



# A generalization of substitution evolution models of nucleotides to genetic motifs

Emmanuel Benard, Christian J. Michel\*

Equipe de Bioinformatique Théorique, FDBT, LSIT (UMR 7005), Université de Strasbourg, Pôle API, Boulevard Sébastien Brant, 67400 Illkirch, France

## ARTICLE INFO

### Article history:

Received 25 March 2011

Received in revised form

15 July 2011

Accepted 18 July 2011

Available online 4 August 2011

### Keywords:

Substitution models

Stochastic models

Genetic motifs

Analytical solutions

Research software

## ABSTRACT

We generalize here the classical stochastic substitution models of nucleotides to genetic motifs of any size. This generalized model gives the analytical occurrence probabilities of genetic motifs as a function of a substitution matrix containing up to three formal parameters (substitution rates) per motif site and of an initial occurrence probability vector of genetic motifs. The evolution direction can be direct (past-present) or inverse (present-past). This extension has been made due to the identification of a Kronecker relation between the nucleotide substitution matrices and the motif substitution matrices. The evolution models for motifs of size 4 (tetranucleotides) and 5 (pentanucleotides) are now included in the SEGM (Stochastic Evolution of Genetic Motifs) web server.

© 2011 Elsevier Ltd. All rights reserved.

## 1. Introduction

We present here a generalization of the classical stochastic substitution models of nucleotides to genetic motifs of any size. The first gene evolution model was proposed by Jukes and Cantor (1969) with 1-parameter substitution (probability  $\alpha$  for all nucleotide substitution types). It was generalized to a 2-parameter substitution model (Kimura, 1980) (probability  $\gamma$  for the nucleotide transitions  $A \leftrightarrow G$  and  $C \leftrightarrow T$ , and probability  $\beta$  for the nucleotide transversions  $A \leftrightarrow C$ ,  $A \leftrightarrow T$ ,  $C \leftrightarrow G$  and  $G \leftrightarrow T$ ) and then, to a 3-parameter substitution model (Kimura, 1981) (probability  $a$  for transitions, probability  $b$  for the transversion type  $A \leftrightarrow T$  and  $C \leftrightarrow G$ , and probability  $c$  for the transversion type  $A \leftrightarrow C$  and  $G \leftrightarrow T$ ). Later, these substitution models were generalized to a greater number of substitution parameters, e.g. a 6-parameter substitution model with equal base frequencies (Zharkikh, 1994).

Nucleotide substitution models were extended to genetic motif substitution models, e.g. Arquès and Michel (1993, 1995) for the pioneer work. The most recent motif substitution models (Michel, 2007a–c), i.e. trinucleotide models with three substitution rates per motif site, are based on a block matrix factorization (Tian and Syan, 2001). However, this approach cannot be used to generalize the substitution models to genetic motifs of any size. Indeed, the construction of large substitution matrices and their eigenvalues and eigenvectors determination are impossible by applying

classical methods of formal calculus with the current software (e.g. Mathematica 8.1 in 2011).

In applied mathematics, Kronecker operators are classically involved in Markov modulated Poisson processes, e.g. in communication (Burman and Smith, 1986; Salvador et al., 2003), etc. Thus, they were also used later in the phylogeny field involving Markov processes, in particular to diagonalize a modulated Markov matrix for the study of the variation in time of site-specific substitution rate (covarion model, (Galtier and Jean-Marie, 2004; Wang et al., 2007; Allman and Rhodes, 2009)), to compute the dinucleotide substitution process on a sequence of a given length (Lunter and Hein, 2004), to deduce the substitution rate of a dinucleotide (codon, respectively) by another dinucleotide (codon, respectively) from the product of substitution rates of nucleotides between these two dinucleotides (codons, respectively) (Mackiewicz et al., 2008; Darot et al., 2006), to study the  $n$ -taxon process in a tree-based model which analyses the evolution of vectors of states, one vector entry being associated with a taxa (Bryant, 2009), etc. By exploring similar mathematical strategies, a Kronecker relation allows to construct large motif substitution matrices from nucleotide substitution matrices per site as well as to determine the eigenvalues and eigenvectors associated with genetic motifs from the elementary eigenvalues and eigenvectors associated with nucleotides per site.

The SEGM (Stochastic Evolution of Genetic Motifs) web server is a web application which was built in 2009 to study evolution of nucleotides, dinucleotides and trinucleotides (Benard and Michel, 2009). Based on this Kronecker property, a new version of this SEGM web server is developed which includes several improved functionalities (in particular a faster computation) and an extension to motifs of size 4 (tetranucleotides) and 5

\* Corresponding author.

E-mail addresses: [benard@dpt-info.u-strasbg.fr](mailto:benard@dpt-info.u-strasbg.fr) (E. Benard), [michel@dpt-info.u-strasbg.fr](mailto:michel@dpt-info.u-strasbg.fr) (C.J. Michel).

(pentanucleotides). This research software extension allows biologists to study evolution of large motifs, for example promoter sites (CAAT box, TATA box, etc.). The theory proposed here is valid for genetic motifs of any size. Thus, application to hexanucleotides is obviously possible, for example, to study evolution of recoding signals stimulating read-through stop codons (Harrell et al., 2002). The current limit is only the PC power (CPU and memory). Indeed, the computation of the evolution probability of one pentanucleotide leads to an analytical solution with  $4^{10}$  terms, i.e. more than one million of terms. The SEGM web server is freely available at <http://lsit-bioinfo.u-strasbg.fr:8080/web/Mathematica/SEGM/SEGM.html>.

## 2. Mathematical model

The generalized model will give here the analytical occurrence probabilities of genetic motifs as a function of a substitution matrix containing up to three formal parameters (substitution rates) per motif site and of an initial occurrence probability vector of genetic motifs. Let us consider a motif of size  $n$  on the genetic alphabet  $\{A,C,G,T\}$ . By convention, a genetic motif is represented by its index  $i$ ,  $1 \leq i \leq 4^n$ , according to the lexicographical order, e.g. if  $n=3$  (trinucleotides), the index  $i=1$  refers to the first motif AAA and  $i=64$  to the last motif TTT. There are  $4^n$  motifs of size  $n$ . For all  $1 \leq i \leq 4^n$ , we denote by  $P_i(t)$ , the occurrence probability of motif  $i$  of size  $n$  at time  $t \geq 0$ .

### 2.1. Stochastic substitution model of genetic motifs of size $n$

The substitution process is handled by a differential equation which determines the occurrence probabilities of the  $4^n$  genetic motifs at time  $t \geq 0$ . The motifs mutate according to constant substitution probabilities. Let us consider two motifs  $ij$  of size  $n$ ,  $1 \leq ij \leq 4^n$ . We denote by  $P_T(j \rightarrow i)$ , the substitution probability of motif  $j$  into motif  $i$  during time  $T$ ,  $T > 0$ . The occurrence probability  $P_i(t+T)$  of motif  $i$  at time  $t+T$  is equal to the sum of probabilities  $P_j(t)$  of the  $4^n$  motifs  $j$  at previous time  $t$  times their substitution probabilities  $P_T(j \rightarrow i)$  into motif  $i$  during  $T$ , i.e.

$$P_i(t+T) = \underbrace{\sum_j P_j(t) P_T(j \rightarrow i)}_{\text{Probability of motif } i \text{ to appear}} \quad (2.1)$$

From Eq. (2.1), the derivative with respect to time  $P'_i(t) = \partial P_i(t) / \partial t$  of the occurrence probability of motif  $i$  at time  $t$  is

$$\begin{aligned} P'_i(t) &= \lim_{T \rightarrow 0} \left[ \frac{P_i(t+T) - P_i(t)}{T} \right] \\ &= \lim_{T \rightarrow 0} \left[ \frac{\sum_j P_j(t) P_T(j \rightarrow i) - P_i(t)}{T} \right] \\ &= \lim_{T \rightarrow 0} \left[ \frac{\sum_{j \neq i} P_j(t) P_T(j \rightarrow i) + P_i(t) P_T(i \rightarrow i) - P_i(t)}{T} \right], \end{aligned}$$

where  $P_T(i \rightarrow i)$  represents the probability that motif  $i$  does not mutate into a different motif  $j \neq i$  during  $T$ . Then,

$$\begin{aligned} P'_i(t) &= \lim_{T \rightarrow 0} \left[ \frac{\sum_{j \neq i} P_j(t) P_T(j \rightarrow i) - P_i(t) (1 - P_T(i \rightarrow i))}{T} \right] \\ &= \lim_{T \rightarrow 0} \left[ \frac{\sum_{j \neq i} P_j(t) P_T(j \rightarrow i) - P_i(t) \sum_{j \neq i} P_T(i \rightarrow j)}{T} \right] \\ &= \sum_{j \neq i} P_j(t) \lim_{T \rightarrow 0} \left[ \frac{P_T(j \rightarrow i)}{T} \right] - P_i(t) \sum_{j \neq i} \lim_{T \rightarrow 0} \left[ \frac{P_T(i \rightarrow j)}{T} \right]. \end{aligned}$$

For all motifs  $ij$ , the instantaneous substitution probability  $P(j \rightarrow i)$  of motif  $j$  into motif  $i$  is assumed to be constant along time. When  $T$  is small enough, there is no more than one motif substitution per motif site. Then, the following approximation

applies

$$P_T(j \rightarrow i) \underset{T \rightarrow 0}{=} P(j \rightarrow i) T$$

and consequently

$$\lim_{T \rightarrow 0} \left( \frac{P_T(j \rightarrow i)}{T} \right) = P(j \rightarrow i).$$

