# Protein Communication System: Evolution and Genomic Structure

Nidhal Bouaynaya and Dan Schonfeld

University of Illinois at Chicago · nbouay1@uic.edu · dans@uic.edu

### Protein Communication Channel

Adaptive evolution has fashioned living organisms as agents of information acquisition, analysis, storage, and transmission. How has this happened? How have living systems evolved to handle the same problems with which we are confronted in this so-called Information Age: problems of information storage and processing, problems of transmission and reliability?

#### **Protein Communication Channel**



#### Analogy and Differences with a Communication Engineering System

Protein Communication System	Communication Engineering System
Set of proteins of the cell	Video
DNA	MPEG
No biological encoder	Encoder
Translation process	Decoder
? (The genomic structure provides the answer)	Objective function: minimize the probability
	of error

The protein communication channel is uniquely characterized by the probability transition matrix,  $Q(k) = \{q_{i,i}(k)\}, 1 \le i, j \le 20$ , at time *k* of the amino acids.

[1-1	6 a	0	2	-	0	α	0	0	0	0	0	0	α	0	0	α	α	α	0	0	• )
0		1-8α	0	0	a	α	0	0	0	0	0	0	0	0	a	2α	0	0	a	α	a
a		0	$1 = 8 \alpha$	2α	0	α	α	0	0	0	0	α	0	0	0	0	0	α	0	α	0
0		0	2α	1 - 8α	0	α	0	0	α	0	0	0	0	α	0	0	0	α	0	0	a
0		α	0	0	1 - 8 a	0	0	α.	0	3α	0	0	0	0	0	α	0	α	0	α	0
α		<u>α</u> τ	<u>a</u>	<u>-</u>	0	1 - 6 a	0	0	0	0	0	0	0	0	20	<u></u>	0	α	4	0	-
0		0	α	0	0	0	1-8α	0	0	a	0	α	α	2α	a	0	0	0	0	α	0
0		0	0	0	10	0	0	1 - 7α	<u></u>	<u>40</u> 2	α	1 11	0	0	<u>0</u> 2	10	α	α	0	0	0
0		0	0	α	0	0	0	<u>a</u>	1 - 8α	0	<u>a</u>	2α	0	α	a	0	α	0	0	0	a
0		0	0	0	α	0	<u>a</u> 3	20	0	1 - 6 a	2	0	20	- <u>#</u>	10	<u>x</u>	0	α	<u>6</u>	0	<u>n</u>
0		0	0	0	0	0	0	3α	α	2α	1 - 9α	0	0	0	a	0	α	α	0	0	0
0		0	0.	0	0	0	α	0.	2α	0	0	1-8α	0	0	0	0.	0.	0	0	α	0
۵		0	0	0	0	0	<u>α</u> 1	0	0	a	0	0	1 - 6α	<u>-</u>	a	α	α	0	0	0	0
0		0	0	α	0	0	2α	0	α	a	0	0	α	1 - 8α	a	0	0	0	0	0	a
0		<u>0</u>	0	0	0	α	<u>a</u>	S.	<u>#</u>	10	<u>0</u> 5	0	20	<u>#</u>	1 - 6 a	α	<u> </u>	0	<u>0</u> 2	0	<u>0</u> 2
10	×	24	0	0	<u>0</u> 2	<u>0</u> 2	0	<u>0</u>	0	<u>0</u> 2	0	<u>0</u>	20	0	α	$1 = \frac{z + \alpha}{3}$	α	0	<u>0</u> 8	<u>a</u>	<u>0</u> 1
α		0	0	0	0	0	0	30	<u>n</u>	0	<u>0</u> 4	$\frac{\alpha}{t}$	α	0	<u>0</u> 1	20	$1 = 6 \alpha$	0	0	0	0
α		0	<u>a</u>	<u>+</u>	<u>0</u> 1	α	0	30	0	20	<u>0</u>	0	0	0	0	0	0	$1 - 6 \alpha$	0	0	0
0		2α	0	0	0	α	0	0	0	a	0	0	0	0	2α	α	0	0	1 - 9 a	0	2 a
0		α	π	0	a	0	α	0	0	0	0	α	0	0	0	α	0	0	0	1 - 8 α	2 a
0		10	0	<u>10</u> 3	0	<u>α</u>	0	0	<u>10</u> 3	a	0	0	0	<u>10</u> 3	$\frac{1 \alpha}{3}$	α	0	0	$\frac{1 \alpha}{3}$	4=	$1 - \frac{12 \alpha}{3}$

