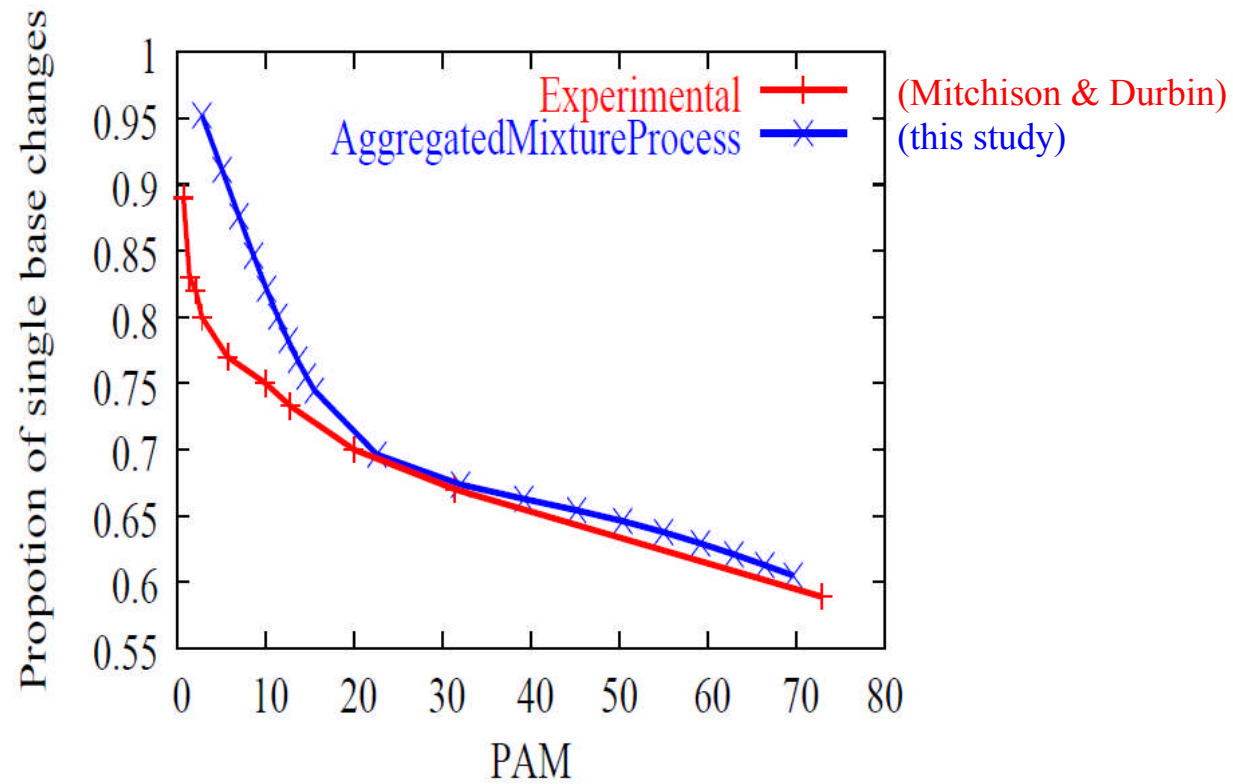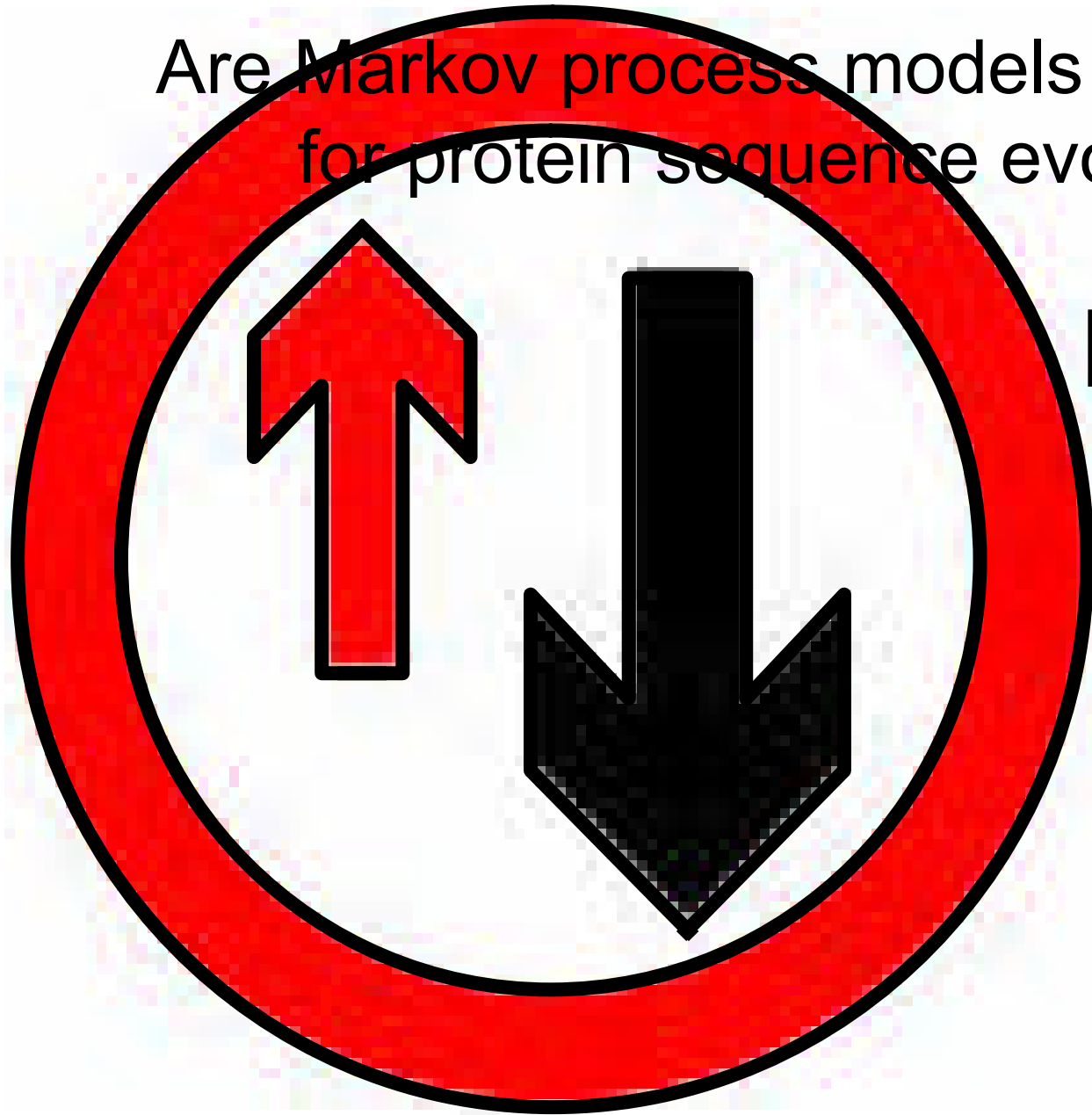
# Benner *et al.* evidence:



# this study:

# Mitchison & Durbin evidence:

Are Markov process models appropriate for protein sequence evolution?

PROCEED WITH CAUTION

# Things to remember from Nick's talk:

⇅ evolution should look the same whether we study it 100MYA or 1MYA or 1YA or today or tomorrow or …

⇅ published evidence of non-Markov protein evolution can be explained by a time-independent codon model-based AMP

⇅ we may proceed with current approaches to sequence evolution based on Markov models!

⇅ possible consequences:  non-Markov evolution of:
  ⇅ protein sequences
  ⇅ purine/pyrimidine (R/Y) encoded DNA
      (nucleotide-based AMP)