Finally, for any motif  $i$ , the derivative  $P'_i(t)$  is

$$\begin{aligned} P'_i(t) &= \sum_{j \neq i} P_j(t) P(j \rightarrow i) - P_i(t) \sum_{j \neq i} P(i \rightarrow j) \\ &= \sum_{j \neq i} P_j(t) P(j \rightarrow i) - P_i(t) (1 - P(i \rightarrow i)) \\ &= \sum_j P_j(t) P(j \rightarrow i) - P_i(t). \end{aligned} \quad (2.2)$$

Let  $P_n(t) = [P_i(t)]_{1 \leq i \leq 4^n}$  be the column vector of size  $4^n$  made of the probabilities  $P_i(t)$  for all  $1 \leq i \leq 4^n$ . From Eq. (2.2), we derive a matrix differential equation which describes the substitution process for genetic motifs

$$\begin{aligned} P'_n(t) &= M_n \cdot P_n(t) - P_n(t) \\ &= (M_n - I_n) \cdot P_n(t), \end{aligned} \quad (2.3)$$

where the symbol  $\cdot$  is the matrix product,  $I_n$  is the identity matrix ( $4^n, 4^n$ ) and  $M_n = [m_{ij}]_{1 \leq ij \leq 4^n}$  is the instantaneous substitution probability matrix whose element  $m_{ij}$  in row  $i$  and column  $j$  refers to the substitution probability of motif  $j$  into motif  $i$

$$m_{ij} = P(j \rightarrow i).$$

The instantaneous substitution probability matrix  $M_n$  is stochastic in column. Indeed, for all  $1 \leq j \leq 4^n$ , the elements of matrix  $M_n$  satisfy  $\sum_{1 \leq i \leq 4^n} m_{ij} = \sum_{1 \leq i \leq 4^n} P(j \rightarrow i) = 1$ . For all  $1 \leq j \leq 4^n$ , the diagonal elements  $m_{jj}$  of  $M_n$  satisfy

$$m_{jj} = 1 - \sum_{1 \leq i \leq 4^n, i \neq j} m_{ij}.$$

Eq. (2.3) is equal to Michel (2007a, Eq. (2)) obtained by a similar approach.

Let  $A_n = M_n - I_n$ . Then, Eq. (2.3) becomes

$$P'_n(t) = A_n \cdot P_n(t). \quad (2.4)$$

If  $A_n$  is diagonalizable, i.e.  $A_n = Q_n \cdot D_n \cdot Q_n^{-1}$  where  $D_n$  is the spectral matrix ( $4^n, 4^n$ ) and  $Q_n$  is its associated eigenvector matrix ( $4^n, 4^n$ ), then Eq. (2.4) becomes

$$P'_n(t) = Q_n \cdot D_n \cdot Q_n^{-1} \cdot P_n(t). \quad (2.5)$$

This differential Eq. (2.5) has the classical solution (Lange, 2005)

$$P_n(t) = Q_n \cdot e^{D_n t} \cdot Q_n^{-1} \cdot P_n(0), \quad (2.6)$$

where  $e^{D_n t}$  is the exponential spectral matrix ( $4^n, 4^n$ ) of matrix  $A_n$ ,  $Q_n$  is its associated eigenvector matrix ( $4^n, 4^n$ ) and  $P_n(0)$  is the vector of the  $4^n$  initial occurrence probabilities of motifs at  $t=0$ .

### 2.2. Substitution matrices of genetic motifs of size $n$

For substitution matrices of genetic motifs of size  $n$  containing up to three substitution parameters per motif site (extension of the 3-parameter substitution model (Kimura, 1981) of nucleotides to any motifs of size  $n$ ), a Kronecker property is identified for constructing these classes of substitution matrices. This property was found after a detailed analysis of the dinucleotide matrix  $\delta$  (Michel, 2007c, Fig. 1) and the trinucleotide matrix  $\delta$  (Michel, 2007b, Fig. B.1).

Let  $k$  be the nucleotide site of a genetic motif of size  $n$ ,  $1 \leq k \leq n$ . For a given site  $k$ , let  $a_k$ ,  $b_k$  and  $c_k$  be the parameter of transitions  $A \leftrightarrow G$  and  $C \leftrightarrow T$ , transversions  $A \leftrightarrow T$  and  $C \leftrightarrow G$  and transversions  $A \leftrightarrow C$  and  $G \leftrightarrow T$ , respectively. Thus, a motif of size  $n$  has  $3n$  substitution parameters.

The substitution matrix  $A_n$  ( $4^n, 4^n$ ) associated with motifs of size  $n$  is a block matrix which is classically constructed recursively by varying  $k=n$  to  $k=1$  as follows (Michel, 2007a, 2007b, 2007c):

$$A_k = \begin{pmatrix} A_{k-1} & c_{n-k+1}I_{k-1} & a_{n-k+1}I_{k-1} & b_{n-k+1}I_{k-1} \\ c_{n-k+1}I_{k-1} & A_{k-1} & b_{n-k+1}I_{k-1} & a_{n-k+1}I_{k-1} \\ a_{n-k+1}I_{k-1} & b_{n-k+1}I_{k-1} & A_{k-1} & c_{n-k+1}I_{k-1} \\ b_{n-k+1}I_{k-1} & a_{n-k+1}I_{k-1} & c_{n-k+1}I_{k-1} & A_{k-1} \end{pmatrix},$$

where  $I_{k-1}$  is the identity matrix ( $4^{k-1}, 4^{k-1}$ ) with  $I_0 = 1$ ,  $A_{k-1}$  is the recursive matrix ( $4^{k-1}, 4^{k-1}$ ) with  $A_0 = -\sum_{k=1}^n (a_k + b_k + c_k)$  and  $a_k, b_k, c_k$ ,  $1 \leq k \leq n$ , are the substitution parameters for the  $k$ th motif site. As the matrix  $A_n$  is real and symmetric,  $A_n$  is diagonalizable, i.e.  $A_n = Q_n \cdot D_n \cdot Q_n^{-1}$  where  $D_n$  is the spectral matrix of  $A_n$  and  $Q_n$  is its associated eigenvector matrix. This property allows the occurrence probabilities  $P_i(t)$  of motifs  $i$  to be determined, i.e. Eq. (2.6).

Let  $N_k$ ,  $1 \leq k \leq n$ , be the nucleotide substitution matrix (4,4) of a site  $k$  of a motif of size  $n$

$$N_k = \begin{pmatrix} d_k & c_k & a_k & b_k \\ c_k & d_k & b_k & a_k \\ a_k & b_k & d_k & c_k \\ b_k & a_k & c_k & d_k \end{pmatrix},$$

with  $d_k = -(a_k + b_k + c_k)$ . As the matrix  $N_k$  is real and symmetric,  $N_k$  is diagonalizable for all  $1 \leq k \leq n$

$$N_k = R \cdot S_k \cdot R^{-1},$$

where the nucleotide spectral matrix  $S_k$  of  $N_k$  is

$$S_k = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & -2(a_k + b_k) & 0 & 0 \\ 0 & 0 & -2(a_k + c_k) & 0 \\ 0 & 0 & 0 & -2(b_k + c_k) \end{pmatrix} \quad (2.7)$$

and its associated nucleotide eigenvectors matrix  $R$  is

$$R = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{pmatrix}. \quad (2.8)$$

**Remark 1.** For the substitution matrix of nucleotides ( $n=1$ ),  $A_1 = N_1 = Q_1 \cdot D_1 \cdot Q_1^{-1} = R \cdot S_1 \cdot R^{-1}$  leading to  $D_1 = S_1$  and  $Q_1 = R$ .

We identify an interesting relation between the matrices  $A_n$  and  $N_k$ .

Indeed, the recursive construction of the substitution matrix  $A_n$  is similar to a Kronecker sum of nucleotide substitution matrices  $N_k$  associated to each site  $k$  of motifs of size  $n$ . For a motif size  $n > 0$  and  $k > 0$ ,

$$A_k = \begin{pmatrix} A_{k-1} & c_{n-k+1}I_{k-1} & a_{n-k+1}I_{k-1} & b_{n-k+1}I_{k-1} \\ c_{n-k+1}I_{k-1} & A_{k-1} & b_{n-k+1}I_{k-1} & a_{n-k+1}I_{k-1} \\ a_{n-k+1}I_{k-1} & b_{n-k+1}I_{k-1} & A_{k-1} & c_{n-k+1}I_{k-1} \\ b_{n-k+1}I_{k-1} & a_{n-k+1}I_{k-1} & c_{n-k+1}I_{k-1} & A_{k-1} \end{pmatrix} \\ = \begin{pmatrix} 0 & c_{n-k+1} & a_{n-k+1} & b_{n-k+1} \\ c_{n-k+1} & 0 & b_{n-k+1} & a_{n-k+1} \\ a_{n-k+1} & b_{n-k+1} & 0 & c_{n-k+1} \\ b_{n-k+1} & a_{n-k+1} & c_{n-k+1} & 0 \end{pmatrix} \otimes I_{k-1} \\ + \begin{pmatrix} A_{k-1} & & & \\ & A_{k-1} & & \\ & & A_{k-1} & \\ & & & A_{k-1} \end{pmatrix},$$

where the diagonal block matrix with  $A_{k-1}$  on the main diagonal is a matrix ( $4^k, 4^k$ ). Then,

$$A_k = \begin{pmatrix} 0 & c_{n-k+1} & a_{n-k+1} & b_{n-k+1} \\ c_{n-k+1} & 0 & b_{n-k+1} & a_{n-k+1} \\ a_{n-k+1} & b_{n-k+1} & 0 & c_{n-k+1} \\ b_{n-k+1} & a_{n-k+1} & c_{n-k+1} & 0 \end{pmatrix} \otimes I_{k-1} + I_1 \otimes A_{k-1},$$

with  $I_1$  the identity matrix (4, 4). Therefore, by definition of the Kronecker sum,

$$A_k = \begin{pmatrix} 0 & c_{n-k+1} & a_{n-k+1} & b_{n-k+1} \\ c_{n-k+1} & 0 & b_{n-k+1} & a_{n-k+1} \\ a_{n-k+1} & b_{n-k+1} & 0 & c_{n-k+1} \\ b_{n-k+1} & a_{n-k+1} & c_{n-k+1} & 0 \end{pmatrix} \oplus A_{k-1}. \quad (2.9)$$

