

## Alignment-free Prediction of Ribonucleases using a Computational Chemistry approach: Comparison with HMM model and Isolation from *Schizosaccharomyces pombe*, Prediction, and Experimental assay of a new sequence

**Guillermo Agüero-Chapín**<sup>1,2</sup>, Humberto González-Díaz,<sup>1,3,4</sup> Gustavo de la Riva,<sup>5</sup> Edrey Rodríguez,<sup>6</sup> Amina Sánchez-Rodríguez,<sup>2</sup> Gianni Podda,<sup>1</sup> Roberto I. Vazquez-Padrón<sup>7</sup>

<sup>1</sup> Dipartimento Farmaco Chimico Tecnologico, Università Degli Studi di Cagliari, 09124, Italy

<sup>2</sup> CAP, Faculty of Chemistry and Pharmacy, IBP, and CBQ, UCLV, Santa Clara, 54830, Cuba

<sup>3</sup> Unit for Bioinformatics & Connectivity Analysis (UBICA), Institute of Industrial Pharmacy, Faculty of Pharmacy, USC, Santiago de Compostela, 15782, Spain

<sup>4</sup> Department of Organic Chemistry, Faculty of Pharmacy, USC, Santiago de Compostela, 15782, Spain

<sup>5</sup> CINVESTAV-LANGEBIO, Irapuato, Guanajuato, 36500, México

<sup>6</sup> Caribbean vitroplants, Santo Domingo, 1464, Dominican Republic

<sup>7</sup> Vascular Biology Institute, School of Medicine, University of Miami, Florida, 33136, USA.

**Key words:** Spectral graph theory / Hidden Markov Model / Ribonucleases / Pac1 / Protein 2D representations

**List of abbreviations:** HP – Hydrophobicity and polarity; RNases – Ribonucleases; QSAR – Quantitative Structure-Activity Relationships; dsRNase – Double-strand-specific ribonuclease; snRNAs – small nucleolar RNA; LDA – Linear Discriminant Analysis; ORF – Open reading frame; MM – Markov Model; HMM – Hidden Markov Model; ROC curve – Receiver Operating Characteristic curve.

### Abstract

The study of type III RNases constitutes an important area in molecular biology. It is known that the *pac1*<sup>+</sup> gene encodes a particular RNase III that shares low amino acid similarity with other genes despite having a double-stranded ribonuclease activity. Bioinformatics methods based on sequence alignment may fail when there is a low amino acid identity percentage between query sequence and others with similar functions (remote homologues) or a similar sequence is not recorded in the database. Quantitative Structure-Activity Relationships (QSAR) applied to protein sequences may allow an alignment-independent prediction of protein function. These sequences QSAR like methods often use 1D sequence numerical parameters as the input to seek sequence-function relationships. However, previous 2D representation of sequences may uncover useful higher-order information. In the work described here we calculated for the first time the Spectral Moments of a Markov Matrix (MMM) associated with a 2D-HP-map of a protein sequence. We used MMMs values to characterize numerically 81 sequences of type III RNases and 133 proteins of a control group. We subsequently developed one MMM-QSAR and one classic Hidden Markov Model (HMM) based on the same data. The MMM-QSAR showed a discrimination power of RNases from other proteins of 97.35% without using alignment, which is a result as good as for the known HMM techniques. We also report for the first time the isolation of a new Pac1 protein (**DO647826**) from *Schizosaccharomyces pombe*, strain 428-4-1. The MMM-QSAR model predicts the new RNase III with the same accuracy as other classical alignment methods. Experimental assay of this protein confirms the predicted activity. The present results suggest that MMM-QSAR models may be used for protein function annotation avoiding sequence alignment with the same accuracy of classic HMM models.

**Corresponding author:** Agüero-Chapin, G.; Centro of Chemical Biactives and Faculty of Chemistry and Pharmacy, Central University of Las Villas, Santa Clara, 54830, Cuba, E-mail: [chapin@uclv.edu.cu](mailto:chapin@uclv.edu.cu)

## 1. Introduction

RNase III is a double-strand-specific ribonuclease (dsRNase) that usually makes staggered cuts in both strands of a double helical RNA, although in some cases it cleaves once in a single-stranded bulge in the helix<sup>1, 2</sup>. The primary biological function of this system is the specific processing of rRNA and mRNA precursors<sup>3-5</sup> but it has also been implicated in other diverse phenomena such as mRNA turnover<sup>6</sup>, conjugative DNA transfer<sup>7</sup>, antisense RNA-mediated regulation and other<sup>8, 9</sup>. For instance, Dicer and Drosha are type III RNases responsible for the generation of short interfering RNAs (siRNAs) from long double-stranded RNAs during RNA interference (RNAi). Also, the cellular processing of shRNAs shares common features with the biogenesis of naturally occurring miRNA such as cleavage by nuclear type III RNase Drosha, export from the nucleus, and processing by a cytoplasmic type III RNase Dicer, and incorporation into the RNA-induced silencing complex (RISC). Each step has a crucial influence on the efficiency of RNAi.<sup>10-13</sup> It involves both RNase proteins in several important biological processes as for instance the function of Dicer on the vascular system regulating embryonic angiogenesis probably by processing miRNAs, which regulate the expression levels of some critical angiogenic regulators in the cell.<sup>14</sup> Recently, RNAi has moved from a purely experimental technique to the stage of potential clinical applications such as possible use the treatment of spinocerebellar ataxia or amyotrophic lateral sclerosis<sup>15</sup>. Many other dsRNases have been characterized from a variety of prokaryotic and eukaryotic sources and RNase III from *Escherichia coli* is an archetype of this class of enzymes<sup>6, 16, 17</sup>. The RNase III family consists of a growing number of enzymes that includes at least 33 bacterial and 22 eukaryotic enzymes<sup>18</sup>. There have been numerous reports of dsRNase activities in eukaryotic cells, some of which exhibited properties consistent with a role in pre-rRNA processing<sup>19-21</sup>.