P: a first-order Markov probability transition matrix between amino acids. Only the terms of the first degree in b (k) are retained. For display clarity, the dependence on the time k has been omitted.



#### where **Q** *E* {PAM, **P**}.

P takes into account all possible mutations between amino acids whether they are accepted or rejected by natural selection. The PAM transition matrix is estimated from protein sequences and hence takes into account the accepted mutations only.

### **Evolution: Constant Point Mutation Rate**

Proposition 1: (Convergence of the amino acid distribution) Consider an initial probability distribution of the amino acids at time 0,  $p_0$  Then, the probability distribution of the amino acids converges, over time, towards a stationary distribution given by  $s_1$  if  $\mathbf{Q} = \mathbf{P}$  and  $s_2$  if  $\mathbf{Q} = \mathbf{P}AM_{250}$ . where

 $\mathbf{s_1} = (0.062, \ 0.0312, \ 0.0312, \ 0.0312, \ 0.0312, \ 0.062, \ 0.0312, \ 0.046, \ 0.0312, \ 0.093, \ 0.0156, \ 0.0312, \ 0.0625, \ 0.0312, \ 0.0312, \ 0.0312, \ 0.0625, \ 0.0312, \ 0.0625, \ 0.0312, \ 0.0625, \ 0.0312, \ 0.0625, \ 0.0312, \ 0.0625, \ 0.0312, \ 0.0625, \ 0.0312, \ 0.0625, \ 0.0312, \ 0.0625, \ 0.0312, \ 0.0625, \ 0.0312, \ 0.0625, \ 0.0312, \ 0.0625,$ 

0.0937, 0.0937, 0.0625, 0.0625, 0.0156, 0.0312)

 $\mathbf{s}_2 = (0.087, \ 0.041, \ 0.042, \ 0.048, \ 0.034, \ 0.039, \ 0.051, \ 0.091, \ 0.033, \ 0.036, \ 0.083, \ 0.08, \ 0.014, \ 0.038, \ 0.014, \ 0$ 

0.053, 0.07, 0.06, 0.0089, 0.028, 0.064).

#### The experimental distribution is

 $\mathbf{r}=(0.0868, 0.0213, 0.059, 0.054, 0.0377, 0.0786, 0.023, 0.0508, 0.0672, 0.077, \ 0.018, 0.0491, 0.041, 0.0393, 0.0491, 0.041, 0.0393, 0.0491, 0.0$ 

0.0426, 0.0737, 0.0606, 0.0688, 0.0131, 0.0377).

Proposition 2: (*Rate of Convergence*) { $p_0Q_k$ }k $\geq 1$  converges at a **geometric rate** with parameter  $m_2$ , where  $\begin{cases} m_2 = 0.53, & \text{if } \mathbf{Q} = PAM_{250}; \\ m_2 \neq 1 - p/2, & \text{if } \mathbf{Q} = \mathbf{P}. \end{cases}$ 

## **Evolution: Time-Varying Point Mutation Rate**

**Theorem 1:** (Weak Ergodicity result) Consider a finite number of PAM matrices denoted by PAM(1),  $\cdots$ , PAM((*N*), where PAM(i) can be PAM<sub>1</sub> or PAM<sub>160</sub> or PAM<sub>250</sub>, etc, for all *i* = 1,  $\cdots$  *N*, Consider the sequence:  $T_{p,k} = t_{p+1}t_{p+2} \cdots t_{p+k}$ , where each ti.  $\mathbb{R}$  {PAM(1),  $\cdots$  PAM((*N*)}. That is at each time *k*, the probability transition matrix is some PAM matrix. Then,  $T_{p,k}$  is weakly ergodic at a uniform geometric rate for all  $p \neq 0$ . So the sequence { $p_{k}$ } thends to a sequence of distributions independently of  $p_{0}$ .