Moreover, by noticing that

$$N_{n-k+1} = \begin{pmatrix} 0 & c_{n-k+1} & a_{n-k+1} & b_{n-k+1} \\ c_{n-k+1} & 0 & b_{n-k+1} & a_{n-k+1} \\ a_{n-k+1} & b_{n-k+1} & 0 & c_{n-k+1} \\ b_{n-k+1} & a_{n-k+1} & c_{n-k+1} & 0 \end{pmatrix} \\ + \begin{pmatrix} d_{n-k+1} & 0 & 0 & 0 \\ 0 & d_{n-k+1} & 0 & 0 \\ 0 & 0 & d_{n-k+1} & 0 \\ 0 & 0 & 0 & d_{n-k+1} \end{pmatrix} \\ = \begin{pmatrix} 0 & c_{n-k+1} & a_{n-k+1} & b_{n-k+1} \\ c_{n-k+1} & 0 & b_{n-k+1} & a_{n-k+1} \\ a_{n-k+1} & b_{n-k+1} & 0 & c_{n-k+1} \\ b_{n-k+1} & a_{n-k+1} & c_{n-k+1} & 0 \end{pmatrix} + d_{n-k+1} \times I_1,$$

with  $d_{n-k+1} = -(a_{n-k+1} + b_{n-k+1} + c_{n-k+1})$ , we can rewrite the recursive Eq. (2.9) with the recursive Kronecker sum equation

$$A_k = (N_{n-k+1} - d_{n-k+1} \times I_1) \oplus A_{k-1}.$$

Thus,

$$A_n = \bigoplus_{k=1}^n (N_k - d_k \times I_1) + A_0 \times I_n \\ = \bigoplus_{k=1}^n N_k - \bigoplus_{k=1}^n (d_k \times I_1) + A_0 \times I_n \\ = \bigoplus_{k=1}^n N_k - \sum_{k=1}^n d_k \times I_n + A_0 \times I_n,$$

with  $A_0 = -\sum_{k=1}^n (a_k + b_k + c_k) = -\sum_{k=1}^n d_k$ , and finally,

$$A_n = \bigoplus_{k=1}^n N_k. \quad (2.10)$$

Appendix A illustrates this recurrence relation with a substitution matrix for dinucleotides.

Classical mathematical results (Laub, 2005) allows the spectral matrix  $D_n$  and the eigenvectors matrix  $Q_n$  to be deduced from  $S_k$  and  $R$ , respectively:

$$\begin{cases} D_n = \bigoplus_{k=1}^n S_k, \\ Q_n = \bigotimes_{k=1}^n R, \\ Q_n^{-1} = \left( \bigotimes_{k=1}^n R \right)^{-1} = \bigotimes_{k=1}^n R^{-1}. \end{cases}$$

Thus, the substitution matrix  $A_n$  can be directly determined from the Kronecker sum of the  $n$  nucleotide spectral matrices  $S_k$  and the Kronecker product of the  $n$  nucleotide eigenvectors matrix  $R$

as follows:

$$A_n = \bigotimes_{k=1}^n R \cdot \bigoplus_{k=1}^n S_k \cdot \bigotimes_{k=1}^n R^{-1}.$$

### 2.3. Analytical solutions giving the occurrence probabilities of genetic motifs of size $n$ at time $t$

By rewriting Eq. (2.6), the occurrence probability  $P_n(t)$  of motifs of size  $n$  at time  $t$  can be expressed as a function of elementary eigenvalues and eigenvectors associated with nucleotides of each site  $k$

$$\begin{aligned} P_n(t) &= \bigotimes_{k=1}^n R \cdot e^{\bigoplus_{k=1}^n S_k t} \cdot \bigotimes_{k=1}^n R^{-1} \cdot P_n(0) \\ &= \bigotimes_{k=1}^n R \cdot \bigotimes_{k=1}^n e^{S_k t} \cdot \bigotimes_{k=1}^n R^{-1} \cdot P_n(0) \\ &= \bigotimes_{k=1}^n (R \cdot e^{S_k t} \cdot R^{-1}) \cdot P_n(0), \end{aligned} \quad (2.11)$$

where  $e^{S_k t}$  is the exponential spectral matrix of matrices  $N_k$ .

**Proposition 1.** Eq. (2.11) gives the occurrence probability  $P_n(t)$  of motifs of size  $n$  at time  $t$  from its past one  $P_n(0)$ . If we express  $P_n(0)$  as a function of  $P_n(t)$  in Eq. (2.11) then equation

$$\tilde{P}_n(t) = \bigotimes_{k=1}^n (R \cdot e^{-S_k t} \cdot R^{-1}) \cdot \tilde{P}_n(0) \quad (2.12)$$

by replacing  $t$  by  $-t$  gives the past probability  $\tilde{P}_n(t)$  of motifs of size  $n$  from its current probability  $\tilde{P}_n(0)$ , i.e. by inverting the direction of the evolution time  $t$ .

### 2.4. Analytical solution giving the occurrence probability of a genetic motif of size $n$ at time $t$

Eq. (2.11) determines the analytical solutions for all the  $4^n$  motifs of size  $n$ . For  $n > 3$  (tetranucleotides, pentanucleotides, etc.), the resolution of this equation system on a current standard PC needs much time and memory. Therefore, for larger motifs, we have also derived an equation allowing to compute directly the occurrence probability  $P_{i_1}(t)$  of a genetic motif  $i_1$  of size  $n$  at time  $t$ . After some algebraic manipulation, we obtain

$$\begin{aligned} P_{i_1}(t) &= \frac{1}{4^n} \sum_{i_2=1}^{4^n} \left( e^{t \times \sum_{k=1}^n L_k[\delta(i_2, k)]} \times \sum_{i_3=1}^{4^n} (P_{i_3}(0) \right. \\ &\quad \left. \times \prod_{k=1}^n (R[\delta(i_1, k), \delta(i_2, k)] \times R[\delta(i_2, k), \delta(i_3, k)]) \right), \end{aligned} \quad (2.13)$$

where  $P_{i_3}(0)$  is the initial occurrence probability of motif  $i_3$  of size  $n$  at  $t=0$ ,  $R$  the nucleotide eigenvectors matrix (2.8),  $\delta(i_x, k) = \lfloor (i_x - 1) / 4^{n-k} \rfloor [4] + 1$ ,  $1 \leq \delta(i_x, k) \leq 4$ , a function associated with the motif  $i_x$  and the site  $k$ ,  $\lfloor x \rfloor$  is the integer part of  $x$ , and  $L_k = [0, -2(a_k + b_k), -2(a_k + c_k), -2(b_k + c_k)]$  is the vector of the four eigenvalues of the nucleotide substitution rates matrix  $N_k$  (see matrix (2.7)) associated with the site  $k$ . The detail of algebraic manipulation is given in Appendix B.

**Remark 2.** The  $4^n$  coefficients of initial occurrence probability  $P_{i_3}(0)$  in Eq. (2.13) are obtained by multiplying the coefficients  $R[\delta(i_1, k), \delta(i_2, k)] \times R[\delta(i_2, k), \delta(i_3, k)]$  of each nucleotide site  $k$ . Moreover, the  $4^n$  eigenvalues in Eq. (2.13) are obtained from the sum of the  $k$  eigenvalues of index  $\delta(i_2, k)$  associated with each nucleotide site  $k$ .

**Remark 3.** As  $R^{-1} = \frac{1}{4}R$  here, the coefficients  $R[\delta(i_1, k), \delta(i_2, k)] \times R[\delta(i_2, k), \delta(i_3, k)]$  in Eq. (2.13) can be seen as coefficients of matrices  $O_k$  defined in Lebre and Michel (2010) by  $O_k[i, j] = R[i, k] \times R^{-1}[k, j]$ .

**Remark 4.** Eq. (2.13) gives a full simplified analytical solution of the occurrence probability  $P_{i_1}(t)$  of a motif  $i_1$  of size  $n$  at time  $t$  composed of  $4^n$  exponents, each exponent being multiplied by a sum of the  $4^n$  initial occurrence probabilities  $P_{i_3}(0)$ , i.e. a total of  $4^{2n}$  terms. Thus, an analytical solution of a pentanucleotide, for example, has more than one million of terms.