One of the best candidates for eukaryotic RNase III homologues is the Pac1 RNase from *Schizosaccharomyces pombe*<sup>22-24</sup>. The Pac1 product is derived from *Schizosaccharomyces pombe pac1*<sup>+</sup> gene expression, which is also involved in the regulation of sexual development<sup>25</sup>, possibly through a mechanism that involves the processing of certain small nucleolar RNAs (snRNAs)<sup>26</sup>. Pac1 works in eukaryotes as dsRNase and shares a functional similarity to RNase III from *E. coli*. This fact was proved either by measuring the ability of Pac1 to degrade double-stranded RNA *in vitro* or by expressing *pac1*<sup>+</sup> in *E. coli*, where it produced an activity that converted dsRNA into acid-soluble products<sup>23</sup>. Despite these observations the Pac1 gene product shows low homology with other RNase III enzymes, particularly with those ones belonging to bacteria. The homology between the different RNase III enzymes varies in the range 20 to 84% depending on their evolutionary distance, suggesting a low level of primary structure conservation<sup>27</sup>. It has been reported that antibodies prepared against Pac1 RNase have failed to react with RNase III<sup>23</sup>. The Pac1 gene product from *Schizosaccharomyces pombe* belongs to subclass II of the RNase III family, which is characterized by the presence of an N-terminal extension and includes fungal RNase III<sup>27, 28</sup>. This contains 363 amino acids (aa) and only its C-terminal 230 residues share 25% amino acid identity with the *Escherichia coli* ribonuclease III<sup>23</sup>.

Methods based on sequence alignment have revealed a low amino acidic identity (20–40 %) for the *pac1*<sup>+</sup> gene product with other typical RNases III, either isolated from bacteria or even from species that are genetically close<sup>27, 29</sup>. However, experimental observations show Pac1 protein to be a dsRNase enzyme. This relatively low degree of conservation probably reflects the species-specificity of RNase III, which prevents genetic complementation between members of the RNase III family<sup>30</sup>.

All of the facts discussed above hinder the prediction of the Pac1 gene product as an RNase III-like enzyme using computational methods based on sequence alignment. In fact, Bioinformatics methods based on sequence alignment may fail in general for cases of low sequence homology between the query and the template sequences deposited in the data base. The lack of function annotation (defined biological function) for the sequences deposited in databases and used as templates for function prediction constitutes another weakness of alignment approaches<sup>31, 32</sup>. Recently, a group of researchers published in PROTEOMICS (2006) a review<sup>33</sup> on the growing importance of machine learning methods for predicting protein functional

class independently of sequence similarity. In this review the authors make reference to various papers on the topic, including their own work<sup>34-45</sup>. These methods often use as the input 1D sequence numerical parameters specifically defined to seek sequence-function relationships. For instance, the so-called pseudo amino acid composition approach<sup>46, 47</sup> based on 1D sequence coupling numbers has been widely used to predict subcellular localization, enzyme family class, structural class, as well as other attributes of proteins based on their sequence similarity<sup>45, 48-74</sup>. Alternatively, some authors generalized molecular indices that are classically used for small molecules<sup>75, 76</sup> to describe protein sequences, such as the generalization of Broto–Moreau indices by Caballero and Fernández *et al.*<sup>77</sup>. On the other hand, many authors have introduced 2D or higher dimension representations of sequences prior to the calculation of numerical parameters. This constitutes an important step in order to uncover useful higher-order information not encoded by 1D sequence parameters<sup>78-98</sup>. In addition, 2D graphs have been used for proteins and DNA sequences by other researchers. For example, Zupan and Randić used spectral-like and zigzag representations. These authors suggested an algorithm for encoding long strings of building blocks (like four DNA bases, twenty natural amino acids, or all 64 possible base triplets) using “zigzag” or “spectrum-like” representations<sup>99</sup>. Hydrophobic cluster analysis (HCA) constitutes another well known technique for the 2D representation of protein sequences<sup>100</sup>. Randić *et al.* ultimately approached protein representations by using 2D schemes based on nucleotide triplet codons or virtual genetic code<sup>101</sup> and we introduced Hydrophobicity-Polarity (HP) 2D Cartesian or lattice-like representations for proteins related to plant metabolism<sup>93</sup>.

In this work, we propose to use the Spectral Moments of a Markov Matrix (MMM) associated to a 2D-HP-graph to numerically characterize protein sequences and seek a QSAR model to predict type III RNAses without alignment. Firstly, we derived Hydrophobicity-Polarity (HP) 2D Cartesian or lattice-like representations (also called maps or graphs) for RNase III and control group protein sequences<sup>93</sup>. We then calculated the MMM values of order  $k$  (symbolized as  $^{SR}\pi_k$ ) to characterize the protein sequence. Spectral Moments for many kinds of graphs have been used before for quantitative structure-activity relationships (QSAR) studies on proteins<sup>102-112</sup>. We subsequently developed a classifier to connect protein sequence information (represented by the  $^{SR}\pi_k$  values) with the classification of sequences as RNase III or not. In general, different kinds of classifiers have been used to derive protein sequence QSAR models<sup>113, 114</sup>. We selected a Linear Discriminant Analysis (LDA), which is a simple but powerful technique<sup>115-121</sup>. The use of this MMM-QSAR model enabled us to predict a novel recombinant Pac1 (rPac1) protein as an RNase III-like enzyme from a new isolate of *Schizosaccharomyces pombe*. Prediction was also supported by profile Hidden Markov Model (HMM) analysis, submission to BLASTp and InterPro<sup>122</sup> servers and demonstrated by experimental evidence.