**Theorem 2:** (Strong Ergodicity Result) Consider a point mutation rate, b (k), which is bounded uniformly on k, i.e.,  $0 \le a \ne b$  (k)  $\ne b \le 1$ . Then the products  $T_{p,k} = P_{p+1} \cdots P_{p+k}$  are strongly ergodic. Thus, the sequence  $\{p_{k}\}_{k \ge 1}$  converges towards the stationary distribution  $s_1$  independently of the initial distribution  $p_0$ . Moreover, the convergence rate is at least geometric.



We show that introns protect coding regions in the DNA sequence from frequent errors in the way hollow uninhabited structures are used by the military to protect important installations, such as aircraft hangars and missile launching facilities, from a bomb attack by serving as a dummy target that resembles the protected structure.

### Genomic Structure: Deterministic Analysis

Under the assumption of a Poisson (m) noise, we obtain the probability of error:

$$P_{e} = 1 - e^{-\lambda KT} \prod_{k=1}^{K} \sum_{n=0}^{T-l_{k}} \frac{\lambda^{n} (T - l_{k})}{n!}$$

Taking the derivative of  $P_e$  e with respect to  $I_{k}$ , we obtain the following coupled system for the optimal exon lengths:

$$l_{i_0} = M \frac{[\prod_{k \neq i_0} \sum_{n=0}^{T-M} \frac{\lambda^n (T-l_k)^n}{n!}] [\sum_{n=1}^{T-M} \frac{\lambda^n (T-l_{i_0})^{n-1}}{(n-1)!}]}{\sum_{j=1}^{F} [\prod_{k \neq j} \sum_{n=0}^{T-M} \frac{\lambda^n (T-l_k)^n}{n!}] [\frac{\sum_{n=1}^{T-M} \lambda^n (T-l_j)^{n-1}}{(n-1)!}]}{\sum_{j=1}^{T-M} \sum_{n=1}^{T-M} \frac{\lambda^n (T-l_j)^{n-1}}{(n-1)!}}]$$

An obvious solution is obtained when  $I_k = M / K$  for all  $k = 1, \dots, K$ .

The asymmetric distribution, which best approximates  $\delta_{MK}$  would have its mode very close to its mean. Amazingly, the exon length distribution of the human genome has its mode almost equal to its mean obtained at about 170 nucleotides!

### **Genomic Structure: Stochastic Analysis**

Probability of Error Analysis Let p(I) be the continuous distribution of the length of exons.

$$P_e = 1 - \left(\int_0^\infty e^{-\lambda l} p(l) \, dl\right)^K$$

Stochastic Optimization Problem:  $p^{*}(l) = \underset{p(l)}{\operatorname{arg\,max}} \int_{0}^{\infty} e^{-\lambda l} p(l) \, dl \quad \text{subject to} \quad \begin{cases} i \end{pmatrix} \int_{0}^{\infty} p(l) \, dl \\ ii \end{pmatrix} \int_{0}^{\infty} p(l) \, dl \ge l_{0} \end{cases}$ 

$$p^{*}(l) = \frac{p_{0}(1+\mu)}{e^{-\lambda l} + \gamma l^{1+\alpha} + \mu}$$

**Random Walk Model**  $X_N = \sum_{i=1}^N l_i, \quad l_i \square f(l) = \alpha l^{-(1+\alpha)}$ 

$$\lim_{x\to\infty} p(x \mid \alpha, \beta, \gamma, \delta) = \frac{C}{x^{1+\alpha}}$$

### **Experimental Results**