**Example 1.** An example of application of Eq. (2.13) is given for determining the analytical solution of the occurrence probability of the dinucleotide AG at time  $t$ . The index of AG is  $i_1 = 3$  among  $4^2 = 16$  dinucleotides ( $n = 2$ ). The closed formula of  $P_3(t)$  obtained is given in Appendix C.

$$\begin{aligned} P_{i_1=3}(t) &= \frac{1}{16} \sum_{i_2=1}^{16} \left( e^{t \times \sum_{k=1}^2 L_k[\delta(i_2, k)]} \times \sum_{i_3=1}^{16} (P_{i_3}(0) \right. \\ &\quad \left. \times \prod_{k=1}^2 (R[\delta(3, k), \delta(i_2, k)] \times R[\delta(i_2, k), \delta(i_3, k)]) \right) \\ &= \frac{1}{16} \sum_{i_2=1}^{16} \left( e^{t \times (L_1[\delta(i_2, 1)] + L_2[\delta(i_2, 2)])} \right. \\ &\quad \left. \times \sum_{i_3=1}^{16} (P_{i_3}(0) \times (R[1, \delta(i_2, 1)] \times R[\delta(i_2, 1), \delta(i_3, 1)]) \right. \\ &\quad \left. \times (R[3, \delta(i_2, 2)] \times R[\delta(i_2, 2), \delta(i_3, 2)]) \right) \\ &= \frac{1}{16} \sum_{i_2=1}^{16} \left( e^{t \times (L_1[\delta(i_2, 1)] + L_2[\delta(i_2, 2)])} \times (P_1(0) \right. \\ &\quad \left. \times (R[1, \delta(i_2, 1)] \times R[\delta(i_2, 1), 1]) \times (R[3, \delta(i_2, 2)] \times R[\delta(i_2, 2), 1]) \right. \\ &\quad \left. + \dots + P_{16}(0) \times (R[1, \delta(i_2, 1)] \times R[\delta(i_2, 1), 4]) \right. \\ &\quad \left. \times (R[3, \delta(i_2, 2)] \times R[\delta(i_2, 2), 4]) \right) \\ &= \frac{1}{16} \left( e^{t \times (L_1[1] + L_2[1])} \times (P_1(0) \times (R[1, 1] \times R[1, 1]) \times (R[3, 1] \right. \\ &\quad \left. \times R[1, 1]) + \dots + P_{16}(0) \times (R[1, 1] \times R[1, 4]) \times (R[3, 1] \right. \\ &\quad \left. \times R[1, 4]) \right) + e^{t \times (L_1[1] + L_2[2])} \times (P_1(0) \times (R[1, 1] \times R[1, 1]) \\ &\quad \times (R[3, 2] \times R[2, 1]) + \dots + P_{16}(0) \times (R[1, 1] \times R[1, 4]) \\ &\quad \times (R[3, 2] \times R[2, 4]) + \dots + e^{t \times (L_1[4] + L_2[4])} \times (P_1(0) \times (R[1, 4] \\ &\quad \times R[4, 1]) \times (R[3, 4] \times R[4, 1]) \\ &\quad + \dots + P_{16}(0) \times (R[1, 4] \times R[4, 4]) \times (R[3, 4] \times R[4, 4])). \end{aligned}$$

## 3. Application: extension of the SEGM web server

### 3.1. Functionalities

The biomathematical model developed here allows to extend the SEGM (Stochastic Evolution of Genetic Motifs) web server (Benard and Michel, 2009) from trinucleotides to tetranucleotides and pentanucleotides, to improve the computation of analytical solutions with a faster calculus and also to add new functionalities, e.g. the result display. SEGM allows the determination of analytical occurrence probabilities  $P(t)$  of genetic motifs of size  $n$  (nucleotides to pentanucleotides) at time  $t$  as a function of substitution parameters  $a_k$  ( $A \leftrightarrow G$  and  $C \leftrightarrow T$ ),  $b_k$  ( $A \leftrightarrow T$  and  $C \leftrightarrow G$ ) and  $c_k$  ( $A \leftrightarrow C$  and  $G \leftrightarrow T$ ) per nucleotide site  $k$  and an initial occurrence probabilities  $P(0)$  of motifs at time  $t=0$ . The evolution direction can be direct (past-present) or inverse (present-past). The results are displayed according to several modes defined by the user: general analytical solutions, numerical solutions, evolution plots and analytical solutions converted in C, Fortran or T<sub>E</sub>X formats in order to facilitate their integration in user-programs. Fig. 1 gives the flowchart of the SEGM web server and an overview of its functionalities.



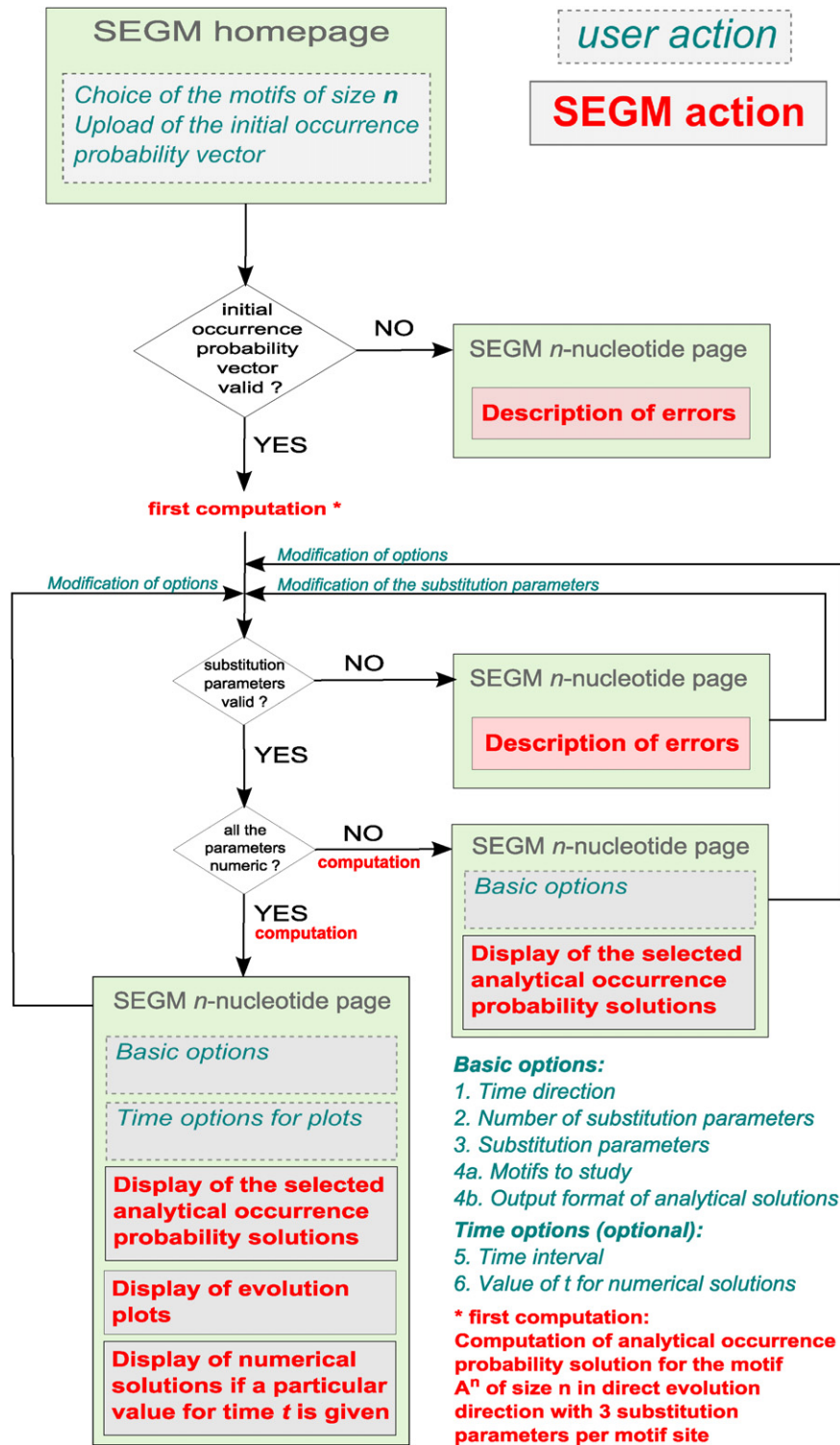


Fig. 1. Flowchart of the SEGM web server.

### 3.2. Size and initial occurrence probabilities of genetic motifs

The size  $n$  and the initial occurrence probability vector  $P(0)$  of studied motifs of size  $n$  at time  $t=0$  are chosen on the homepage of the SEGM (Fig. 2). The initial vector  $P(0)$  is uploaded thanks to a xls file containing the  $4^n$  initial occurrence probabilities of studied motifs of size  $n$  at time  $t=0$ . A link gives to the user the possibility of downloading a template for the xls file. After submission, the vector  $P(0)$  is checked by SEGM: its values must be numerical and positive,

and its sum must be equal to 1. If errors are detected, a description of these errors is given (Fig. 1). Otherwise, the corresponding  $n$ -nucleotide page is displayed and the result of a first computation using default options and parameters gives the analytical occurrence probability  $P(t)$  of the motif  $A^n$  of size  $n$  in the direct evolution direction and with three formal substitution parameters  $a_k$ ,  $b_k$  and  $c_k$  per motif site. The user can modify these default options and the parameters to get new results with the same initial vector  $P(0)$ . To study evolution of motifs of different size  $n$  or with different initial



# Stochastic Evolution of Genetic Motifs

Emmanuel Benard and Christian J. Michel

Theoretical Bioinformatics, LSiiT/CNRS UMR7005 - University of Strasbourg

1. Choose the motif size:

Dinucleotides (2) ▾

2. Upload the initial occurrence probabilities file:

Enter a XLS file containing the 16 initial occurrence probabilities of motifs of size 2:

Example of a valid XLS file containing 16 initial occurrence probabilities of Dinucleotides available [here](#)

Fig. 2. Homepage of the SEGM web server: choice of the size of the studied genetic motifs and upload of the initial occurrence probability vector of genetic motifs.

1. Evolutionary time sens:

Direct (past -> present) ▾

2. Number of substitution parameters per motif site:

3 parameters: 1 transition rate ( $A \leftrightarrow G = C \leftrightarrow T$ ), 1 transversion I rate ( $A \leftrightarrow T = C \leftrightarrow G$ ), 1 transversion II rate ( $A \leftrightarrow C = G \leftrightarrow T$ ).

2 parameters: 1 transition rate ( $A \leftrightarrow G = C \leftrightarrow T$ ), 1 transversion rate ( $A \leftrightarrow T = A \leftrightarrow C = C \leftrightarrow G = G \leftrightarrow T$ ).  
 $u[x]=a[x]$ ,  $v[x]/2=b[x]=c[x]$

1 parameter: 1 substitution rate ( $A \leftrightarrow C = A \leftrightarrow G = A \leftrightarrow T = C \leftrightarrow G = C \leftrightarrow T = G \leftrightarrow T$ ).  
 $p[x]/3=a[x]=b[x]=c[x]$

3 parameters ▾ [More about mutation matrices and substitution parameters](#)

Fig. 3. Option 1: Choice of the evolutionary time direction. Option 2: Choice of the number of substitution parameters.

probabilities  $P(0)$  of motifs, he must go back to the homepage of SEGM and upload a new xls file (Fig. 2).