## 2. Materials and methods

### 2.1 Computational methods.

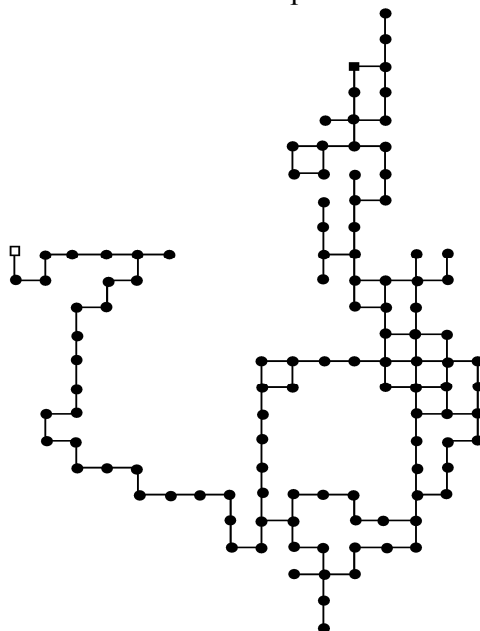
A Markov Model (MM), also called MARCH-INSIDE, was used to codify information about 81 RNase III protein sequences belonging to prokaryote and eukaryote species downloaded from the GenBank database. Briefly, our methodology considers as states of the Markov Chain (MC) any atom, nucleotide or amino acid (aa) depending on the kind of molecule to be described<sup>123, 124</sup>. Therefore, MM deals with the calculation of the probabilities ( $^k p_{ij}$ ) with which the charge distribution of aa moves from any aa in the vicinity  $i$  at time  $t_0$  to another aa  $j$  along the protein backbone in discrete time periods until a stationary state is achieved<sup>125, 126</sup>.

Each RNase III sequence was labelled by its accession number; see Table I in the supplementary material (SM). The control group consists of 133 proteins, which were selected from 2184 high-resolution proteins in a structurally non-redundant subset of the Protein Data Bank (PDB); most of the data were published by other authors to distinguish enzymes and non-enzymes without alignment<sup>127</sup> (see Table II in the SM). Many researchers have demonstrated the possibility of predicting protein function from sequences<sup>128</sup> and we used 2D-HP graphs to encode information about RNase III amino acid sequences<sup>93</sup>. We then calculated for the

first time the  $^{HP}\pi_k$  values for these graphs. As can be seen from the discussion above, we selected  $^{HP}\pi_k$  based on the utility of other non-stochastic Spectral Moments<sup>103-112</sup> as well as other MMMs and other stochastic parameters<sup>102, 129-131</sup>.

It is important to point out that this 2D graphical representation for proteins is similar to those previously reported for DNA<sup>89, 91, 132</sup> but the 20 different amino acids are regrouped into HP classes instead of using 4 types of bases. These four groups characterize the HP physicochemical nature of the amino acids as polar, non-polar, acidic or basic<sup>133</sup>. The 2D-HP graph for the deduced amino acid sequence of rPac1 protein, obtained from *Schizosaccharomyces pombe* strain 428-4-1 (uploaded by our group with accession number **DQ647826**), is shown in **Figure 1**. It is worth noting that 363 amino acids are rearranged in a 2D space compacting protein representation. Each amino acid in the sequence is placed in a Cartesian 2D space starting with the first monomer at the (0, 0) coordinates. The coordinates of the successive amino acids are calculated as follows:

- a) Increase by +1 the abscissa axis coordinate for an acid amino acid (rightwards-step) or:
- b) Decrease by -1 the abscissa axis coordinate for a basic amino acid (leftwards-step) or:
- c) Increase by +1 the ordinate axis coordinate for a polar amino acid (upwards-step) or:
- d) Decrease by -1 the ordinate axis coordinate for a non-polar amino acid (downwards-step).



**Figure 1.** 2D Cartesian representation for amino acid sequence of rPac1 protein from *Schizosaccharomyces pombe* strain 428-4-1; GenBank Accession number **DQ647826**. Note that a node may contain more than one amino acid, which ensures graph compactness.

## 2.2 2D-HP graph MMMs used as sequence numerical descriptors.

After the representation of the sequences we assigned to each graph a stochastic matrix  $^1\Pi$ . Note that the number of nodes (n) in the graph is equal to the number of rows and columns in  $^1\Pi$  but may be equal or even smaller than the number of amino acids or DNA bases in the sequence. The elements of  $^1\Pi$  are the probabilities  $^1p_{ij}$  of reaching a node  $n_i$  with charge  $Q_i$  moving through a walk of length  $k = 1$  from another node  $n_j$  with charge  $Q_j$ <sup>134</sup>:

$$p_{ij} = \frac{Q_j}{\sum_{m=l}^n \alpha_{il} \cdot Q_l} \quad (1)$$

Where  $\alpha_{ij}$  equals 1 if the nodes  $n_i$  and  $n_j$  are adjacent in the graph and equal to 0 otherwise.  $Q_j$  is equal to the sum of the electrostatic charges of all amino acids placed at this node. It then becomes straightforward to

carry out the calculation of the spectral moments of  ${}^1\Pi$  in order to numerically characterize the protein sequence:

$$MMM_k = {}^{SR}\pi_k = \sum_{i=j}^n p_{ij} = Tr\left[({}^1\Pi)^k\right] \quad (2)$$

Where Tr is called the trace and indicates that we sum all the values in the main diagonal of the matrices  ${}^k\Pi = ({}^1\Pi)^k$ , which are the natural powers of  ${}^1\Pi$ . The present class of MMMs encodes in a stochastic manner the distribution of the amino acid properties (charge) through all of the nodes placed at different distances in the 2D-HP lattice. Expansion of expression (2) for  $k = 0$  gives the order zero  $MMM_0$  ( ${}^{HP}\pi_0$ ); for  $k = 1$  the short-range  $MMM_1$  ( ${}^{HP}\pi_1$ ), for  $k = 2$  the middle-range  $MMM_2$  ( ${}^{HP}\pi_2$ ), and for  $k = 3$  the long-range MMMs. This extension is illustrated for the linear graph  $n_1$ - $n_2$ - $n_3$ , which is characteristic of the sequence (Asp-Glu-Asp-Lys); please note that the central node contains both Glu and Asp:

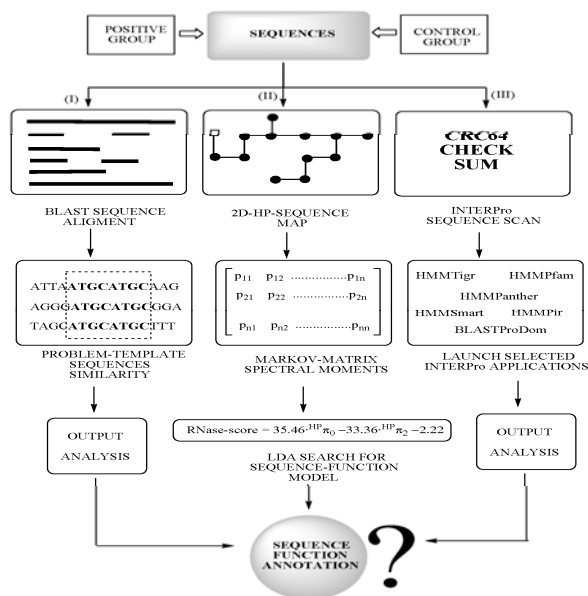
$${}^{HP}\pi_0 = Tr\left[({}^1\Pi)^0\right] = Tr\left[\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}\right] = 3 \quad (2a)$$

$${}^{HP}\pi_1 = Tr\left[({}^1\Pi)^1\right] = Tr\left[\begin{pmatrix} {}^1p_{11} & {}^1p_{12} & 0 \\ {}^1p_{21} & {}^1p_{22} & {}^1p_{23} \\ 0 & {}^1p_{32} & {}^1p_{33} \end{pmatrix}\right] = {}^1p_{11} + {}^1p_{22} + {}^1p_{33} \quad (2b)$$

$${}^{HP}\pi_2 = Tr\left[({}^1\Pi)^2\right] = Tr\left[\begin{pmatrix} {}^1p_{11} & {}^1p_{12} & 0 \\ {}^1p_{21} & {}^1p_{22} & {}^1p_{23} \\ 0 & {}^1p_{32} & {}^1p_{33} \end{pmatrix} \cdot \begin{pmatrix} {}^1p_{11} & {}^1p_{12} & 0 \\ {}^1p_{21} & {}^1p_{22} & {}^1p_{23} \\ 0 & {}^1p_{32} & {}^1p_{33} \end{pmatrix}\right] = {}^2p_{11} + {}^2p_{22} + {}^2p_{33} \quad (2c)$$

$${}^{HP}\pi_2 = Tr\left[({}^1\Pi)^2\right] = Tr\left[\begin{pmatrix} {}^1p_{11} & {}^1p_{12} & 0 \\ {}^1p_{21} & {}^1p_{22} & {}^1p_{23} \\ 0 & {}^1p_{32} & {}^1p_{33} \end{pmatrix} \cdot \begin{pmatrix} {}^1p_{11} & {}^1p_{12} & 0 \\ {}^1p_{21} & {}^1p_{22} & {}^1p_{23} \\ 0 & {}^1p_{32} & {}^1p_{33} \end{pmatrix} \cdot \begin{pmatrix} {}^1p_{11} & {}^1p_{12} & 0 \\ {}^1p_{21} & {}^1p_{22} & {}^1p_{23} \\ 0 & {}^1p_{32} & {}^1p_{33} \end{pmatrix}\right] = {}^3p_{11} + {}^3p_{22} + {}^3p_{33} \quad (2d)$$

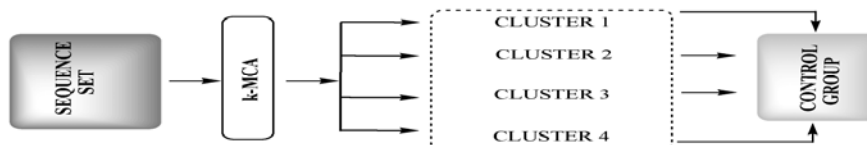
All calculations of  ${}^{HP}\pi_k$  values for protein sequences of both groups were carried out with our in-house software BIOMARKS *version 1.0* ®, including sequence representation <sup>135</sup>. We proceeded to upload a row data table with eleven  ${}^{HP}\pi_k$  values for each sequence ( $k = 0, 1, 2, \dots, 10$ ) and grouping variable RNaseIII-score = 1 (for RNAses) and -1 (for control group sequences) to statistical analysis software <sup>136</sup>. The overall methodology is represented schematically in order to improve the understanding of our approach (see **Figure 2**).



**Figure 2.** Schematic representation of the steps given in this work.

### 2.3 Statistical analysis. K-Means cluster analysis.

The negative group was selected from 2184 proteins with diverse functions (enzymes and non-enzymes) recorded in the PDB, as mentioned before. Our negative subset was designed according to K-Means cluster analysis (k-MCA) <sup>137</sup>. The method consists of carrying out a partition of the starting group made up by a non-RNase III series of proteins into several statistically representative clusters of sequences. Thus, one may select the members to conform to the negative subset from all of these clusters. This procedure ensures that the main protein classes (as determined by the clusters derived from k-MCA) will be represented in the model control group, thus allowing the representation of the entire ‘experimental universe’. The spectral moment series was explored as clustering variables in order to carry out k-MCA. The procedure described above is represented graphically in **Figure 3**, where a cluster analysis was carried out to select a representative sample for the control group.