**Example 2.** The choice of the initial vector  $P(0)$  depends on the evolutionary problem studied. For example, in order to study “primitive” dinucleotides at donor site, precisely their past occurrence probabilities  $P(t)$  by inverting the time direction, then the initial vector  $P(0)$  could be the 16 dinucleotides probabilities at donor site at current time, e.g. probabilities obtained from the ICE (Information for the Coordinates of Exons) database from current genes (see Benard and Michel, 2009, Table 2). Another example, suppose that a DNA sequence is a series of A, e.g. a poly(A) tail to an RNA molecule. Then, the 256 occurrence probabilities  $P(t)$  of tetranucleotides in this sequence subjected to substitutions can be studied with an initial vector  $P(0)$  associated to this sequence, i.e. precisely  $P_1(0) = 1$  for the motif AAAA and  $P_i(0) = 0$  for the 255 other motifs  $i \neq 1$ .

### 3.3. Basic options

#### 3.3.1. Time direction

After the submission of the initial occurrence probability vector  $P(0)$ , the first option is the choice of the evolutionary time direction (Fig. 3). The determination of analytical occurrence probabilities  $P(t)$  of motifs can be carried out in direct time direction (past-present) using Eq. (2.11) or in inverse time direction (present-past) using Eq. (2.12). By default, solutions are calculated in direct time direction.

#### 3.3.2. Number of substitution parameters per motif site

Option 2 permits to choose the number of substitution parameters per motif site (Fig. 3). The biomathematical model of SEGM is an extension of the 3-parameter substitution model (Kimura, 1981) of nucleotides to motifs based on three types of substitutions for each motif site  $k$ : transitions  $a_k$  ( $A \leftrightarrow G$  and  $C \leftrightarrow T$ ), transversions  $b_k$  ( $A \leftrightarrow T$  and  $C \leftrightarrow G$ ) and transversions  $c_k$

( $A \leftrightarrow C$  and  $G \leftrightarrow T$ ). It is the model by default. SEGM can also study particular cases extending the 2-parameter nucleotide substitution model (Kimura, 1980) and the 1-parameter nucleotide substitution model (Jukes and Cantor, 1969) to motifs. For the model of two substitution parameters per motif site, the parameters are the transitions  $u_k = a_k$  ( $A \leftrightarrow G$  and  $C \leftrightarrow T$ ) and the transversions  $v_k/2 = b_k = c_k$  ( $A \leftrightarrow C$ ,  $A \leftrightarrow T$ ,  $C \leftrightarrow G$  and  $G \leftrightarrow T$ ). For the model of one substitution parameter per motif site, the parameters are defined by  $p_k/3 = a_k = b_k = c_k$ . The different models and their associated substitution matrices are explained in a pdf file accessible in Option 2 (Fig. 3).

### 3.3.3. Substitution parameters

The substitution parameters can be left formal or set with numerical values in Option 3 (Fig. 4). Numerical values must be positive and their sum must be less or equal to 1. After submission, a description of the encountered errors is displayed if these conditions are not verified (Fig. 1). By default, the substitution parameters are left formal to derive formal analytical solutions.

**Remark 5.** The motif sites in SEGM are indexed from 0 to  $n-1$ ,  $n$  being the size of the studied motifs.

### 3.3.4. Choice of the studied motifs

Option 4a in this version of SEGM offers the possibility to study evolution up to four large motifs simultaneously. The user selects one motif per list (Fig. 5).

### 3.3.5. Output format of the analytical solutions

By default, the analytical occurrence probabilities  $P(t)$  of genetic motifs are displayed in a readable text format. In order to facilitate their integration in external user-programs, Option 4b allows others formats: C, Fortran and TeX (Fig. 5).

## 3.4. Results

The analytical occurrence probabilities  $P(t)$  of genetic motifs are given with formal substitution parameters (Fig. 6) in the first computation (Fig. 1) or when the substitution parameters

**3. Substitution parameters:**  
 Enter values for the substitution parameters.  
 Non numerical or rational values will be replaced by the name of the corresponding parameter.  
 All the substitution parameters must have a numerical value to get plots.  
 All the substitution parameters and their sum must be  $\geq 0$  and  $< 1$ .

Site 0: a[0]:  b[0]:  c[0]:

Site 1: a[1]:  b[1]:  c[1]:

**Substitution parameters statut:**  
 Parameters sum = "a0" + "a1" + "b0" + "b1" + "c0" + "c1"

Fig. 4. Option 3: Substitution parameters which can be formal or numerical.

**4. Choice of the probabilities to study and plot:**  
 Choose up to 4 analytical solutions.  
 By default, only the analytical solution of the motif AA is displayed and plotted.

motif AA - [-----] - [-----] - [-----]

**4b. Choice of the analytical solutions output format:**  
 The analytical solutions can be displayed in 4 formats: Standard, C, Fortran and TeX.  
 By default, the analytical solutions are displayed in Standard format.

Standard -

**SUBMIT**

Fig. 5. Option 4a: Choice of the studied motifs. Option 4b: Choice of the output format for the analytical occurrence probabilities  $P(t)$  of genetic motifs at time  $t$ .

## Results

### Analytical solutions (Standard format):

ProbAA(t)	$0.0625 - 0.0150523/E^{2(a_0 + b_0)t} - 0.0100126/E^{2(a_1 + b_1)t} - 0.0164957/E^{2(a_0 + a_1 + b_0 + b_1)t} - 0.0164097/E^{2(a_0 + c_0)t} + 0.0104875/E^{2(b_0 + c_0)t} + 0.00390817/E^{2(a_0 + a_1 + b_1 + c_0)t} + 0.00302624/E^{2(a_1 + b_0 + b_1 + c_0)t} + 0.00532176/E^{2(a_1 + c_1)t} + 0.00707809/E^{2(a_0 + a_1 + b_0 + c_1)t} + 0.00429478/E^{2(b_1 + c_1)t} - 0.00876884/E^{2(a_0 + b_0 + b_1 + c_1)t} + 0.00763192/E^{2(a_0 + a_1 + c_0 + c_1)t} + 0.00332749/E^{2(a_1 + b_0 + c_0 + c_1)t} - 0.0103474/E^{2(a_0 + b_1 + c_0 + c_1)t} - 0.00828515/E^{2(b_0 + b_1 + c_0 + c_1)t}$
-----------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Fig. 6. Example of analytical occurrence probability  $P_1(t)$  of the dinucleotide AA at time  $t$  with three formal substitution parameters per site and a particular initial occurrence probability vector of genetic motifs.

are left formal (Fig. 4). When all the substitution parameters are set with numerical values (Option 3, Fig. 4), evolution plots of the chosen motifs (Option 4a, Fig. 5) are displayed: a comparative evolution plot containing the evolution curves of the chosen motifs and a plot drawing their sum curve (Fig. 7). In this case, two additional options regarding evolution time  $t$  are available.

The time interval of plots (Fig. 7) can be modified thanks to Option 5 (Fig. 8). The times  $t_{min}$  and  $t_{max}$  must be always positive even if the inverse time direction is chosen. By default, the time  $t$  varies from  $t_{min}=0$  to  $t_{max}=5$  whatever the time direction chosen.

A particular value of time  $t$  can be set in Option 6 (Fig. 8) to have the values of the occurrence probabilities  $P(t)$  of studied motifs (Fig. 9). As for the evolution plots, a value of the probability sum of studied motifs is also given.

#### 4. Discussion

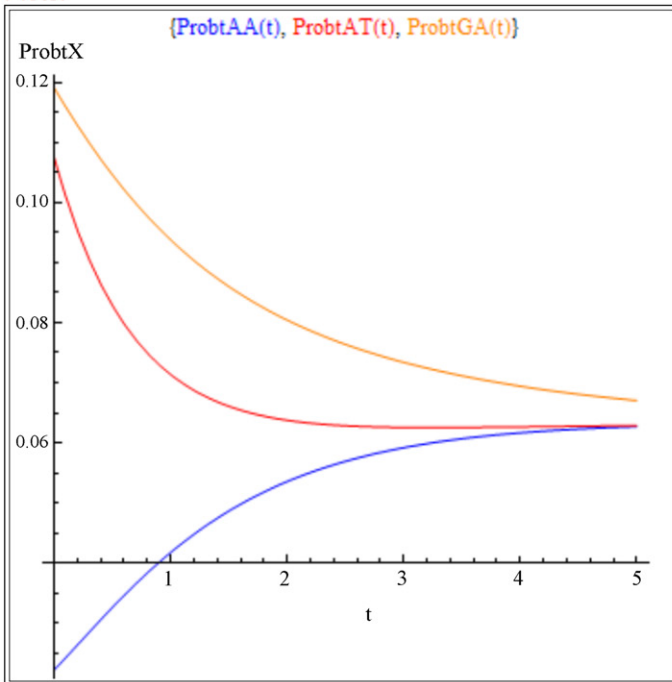
We proposed here a generalization of stochastic substitution models of nucleotides to genetic motifs of any size. This generalized

##### Numerical solutions:

ProbtAA(t) with t = 2.5	0.0569366
ProbtAT(t) with t = 2.5	0.0628471
ProbtGA(t) with t = 2.5	0.0764011
Sum: ProbtAA(t) + ProbtAT(t) + ProbtGA(t) with t = 2.5	0.196185

Fig. 9. Example of numerical solutions for the three dinucleotides AA, AT and GA and their sum at time  $t=2.5$  for particular substitution parameters and a given initial occurrence probability vector of genetic motifs.