**Figure 3.** k-MCA procedure for control group design.

## 2.4 Linear Discriminant Analysis.

LDA forward stepwise analysis was carried out for variable selection to build up the model<sup>115-121</sup>. All of the variables included in the model were standardized in order to bring them onto the same scale. Subsequently, a standardized linear discriminant equation that allows comparison of their coefficients was obtained<sup>138</sup>. The square of Mahalanobis's distance ( $D^2$ ) and Wilk's ( $\lambda$ ) statistic ( $\lambda = 0$  perfect discrimination, being  $0 < \lambda < 1$ ) were examined in order to assess the discriminatory power of the model. Pac1 protein was submitted to BLASTp to show graphically the similarity of the sequence compared to other RNases III. Each sequence presented in this study was also submitted to the InterPro server<sup>122</sup> in order to compare our methodology with other classical sources of predictive functional annotation. InterPro consists of a database of protein families, domains and functional sites in which identifiable features found in known proteins can be applied to unknown protein sequences.

## 3. Experimental Section

**3.1 Strains and culture media.** The *Schizosaccharomyces pombe* strain 428-4-1 was routinely grown in Yeast Extract (YEB) medium at 30 °C during 12 h. Bacterial strain *Escherichia coli* DH5 $\alpha$  was grown in Luria Broth (LB). Transformed bacteria were recovered in the same LB medium but supplemented with carbenicillin at 100  $\mu$ g/mL. Media were also supplemented with bacteriological agar when required.

**3.2 Total DNA extraction.** A colony from *Schizosaccharomyces pombe* strain 428-4-1 was inoculated in 5 mL of YEB medium and grown at 30 °C during 12 hours until OD<sub>600</sub> = 0.5. From this culture, 250  $\mu$ L was transferred to 50 mL of the same medium and grown overnight at the same temperature. When the OD<sub>600</sub> = 0.8, cells were collected by centrifugation and broken using small glass pearls. A cellular pellet was re-suspended in 500  $\mu$ L of sterile water at 50 °C and the extract was separated from cellular debris by centrifugation. Total DNA was purified using a total DNA extraction kit (Qiagen GmbH, Germany). Total DNA solution was measured at 260 nm in a GENESYS 10 spectrophotometer, reaching a concentration of 3.8  $\mu$ g/ $\mu$ L. The solution was also run on agarose gel (0.8%) and high integrity was seen.

**3.3 Primer design.** Forward (PAC5') 5'-ccc**ATGGGACGGTTTAAGAGGCATC**-3' and reverse (PAC3') 5'-gtggg**gttaac**cgggcaaac**TTAG**-3' primers were designed based on the previously reported *pac1*<sup>+</sup> coding sequence from *Schizosaccharomyces pombe* mutant *snm1-1*. The primer sequences show the restriction sites NcoI and KpnI introduced at the 3' and 5' ends, i.e. the first ATG and the stop TTA codon. The coding regions are shown in capital letters<sup>139</sup>.

**3.4 PCR amplifications.** Amplification of the *pac1*<sup>+</sup> gene from *Schizosaccharomyces pombe* was performed by standard PCR from its total DNA. The reaction mixture containing 10 ng of template, 1mM of each dNTP, 1.5 mM MgCl<sub>2</sub>, 2  $\mu$ M of each PAC5' and PAC3' primers, 1x buffer Taq Pol (Gibco BRL) and 2.5 U Taq Pol (Gibco) was completed to a total volume of 50  $\mu$ L. The PCR was carried out using a thermocycler (Perkin-Elmer 2400) programmed as follows: 5 minutes initial template denaturation at 94 °C, cycle steps: 1 minutes template denaturation at 94 °C, 2 minutes primer annealing at 45 °C, 2 minutes primer extension at 72 °C for 30 cycles; plus a final extension step at 72 °C for 5 minutes<sup>29, 30, 139</sup>. PCR reaction showed a band coinciding with the size of the reported *pac1*<sup>+</sup> ORF<sup>139</sup>.

**3.5 Plasmid construction and sequencing.** The PCR amplification product was purified using a GEL Band Purification kit (*AmershamPharmaciaBiotech*) and ligated to pMOS-Blue T-vector (*AmershamPharmaciaBiotech*). The ligation was transformed into electrocompetent *E. coli* DH5 $\alpha$  by electroporation in 0.2 mm cuvettes using a Gene Pulser Machine (BioRad) (12.5 kV, 25  $\mu$ F, 1000  $\omega$ ). The transformation was plated onto LB medium supplemented with 40  $\mu$ L of 20  $\mu$ g/mL X-gal solution and 4  $\mu$ L of isopropylthio- $\beta$ -D-galactoside from a 200  $\mu$ g/mL IPTG solution per plate and allowed to grow overnight at 37 °C. White colonies – presumably carrying the recombinant *pac1* gene inserted in pMOS-Blue T-vector,

named pRSPac1 – were selected and plasmid DNA extracted for analysis of the cloned fragment by restriction enzymes. Sequencing of the cloned fragment was performed using an ABI 3700 sequencer (Applied Biosystems)<sup>140</sup> and this showed a product of 1.111 Kb.

**3.6 Purification of recombinant Pac1.** A single colony of *E. coli* DH5 $\alpha$  with pRSPac1 was grown overnight at 30 °C in 5 mL of LB medium supplemented with carbencillin at 100  $\mu$ g/mL. 250  $\mu$ L of culture was then inoculated to 250 mL of the same medium supplemented with carbenicillin (100 $\mu$ g/mL) and grown under the same culture conditions until OD<sub>600</sub> = 0.8; at this point 50  $\mu$ L of 200  $\mu$ g/mL IPTG solution was added to the culture. Three hours after induction, cells were harvested by centrifugation and washed with 15 mL of 50 mM tris-HCl (pH 8), 100 mM NaCl and 1 mM EDTA. Cells were collected by centrifugation and stored at –70 °C overnight. Around 3 g of frozen cells were resuspended in 15 mL of lysis buffer (1% NP40, 0.5% sodium deoxycholate, 0.1 M NaCl, 30 mM Tris-HCl (pH 8), 1mM EDTA); 5 mM MgCl<sub>2</sub> and DNase I (10  $\mu$ g/mL) were added. The cell suspension was incubated on ice for 10 minutes. Inclusion bodies were collected by washing four times with lysis buffer and twice with 50 mM Tris-HCl 5 mM (pH 8), 1 mM DTT. Finally, the sample was dissolved in 5 mL of loading buffer and boiled on a water bath for 10 minutes. The total volume of extract was divided into five preparative PAGE electrophoresis samples containing 1 mL of protein extract, which were run in 12% gel. The component corresponding to 45.5 kDa recombinant Pac1 protein was visualized by staining with an aqueous solution of 0.05% Coomassie brilliant blue R250. In each case the recombinant protein was excised from polyacrylamide gel, recovered by electroelution, combined and concentrated using with a Centricon-10 (Amicon) to 0.5 mL and diluted to 1.5 mL with storage buffer to a final composition of 500 mM NaCl, 20 mM sodium phosphate (ph 7.4), 67 mM imidazole, 1 mM DTT, 1 mM EDTA and 30% glycerol. The recPac1 preparation was stored at –20 °C<sup>29, 30, 139</sup>.

**3.7 Synthesis and preparation of complementary RNA strands.** The enzymatic assay of recombinant Pac1 was carried out according to the optimized conditions described by Rotondo and Frendewey<sup>29</sup>. In a previous experiment (data not shown) we amplified by PCR a fragment corresponding to the fourth intron of *Schizosaccharomyces pombe*  $\beta$ -tubuline from its total DNA and inserted the amplified fragment into pBluescript II KS (–) for further *in vitro* transcription purposes. The integrity of the amplified sequence and transcriptional fusion was tested by sequencing. We reproduced exactly the described assay to compare the activity of our recombinant enzyme with the results from other reports. This construction was used as a template for the PCR of fragments corresponding to transcriptional-fusion suitable for the synthesis of both complementary strands of dsRNA substrate for an *in vitro* transcription reaction. For this purpose the following primers were synthesized:

a) 5'- gctcgaattaaccctcactaag↓ggaacGTAGGTTTTTTTGCTTTC-3' (T3 promoter in lower case, 5' end of the *Schizosaccharomyces pombe*  $\beta$ -tubuline fourth intron in upper case).

b) 5'-ggtacctaatacactactatag↓ggagaCTACAGTCGTCAGTAC-3' (T7 promoter in lower case, complement of the 3' end of the *Schizosaccharomyces pombe*  $\beta$ -tubuline fourth intron in upper case).

The arrows indicate the transcription initiation site. The PCR products were purified and 50 ng of each was used to synthesize both complementary strands of the dsRNA Pac1 substrate. The transcription reactions were prepared in a final volume of 20  $\mu$ L containing 40 mM Tris-HCl (pH 7.9), 6 mM MgCl<sub>2</sub>, 2 mM spermedine, 10 mM DTT, 0.5 mM of each ribonucleoside (*AmershamPharmaciaBiotech*), 50  $\mu$ Ci [ $\alpha$ <sup>32</sup> P] UTP (800 Ci/mmol), 20 U RNAsin (*Promega*) and 20 U T3 or T7 RNA Polymerase (*Amersham Pharmacia Biotech*). In the case of the transcription reaction driven by the T3 promoter, the addition of 50 mM NaCl to the reaction mixture was required. In all cases the reactions were prepared on ice were then incubated at 37 °C during 10 minutes. The resulting transcripts were treated with DNase I (*Promega*), phenol extraction and precipitation with 2.5 V/V of absolute ethanol was carried out and the samples were stored overnight at –70 °C. The complementary RNA strands were collected by centrifugation at 16 000 g during 10 minutes at 4 °C. Finally, the pellets were washed with 70% ethanol, dried and re-suspended in diethyl pyrocarbonate treated with distilled water and stored at –70 °C.

**3.8 Preparation of dsRNA substrate for Pac1 enzymatic assay.** Equimolar quantities of both complementary strands were mixed in diethyl pyrocarbonate and treated distilled water to give a final volume of 50  $\mu$ L. The mixture was heated during 10 minutes at 100 °C in a water bath. The whole bath was



then firmly closed and placed into thermal box overnight to allow annealing of both complementary strands into the dsRNA substrate. The unpaired ends and RNA strands were removed by RNase A (Promega) treatment. The dsRNA substrate was purified (PAGE-TBE 15% gel) and stored in diethyl pyrocarbonate (DEPC) treated distilled water at  $-70\text{ }^{\circ}\text{C}$ . The substrate for the Pac1 assay consisted of 101 bp dsRNA, identical to the substrate used by Rotondo and Frendewey<sup>29</sup>.

**3.9 Enzymatic assay of recombinant Pac1.** The Pac1 assay was carried out using the following conditions: 30 mM Tris-HCl (pH7.6), 1 mM DTT, 5 mM of  $\text{MgCl}_2$ , 10 nM of dsRNA substrate and different quantities (0, 1, 10, 100 nM) of purified recombinant Pac1 enzyme. Enzymatic reactions were completed on ice and started by the addition of 0.1V of 50 mM  $\text{MgCl}_2$ , incubated at  $30\text{ }^{\circ}\text{C}$  for 10 minutes and stopped by the addition of 500  $\mu\text{L}$  of 5% ice-cooled TCA followed by 15 minutes on ice. The aliquots were centrifuged at 16 000 g during 5 minutes in a Spin-X filter unit (Costar). The soluble fractions (filtrate) were quantified by liquid scintillation counting. The counting data represent the amount of acid-precipitable polynucleotide phosphorus (dsRNA) substrate transformed into acid soluble cleavage products by Pac1 enzyme. The procedure was repeated three times with three repetitions per experiment<sup>29, 30, 139</sup>.