##### Plots:



##### Plot sum:

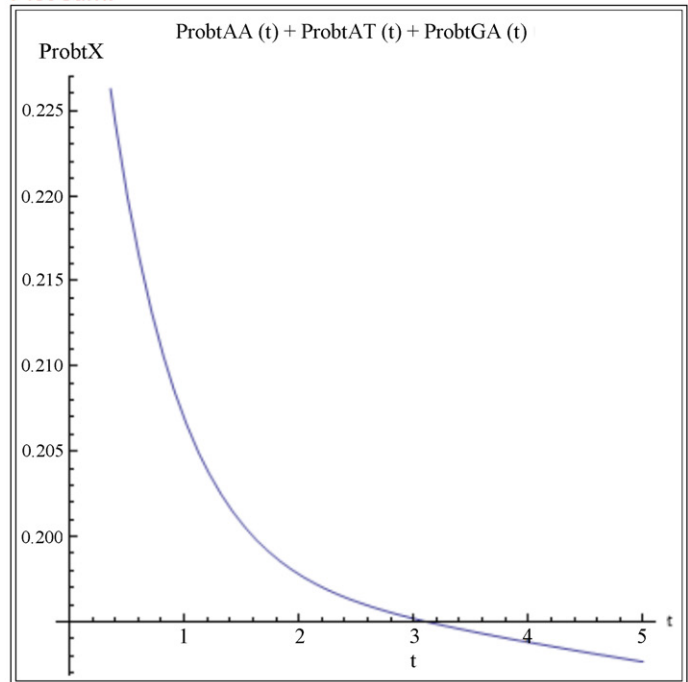


Fig. 7. Example of evolution curves for the three dinucleotides AA, AT and GA (left) and their sum (right).

5. Time interval for plots (optional):  
 $0 \leq t_{min} < t_{max}$

t min:  t max:

---

6. Time value for numerical solutions (optional):  
 $0 \leq t$

t:

Fig. 8. Option 5: Choice of the time interval for plots. Option 6: Choice of a time value for numerical solutions.



model is based on the identification of a Kronecker relation between the nucleotide substitution matrices and the motif substitution matrices. It gives the analytical occurrence probabilities of genetic motifs as a function of a substitution matrix containing up to three formal parameters (substitution rates) per motif site and an initial occurrence probabilities of genetic motifs. This biomathematical model was included in a new version of the SEGM web server offering now several improved functionalities, in particular a faster computation, and an extension to genetic motifs of size 4 (tetranucleotides) and 5 (pentanucleotides). The current limit of the computer implementation of this model is the PC power (CPU and memory). We are currently investigating a parallel approach based on GPU (Mathematica 8.1 integrates GPU programming) in order to allow the evolution analysis of genetic motifs greater than five nucleotides.

### Acknowledgments

We thank the reviewers for their advice.

### Appendix A. Substitution matrix for dinucleotides with the Kronecker method

Construction of the dinucleotide substitution matrix  $A_2$  (16,16) from the Kronecker sum of the two matrices  $N_1$  and  $N_2$  of size (4,4) associated to the nucleotide substitution matrices at dinucleotide sites 1 and 2, respectively, using Eq. (2.10):

$$N_1 \oplus N_2 = \begin{pmatrix} d_1 & c_1 & a_1 & b_1 \\ c_1 & d_1 & b_1 & a_1 \\ a_1 & b_1 & d_1 & c_1 \\ b_1 & a_1 & c_1 & d_1 \end{pmatrix} \oplus \begin{pmatrix} d_2 & c_2 & a_2 & b_2 \\ c_2 & d_2 & b_2 & a_2 \\ a_2 & b_2 & d_2 & c_2 \\ b_2 & a_2 & c_2 & d_2 \end{pmatrix},$$

with  $d_1 = -(a_1 + b_1 + c_1)$  and  $d_2 = -(a_2 + b_2 + c_2)$ . Then,

$$N_1 \oplus N_2 = \begin{pmatrix} d_1 & c_1 & a_1 & b_1 \\ c_1 & d_1 & b_1 & a_1 \\ a_1 & b_1 & d_1 & c_1 \\ b_1 & a_1 & c_1 & d_1 \end{pmatrix} \otimes \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} + \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \otimes \begin{pmatrix} d_2 & c_2 & a_2 & b_2 \\ c_2 & d_2 & b_2 & a_2 \\ a_2 & b_2 & d_2 & c_2 \\ b_2 & a_2 & c_2 & d_2 \end{pmatrix} = \begin{pmatrix} d_1 I_1 & c_1 I_1 & a_1 I_1 & b_1 I_1 \\ c_1 I_1 & d_1 I_1 & b_1 I_1 & a_1 I_1 \\ a_1 I_1 & b_1 I_1 & d_1 I_1 & c_1 I_1 \\ b_1 I_1 & a_1 I_1 & c_1 I_1 & d_1 I_1 \end{pmatrix} + \begin{pmatrix} N_2 & 0 & 0 & 0 \\ 0 & N_2 & 0 & 0 \\ 0 & 0 & N_2 & 0 \\ 0 & 0 & 0 & N_2 \end{pmatrix},$$

where  $I_1$  is the identity matrix (4,4). Then,

$$N_1 \oplus N_2 = \begin{pmatrix} d_1 I_1 + N_2 & c_1 I_1 & a_1 I_1 & b_1 I_1 \\ c_1 I_1 & d_1 I_1 + N_2 & b_1 I_1 & a_1 I_1 \\ a_1 I_1 & b_1 I_1 & d_1 I_1 + N_2 & c_1 I_1 \\ b_1 I_1 & a_1 I_1 & c_1 I_1 & d_1 I_1 + N_2 \end{pmatrix}$$

$$= \begin{pmatrix} d & c_2 & a_2 & b_2 & c_1 & 0 & 0 & 0 & a_1 & 0 & 0 & 0 & b_1 & 0 & 0 & 0 \\ c_2 & d & b_2 & a_2 & 0 & c_1 & 0 & 0 & 0 & a_1 & 0 & 0 & 0 & b_1 & 0 & 0 \\ a_2 & b_2 & d & c_2 & 0 & 0 & c_1 & 0 & 0 & 0 & a_1 & 0 & 0 & 0 & b_1 & 0 \\ b_2 & a_2 & c_2 & d & 0 & 0 & 0 & c_1 & 0 & 0 & 0 & a_1 & 0 & 0 & 0 & b_1 \\ c_1 & 0 & 0 & 0 & d & c_2 & a_2 & b_2 & b_1 & 0 & 0 & 0 & a_1 & 0 & 0 & 0 \\ 0 & c_1 & 0 & 0 & c_2 & d & b_2 & a_2 & 0 & b_1 & 0 & 0 & 0 & a_1 & 0 & 0 \\ 0 & 0 & c_1 & 0 & a_2 & b_2 & d & c_2 & 0 & 0 & b_1 & 0 & 0 & 0 & a_1 & 0 \\ 0 & 0 & 0 & c_1 & b_2 & a_2 & c_2 & d & 0 & 0 & 0 & b_1 & 0 & 0 & 0 & a_1 \\ a_1 & 0 & 0 & 0 & b_1 & 0 & 0 & 0 & d & c_2 & a_2 & b_2 & c_1 & 0 & 0 & 0 \\ 0 & a_1 & 0 & 0 & 0 & b_1 & 0 & 0 & c_2 & d & b_2 & a_2 & 0 & c_1 & 0 & 0 \\ 0 & 0 & a_1 & 0 & 0 & 0 & b_1 & 0 & a_2 & b_2 & d & c_2 & 0 & 0 & c_1 & 0 \\ 0 & 0 & 0 & a_1 & 0 & 0 & 0 & b_1 & c_2 & a_2 & c_2 & d & 0 & 0 & 0 & c_1 \\ b_1 & 0 & 0 & 0 & a_1 & 0 & 0 & 0 & c_1 & 0 & 0 & 0 & d & c_2 & a_2 & b_2 \\ 0 & b_1 & 0 & 0 & 0 & a_1 & 0 & 0 & 0 & c_1 & 0 & 0 & c_2 & d & b_2 & a_2 \\ 0 & 0 & b_1 & 0 & 0 & 0 & a_1 & 0 & 0 & 0 & c_1 & 0 & a_2 & b_2 & d & c_2 \\ 0 & 0 & 0 & b_1 & 0 & 0 & 0 & a_1 & 0 & 0 & 0 & c_1 & b_2 & a_2 & c_2 & d \end{pmatrix} = A_2,$$

with  $d = -(a_1 + b_1 + c_1 + a_2 + b_2 + c_2)$ .

### Appendix B. Detailed derivation of the analytical occurrence probability of a motif $i_1$ of size $n$ at time $t$

In order to avoid the computation of the occurrence probabilities of the  $4^n$  motifs of size  $n$ , we also derived from Eq. (2.11) a formula that determines the analytical occurrence probability  $P_{i_1}(t)$  at time  $t$  of a motif  $i_1$  of size  $n$ . Indeed,

$$P_{i_1}(t) = P_n(t)[i_1] = \left( \bigotimes_{k=1}^n (R \cdot e^{S_k t} \cdot R^{-1}) \cdot P_n(0) \right) [i_1] = \left( \bigotimes_{k=1}^n R \right) [i_1] \cdot \left( \bigotimes_{k=1}^n e^{S_k t} \right) \cdot \left( \bigotimes_{k=1}^n R^{-1} \right) \cdot P_n(0). \quad (B.1)$$

Thus, the computation of the analytical occurrence probability of the motif  $i_1$  does not require the construction of the whole eigenvectors matrix  $Q = \bigotimes_{k=1}^n R$  but only the  $i_1$ th row of  $Q[i_1, \cdot] = \left( \bigotimes_{k=1}^n R \right) [i_1, \cdot]$ .