## 4. Results and discussion

### 4.1 MMM-QSAR model to predict type III RNAses without alignment.

Many different parameters can be used to encode protein sequence information and further assign or predict the function or physical properties of proteins and their mutants<sup>141, 142</sup>. The present approach involves the calculation of different sequence parameters based on MMs, which can be applied to different kinds of molecular graphs<sup>131</sup> including DNA, RNA and proteins<sup>93, 143</sup>. MMs have been applied successfully to Genomics and Proteomics and represent an important tool for analyzing biological sequence data. In particular, MMs have been used for protein folding recognition<sup>144</sup> and the prediction of protein signal sequences<sup>145, 146</sup>. MMs have also been applied to predict alpha turns<sup>147</sup>, beta turns<sup>148</sup>, as well as other tight turns and their types<sup>149</sup>. Particularly, MMs have been further used to predict the specificity of GalNac-transferase<sup>150</sup> and cleavage sites in proteins by proteases<sup>151-154</sup>, greatly stimulating the development for drug design against AIDS and SARS<sup>155-163</sup>. In this work we calculated MMMs ( ${}^{HP}\pi_k$ ) of the stochastic matrix that describe the distribution of the amino acids of the protein sequence in the 2D-HP graph. This calculation was carried out for two groups of protein sequences, one made up of RNase III-like enzymes and the other formed by heterogeneous proteins. This last group contains 133 members and these were selected as follows:

Original data were submitted to k-Means cluster analysis as described previously. The k-MCA divided the data into four clusters containing 439, 684, 592 and 469 members, respectively. Selection was based on the distance from each member with respect to the cluster centre (Euclidean distance). We selected the closer cases to the centre in order to ensure the inclusion of representative members of each cluster in the control group. Depending on the cluster size, a proportional number of proteins were set; 27 cases were taken from the first cluster, 42 from the second, 36 from the third and 28 from the fourth to give a total of 133 members in the control group. We always bore in mind the principle of discriminant analysis in terms of balancing the size of the control group with respect to the RNase III group.

A simple MMM-QSAR was then developed to classify a novel sequence as RNase III or not. The best equation found for this purpose was:

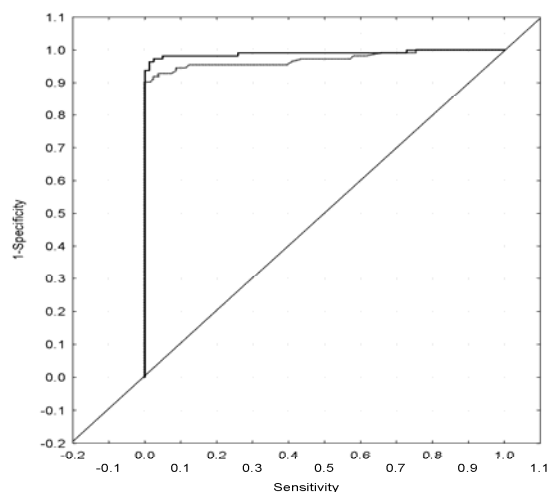
$$\text{RNaseIII-score} = 35.46 \times {}^{HP}\pi_0 - 33.36 \times {}^{HP}\pi_2 - 2.22 \quad (3)$$

The statistical parameters for the above equation were Wilk's statistic ( $\lambda = 0.18$ ), Mahalanobis's distance ( $D^2 = 16.36$ ) and error level (p-level  $< 0.001$ )<sup>164</sup>. This discriminant function misclassified only four cases out of 214 proteins used in both the training and validation series, reaching a high level of accuracy of 98.13%. More specifically, the model classified correctly 77/81 (95.06%) of RNase III-like enzymes and 100% of the control group. The respective classification matrices for training and cross-validation are depicted in **Table 1**.

**Table 1.** Classification results derived from the model for training and validation series

MMM Training				MMM Validation			
Total%	97.35	RNases	Control	RNases	Control	100	Total%
RNases	93.44	<b>57</b>	4	<b>20</b>	0	100	RNases
Control	100	0	<b>90</b>	0	<b>43</b>	100	Control
MMM All sequences				HMM classic			
Total%	98.1	RNases	Control	RNases	Control	97.50	Total%
RNases	95.1	77	4	<b>80</b>	1	98.75	RNases
Control	100	0	133	5	<b>128</b>	96.24	Control

A validation procedure was subsequently performed in order to assess the model predictability. This validation was carried out with an external series of 20 RNase III-like proteins and a further 43 diverse proteins (see **Table 1**). The present model showed an average predictability of 100% for each group, which is remarkable in comparison to results obtained by other researchers on using the LDA method in QSAR studies<sup>165-168</sup>. These results are consistent with those obtained in our previous report, in which we used 2D coupling numbers as sequence descriptors for function annotation of plant metabolism enzymes<sup>93</sup>. In addition, we carried out a classification analysis with all of the proteins included. These results provide further evidence of the robustness of the results obtained. The Receiver Operating Characteristic (ROC) curve was also constructed for the training and validation series. Notably, the curve presented a pronounced curvature (convexity) with respect to the  $y = x$  line for both series (see **Figure 4**). This result confirms that the present model is a significant classifier, having areas of 0.99 (training) and 0.97 (validation) – i.e. markedly higher than 0.5, which is the value for a random classifier<sup>169</sup>.



**Figure 4.** Receiver Operating Characteristic curve (ROC-curve) for training (dark line), validation (dot line) and random classifier (light line) with areas under curve of 0.99, 0.97, and 0.5, respectively.

## 4.2 Isolation, prediction and assay of a novel Pac1 from *Schizosaccharomyces pombe* strain 428-4-1

### 4.2.1 Isolation.

In this work we isolated, cloned and expressed a new Pac1 DNA sequence from *Schizosaccharomyces pombe* strain 428-4-1, its nucleotide and amino acid sequence was recorded on the GenBank database with accession number **DQ647826**. The theoretical prediction of its translated ORF as an RNase III-like enzyme was performed by the present alignment-independent approach instead of traditional alignment methods. The theoretical prediction of rPac1 as a double-stranded RNase was confirmed experimentally by *in vitro* assays.