Moreover, the matrix  $Q(4^n, 4^n)$  is constructed from a Kronecker product of  $n$  matrices  $R(4,4)$ . Then, the  $i_1$ th row of  $Q$  can be determined from a Kronecker product of  $n$  rows of  $R$  of indexes corresponding to the indexes  $\delta(i_1, k)$ ,  $1 \leq \delta(i_1, k) \leq 4$ , of nucleotides of each site  $k$ ,  $1 \leq k \leq n$ , in the motif  $i_1$ . For example, the row of  $Q$  corresponding to the occurrence probability of the motif GCT of size  $n=3$  and of index  $i_1=40$  is the result of the Kronecker product of the following three rows of  $R$  of indexes  $\delta(i_1, k)$ : row  $\delta(i_1, 1)=3$  corresponding to the nucleotide G for the site  $k=1$ , row  $\delta(i_1, 2)=2$  corresponding to the nucleotide C for the site  $k=2$  and row  $\delta(i_1, 3)=4$  corresponding to the nucleotide T for the site  $k=3$ . Thus, the corresponding row of index  $i_1=40$  in  $Q$  is obtained by  $Q[40, \cdot] = R[\delta(i_1, 1), \cdot] \otimes R[\delta(i_1, 2), \cdot] \otimes R[\delta(i_1, 3), \cdot] = R[3, \cdot] \otimes R[2, \cdot] \otimes R[4, \cdot]$ . The indexes  $\delta(i_1, k)$  of nucleotides of each site  $k$  of a motif  $i_1$  of size  $n$  are calculated by the formula

$$\delta(i_1, k) = \left\lfloor \frac{i_1 - 1}{4^{n-k}} \right\rfloor [4] + 1,$$

where  $\lfloor (i_1 - 1)/4^{n-k} \rfloor$  is the integer part of  $(i_1 - 1)/4^{n-k}$  and  $[ \cdot ]$  is the modulo function. As an illustration with the previous example, i.e. the motif GCT of size  $n=3$  and index  $i_1=40$ ,  $\delta(i_1, 1) = \lfloor (40 - 1)/4^{3-1} \rfloor [4] + 1 = 3$ .

Thus, Eq. (B.1) can be rewritten

$$P_{i_1}(t) = \left( \bigotimes_{k=1}^n R[\delta(i_1, k), \cdot] \right) \cdot \left( \bigotimes_{k=1}^n e^{S_k t} \right) \cdot \left( \bigotimes_{k=1}^n R^{-1} \right) \cdot P_n(0). \quad (B.2)$$

Eq. (B.2) can still be simplified. Indeed, the matrix product  $\left( \bigotimes_{k=1}^n R[\delta(i_1, k), \cdot] \right) \cdot \left( \bigotimes_{k=1}^n e^{S_k t} \right)$  uses a row vector  $Q[i_1, \cdot]$  and a

diagonal matrix  $e^{S^t}$  whose diagonal elements are exponents of eigenvalues of the substitution rates matrix  $A_n$  ( $4^n, 4^n$ ). This matrix  $e^{S^t}$  can also be obtained by a Kronecker product of  $n$  matrices  $e^{S_k t}$ . The matrix product  $Q[i_1, \cdot] \cdot e^{S^t}$  can then be replaced by a scalar product between the row vector  $Q[i_1, \cdot]$  and a row vector  $e^{L^t}$  composed of the diagonal elements of  $e^{S^t}$ . Let  $L_k$  be such a row vector associated with the site  $k$  of a motif of size  $n$  and containing the diagonal elements of  $S_k$ , i.e. the eigenvalues of the nucleotide substitution matrix  $N_k$  associated with a motif site  $k$ , i.e.  $L_k = [0, -2(a_k + b_k), -2(a_k + c_k), -2(b_k + c_k)]$  (see matrix (2.7)). By using the row vector  $e^{L_k t}$ , Eq. (B.2) can be rewritten

$$P_{i_1}(t) = \left( \bigotimes_{k=1}^n R[\delta(i_1, k)] \right) \times \left( \bigotimes_{k=1}^n e^{L_k t} \right) \cdot \left( \bigotimes_{k=1}^n R^{-1} \right) \cdot P_n(0). \quad (\text{B.3})$$

As  $\bigotimes_{k=1}^n e^{L_k t}$  is also a row vector, Eq. (B.3) divides by  $4^n$  the number of operations of Eq. (B.2).

Let  $U_{i_1}$  be the row vector associated with the motif  $i_1$  and defined by  $U_{i_1} = (\bigotimes_{k=1}^n R[\delta(i_1, k)]) \times (\bigotimes_{k=1}^n e^{L_k t})$  with its  $i$ th element  $U_{i_1}[i] = (\bigotimes_{k=1}^n R)[i_1, i] \times e^{L_k t}$  where  $e^{L_k t}$  is the  $i$ th element of the row vector  $\bigotimes_{k=1}^n e^{L_k t}$ . Let  $V_{i_1}$  be the row vector associated with the motif  $i_1$  and defined by  $V_{i_1} = U_{i_1} \cdot (\bigotimes_{k=1}^n R^{-1})$  with its  $i$ th element  $V_{i_1}[i] = \sum_{j=1}^{4^n} U_{i_1}[j] \times (\bigotimes_{k=1}^n R^{-1})[j, i]$ . Then,

$$V_{i_1}[i] = \sum_{j=1}^{4^n} \left( \bigotimes_{k=1}^n R \right) [i_1, j] \times e^{L_k t} \times \left( \bigotimes_{k=1}^n R^{-1} \right) [j, i].$$

From Eq. (B.3),

$$\begin{aligned} P_{i_1}(t) &= V_{i_1} \cdot P_n(0) = \sum_{i_2=1}^{4^n} V_{i_1}[i_2] \times P_{i_2}(0) \\ &= \sum_{i_2=1}^{4^n} \left( \sum_{i_3=1}^{4^n} \left( \bigotimes_{k=1}^n R \right) [i_1, i_3] \times e^{L_k t} \left( \bigotimes_{k=1}^n R^{-1} \right) [i_3, i_2] \right) \times P_{i_2}(0). \end{aligned}$$

By expressing  $P_{i_1}(t)$  as a sum of  $4^n$  exponents of eigenvalues  $e^{L_{i_2} t}$ , each one associated with a sum of  $4^n$  initial occurrence probabilities, we obtain

$$\begin{aligned} P_{i_1}(t) &= \sum_{i_2=1}^{4^n} e^{L_{i_2} t} \times \sum_{i_3=1}^{4^n} \left( P_{i_3}(0) \times \left( \bigotimes_{k=1}^n R \right) [i_1, i_2] \times \left( \bigotimes_{k=1}^n R^{-1} \right) [i_2, i_3] \right) \\ &= \sum_{i_2=1}^{4^n} e^{t \times \sum_{k=1}^n L_k[\delta(i_2, k)]} \times \sum_{i_3=1}^{4^n} \left( P_{i_3}(0) \times \prod_{k=1}^n (R[\delta(i_1, k), \delta(i_2, k)]) \right. \\ &\quad \left. \times \prod_{k=1}^n (R^{-1}[\delta(i_2, k), \delta(i_3, k)]) \right) \\ &= \sum_{i_2=1}^{4^n} e^{t \times \sum_{k=1}^n L_k[\delta(i_2, k)]} \times \sum_{i_3=1}^{4^n} \left( P_{i_3}(0) \times \prod_{k=1}^n (R[\delta(i_1, k), \delta(i_2, k)]) \right. \\ &\quad \left. \times R^{-1}[\delta(i_2, k), \delta(i_3, k)] \right), \quad (\text{B.4}) \end{aligned}$$

with  $\delta(i_x, k) = \lfloor (i_x - 1) / 4^{n-k} \rfloor [4] + 1$ . As  $R^{-1} = \frac{1}{4} R$  (Remark 3), Eq. (B.4) simplifies

$$\begin{aligned} P_{i_1}(t) &= \frac{1}{4^n} \sum_{i_2=1}^{4^n} e^{t \times \sum_{k=1}^n L_k[\delta(i_2, k)]} \\ &\quad \times \sum_{i_3=1}^{4^n} \left( P_{i_3}(0) \times \prod_{k=1}^n (R[\delta(i_1, k), \delta(i_2, k)] \times R[\delta(i_2, k), \delta(i_3, k)]) \right). \end{aligned}$$

## Appendix C. Analytical solution of the occurrence probability of the dinucleotide AG at time $t$

$$\begin{aligned} P_3(t) &= \frac{1}{16} [e^0 (P_1(0) + P_2(0) + P_3(0) + P_4(0) + P_5(0) + P_6(0) + P_7(0) \\ &\quad + P_8(0) + P_9(0) + P_{10}(0) + P_{11}(0) + P_{12}(0) + P_{13}(0) + P_{14}(0) + P_{15}(0) \\ &\quad + P_{16}(0)) \\ &\quad + e^{-2(a_2 + b_2)t} (-P_1(0) - P_2(0) + P_3(0) + P_4(0) - P_5(0) - P_6(0) + P_7(0) \\ &\quad + P_8(0) - P_9(0) - P_{10}(0) + P_{11}(0) + P_{12}(0) - P_{13}(0) - P_{14}(0) + P_{15}(0) \\ &\quad + P_{16}(0)) \\ &\quad + e^{-2(a_2 + c_2)t} (-P_1(0) + P_2(0) + P_3(0) - P_4(0) - P_5(0) + P_6(0) + P_7(0) \\ &\quad - P_8(0) - P_9(0) + P_{10}(0) + P_{11}(0) - P_{12}(0) - P_{13}(0) + P_{14}(0) + P_{15}(0) \\ &\quad - P_{16}(0)) \\ &\quad + e^{-2(b_2 + c_2)t} (P_1(0) - P_2(0) + P_3(0) - P_4(0) + P_5(0) - P_6(0) + P_7(0) \\ &\quad - P_8(0) + P_9(0) - P_{10}(0) + P_{11}(0) - P_{12}(0) + P_{13}(0) - P_{14}(0) + P_{15}(0) \\ &\quad - P_{16}(0)) \\ &\quad + e^{-2(a_1 + b_1)t} (P_1(0) + P_2(0) + P_3(0) + P_4(0) + P_5(0) + P_6(0) + P_7(0) \\ &\quad + P_8(0) - P_9(0) - P_{10}(0) - P_{11}(0) - P_{12}(0) - P_{13}(0) - P_{14}(0) - P_{15}(0) - P_{16}(0)) \\ &\quad + e^{-2(a_1 + b_1 + a_2 + b_2)t} (-P_1(0) - P_2(0) + P_3(0) + P_4(0) - P_5(0) - P_6(0) \\ &\quad + P_7(0) + P_8(0) + P_9(0) + P_{10}(0) - P_{11}(0) - P_{12}(0) + P_{13}(0) + P_{14}(0) \\ &\quad - P_{15}(0) - P_{16}(0)) \\ &\quad + e^{-2(a_1 + b_1 + a_2 + c_2)t} (-P_1(0) + P_2(0) + P_3(0) - P_4(0) - P_5(0) + P_6(0) \\ &\quad + P_7(0) - P_8(0) + P_9(0) - P_{10}(0) - P_{11}(0) + P_{12}(0) + P_{13}(0) - P_{14}(0) \\ &\quad - P_{15}(0) + P_{16}(0)) \\ &\quad + e^{-2(a_1 + b_1 + b_2 + c_2)t} (P_1(0) - P_2(0) + P_3(0) - P_4(0) + P_5(0) - P_6(0) \\ &\quad + P_7(0) - P_8(0) - P_9(0) + P_{10}(0) - P_{11}(0) + P_{12}(0) - P_{13}(0) + P_{14}(0) \\ &\quad - P_{15}(0) + P_{16}(0)) \\ &\quad + e^{-2(a_1 + c_1)t} (P_1(0) + P_2(0) + P_3(0) + P_4(0) - P_5(0) - P_6(0) - P_7(0) \\ &\quad - P_8(0) - P_9(0) - P_{10}(0) - P_{11}(0) - P_{12}(0) + P_{13}(0) + P_{14}(0) + P_{15}(0) \\ &\quad + P_{16}(0)) \\ &\quad + e^{-2(a_1 + c_1 + a_2 + b_2)t} (-P_1(0) - P_2(0) + P_3(0) + P_4(0) + P_5(0) + P_6(0) \\ &\quad - P_7(0) - P_8(0) + P_9(0) + P_{10}(0) - P_{11}(0) - P_{12}(0) - P_{13}(0) - P_{14}(0) \\ &\quad + P_{15}(0) + P_{16}(0)) \\ &\quad + e^{-2(a_1 + c_1 + a_2 + c_2)t} (-P_1(0) + P_2(0) + P_3(0) - P_4(0) + P_5(0) - P_6(0) \\ &\quad - P_7(0) + P_8(0) + P_9(0) - P_{10}(0) - P_{11}(0) + P_{12}(0) - P_{13}(0) + P_{14}(0) \\ &\quad + P_{15}(0) - P_{16}(0)) \\ &\quad + e^{-2(a_1 + c_1 + b_2 + c_2)t} (P_1(0) - P_2(0) + P_3(0) - P_4(0) - P_5(0) + P_6(0) - P_7(0) \\ &\quad + P_8(0) - P_9(0) + P_{10}(0) - P_{11}(0) + P_{12}(0) + P_{13}(0) - P_{14}(0) \\ &\quad + P_{15}(0) - P_{16}(0)) \\ &\quad + e^{-2(b_1 + c_1)t} (P_1(0) + P_2(0) + P_3(0) + P_4(0) - P_5(0) - P_6(0) - P_7(0) \\ &\quad - P_8(0) + P_9(0) + P_{10}(0) + P_{11}(0) + P_{12}(0) - P_{13}(0) - P_{14}(0) - P_{15}(0) \\ &\quad - P_{16}(0)) \\ &\quad + e^{-2(b_1 + c_1 + a_2 + b_2)t} (-P_1(0) - P_2(0) + P_3(0) + P_4(0) + P_5(0) + P_6(0) \\ &\quad - P_7(0) - P_8(0) - P_9(0) - P_{10}(0) + P_{11}(0) + P_{12}(0) + P_{13}(0) + P_{14}(0) \\ &\quad - P_{15}(0) - P_{16}(0)) \\ &\quad + e^{-2(b_1 + c_1 + a_2 + c_2)t} (-P_1(0) + P_2(0) + P_3(0) - P_4(0) + P_5(0) - P_6(0) \\ &\quad - P_7(0) + P_8(0) - P_9(0) + P_{10}(0) + P_{11}(0) - P_{12}(0) + P_{13}(0) - P_{14}(0) \\ &\quad - P_{15}(0) + P_{16}(0)) \\ &\quad + e^{-2(b_1 + c_1 + b_2 + c_2)t} (P_1(0) - P_2(0) + P_3(0) - P_4(0) - P_5(0) + P_6(0) - P_7(0) \\ &\quad + P_8(0) + P_9(0) - P_{10}(0) + P_{11}(0) - P_{12}(0) - P_{13}(0) + P_{14}(0) - P_{15}(0) + P_{16}(0))]. \end{aligned}$$

## References

- Allman, E.S., Rhodes, J.A., 2009. The identifiability of covarion models in phylogenetics. *IEEE/ACM Trans. Comput. Biol. Bioinformatics* 6, 76–88.
- Arquès, D.G., Michel, C.J., 1993. Analytical expression of the purine/pyrimidine codon probability after and before random mutations. *Bull. Math. Biol.* 55, 1025–1038.
- Arquès, D.G., Michel, C.J., 1995. Analytical solutions of the dinucleotide probability after and before random mutations. *J. Theor. Biol.* 175, 533–544.

- Benard, E., Michel, C.J., 2009. Computation of direct and inverse mutations with the SEGM web server (Stochastic Evolution of Genetic Motifs): an application to splice sites of human genome introns. *Comput. Biol. Chem.* 33, 245–252.
- Bryant, D., 2009. Hadamard phylogenetic methods and the n-taxon process. *Bull. Math. Biol.* 71, 339–351.
- Burman, D.Y., Smith, D.R., 1986. An asymptotic analysis of a queueing system with Markov-modulated arrivals. *Oper. Res.* 34, 105–119.
- Darot, J., Yeang, C.-H., Haussler, D., 2006. Detecting the dependent evolution of biosequences. *Res. Comput. Mol. Biol. Lect. Notes Comput. Sci* 3909, 595–609.
- Galtier, N., Jean-Marie, A., 2004. Markov-modulated markov chains and the covarion process of molecular evolution. *J. Comput. Biol.* 11, 727–733.
- Harrell, A., Melcher, U., Atkins, J.F., 2002. Predominance of six different hexanucleotide recoding signals 3' of read-through stop codons. *Nucleic Acids Res.* 30, 2011–2017.
- Jukes, T.H., Cantor, C.R., 1969. Evolution of protein molecules. In: Munro, H.N. (Ed.), *Mammalian Protein Metabolism*. Academic Press, New York, pp. 21–132.
- Kimura, M., 1980. A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences. *J. Mol. Evol.* 16, 111–120.
- Kimura, M., 1981. Estimation of evolutionary distances between homologous nucleotide sequences. *Proc. Natl. Acad. Sci. USA* 78, 454–458.
- Lange, K., 2005. *Applied Probability*. Springer-Verlag, New York.
- Laub, A.J., 2005. *Matrix Analysis for Scientists and Engineers*. SIAM Publications, Philadelphia, PA.
- Lebre, S., Michel, C.J., 2010. A stochastic evolution model for residue insertion-deletion independent from substitution. *Comput. Biol. Chem.* 34, 259–267.
- Lunter, G., Hein, J., 2004. A nucleotide substitution model with nearest-neighbour interactions. *Bioinformatics* 20, i216–223.
- Mackiewicz, P., Biecek, P., Mackiewicz, D., Kiraga, J., Baczkowski, K., Sobczynski, M., Cebzat, S., 2008. Optimisation of asymmetric mutational pressure and selection pressure around the universal genetic code. In: Bubak, M., et al. (Eds.), *ICCS 2008, Part III, LNCS 5103*, pp. 100–109.
- Michel, C.J., 2007a. An analytical model of gene evolution with 9 mutation parameters: an application to the amino acids coded by the common circular code. *Bull. Math. Biol.* 69, 677–698.
- Michel, C.J., 2007b. Codon phylogenetic distance. *Comput. Biol. Chem.* 31, 36–43.
- Michel, C.J., 2007c. Evolution probabilities and phylogenetic distance of dinucleotides. *J. Theor. Biol.* 249, 271–277.
- Salvador, P., Valadas, R., Pacheco, A., 2003. Multiscale fitting procedure using Markov modulated Poisson processes. *Telecommun. Syst.* 23, 123–148.
- Tian, Y., Styán, G.P.H., 2001. How to establish block-matrix factorizations. *Electron. J. Linear Algebra* 8, 115–127.
- Wang, H.-C., Spencer, M., Susko, E., Roger, A.J., 2007. Testing for covarion-like evolution in protein sequences. *Mol. Biol. Evol.* 24, 294–305.
- Zharkikh, A., 1994. Estimation of evolutionary distances between nucleotide sequences. *J. Mol. Evol.* 39, 315–329.