### 4.2.2 Prediction.

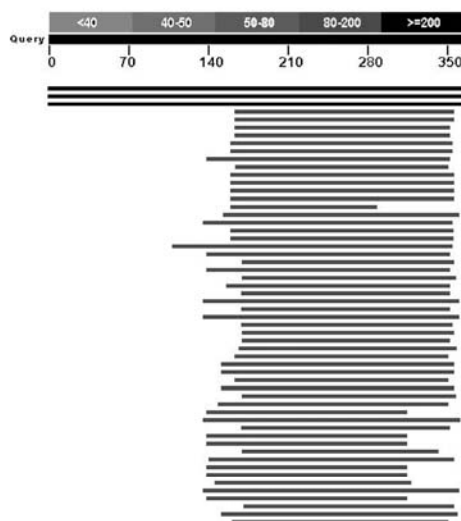
Our Pac1 protein sequence was analyzed using the MMM-QSAR methodology with the aim of recognizing the rPac1 gene product as a eukaryotic RNase III homologue. The sequence was represented in a Cartesian 2D system and calculated including the whole data set. This particular case was included in the validation subset in order to make a prediction. The MMM-QSAR model even very simple (two variables) allowed the correctly classification of the rPac1 product as an RNase III-like enzyme with the maximum probability ( $p = 1$ ). In order to make a graphical comparison between our methodology and alignment methods like BLASTp<sup>170-173</sup>, several representative RNase III protein sequences from prokaryotes and eukaryotes were selected together with rPac1 for representation in a 2D-mapping system (see **Figure 5**).



**Figure 5.** 2D-HP map superposition of RNases from prokaryotes (dark grey), eukaryotes (in light grey) and rPac1 **DQ647826** from *Schizosaccharomyces pombe* strain 428-4-1 (in black).

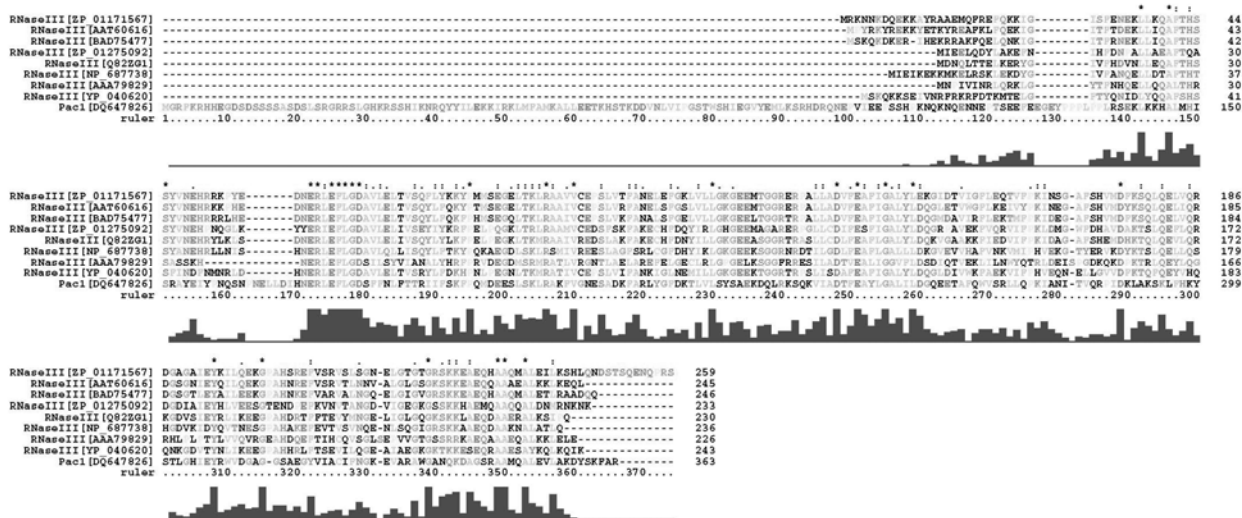
The 2D-HP map protein representation revealed a significant separation for the groups consisting of dsRNases from prokaryotes (in dark grey) and eukaryotes (in light grey). The rPac1 protein (in black) is placed between the two groups, acting as a sort of link between the RNase III families. This representation possibly supports evolutionary relationships between double-stranded RNase protein sequences. Since the Cartesian 2D protein representation is mainly based on amino acid composition, we can highlight a major region from rPac1 matching eukaryote sequences (in light grey) and another small region that lies within the prokaryote region (in dark grey). There is also a non-matching region specific for rPac1 in *Schizosaccharomyces pombe* that does not exist in other eukaryotes. However, matching regions in the graph made a significant contribution to calculation of the spectral moments, thus allowing successful recognition of rPac1 as RNase III.

A BLASTp analysis was carried out on the translated rPac1 DNA sequence (see **Figure 6**). This method recognized successfully our query sequence as a Pac1 ribonuclease, reaching up to 98% of amino acid identity with others already recorded from *Schizosaccharomyces pombe* strains. Although this analysis showed lower scores (close to 80%) in comparison to other typical dsRNases, the approach still enabled protein query recognition as RNase III. With the aim of comparing different methods, it is possible to set an equivalence for the score value (80%) from BLASTp with our predicted probability,  $p = 1$ , for rPac1 to act as an RNase III-like enzyme. BLASTp also revealed low amino acid identity (< 40%) toward the C-terminal portion despite this representing the highest conserved region in the four existing RNase III subclasses. On the other hand, as mentioned previously, each sequence included in the study was submitted to InterProt. All cases (100%) from the RNase III group matched significantly with RNase III domains (IPR000999), allowing the total recognition as dsRNases (see **Tables ISM**). In the case of the control group, six cases did not have InterProt identification and three of them did not have any hits reported (95, 50% of predictability) (see **Table IISM**).



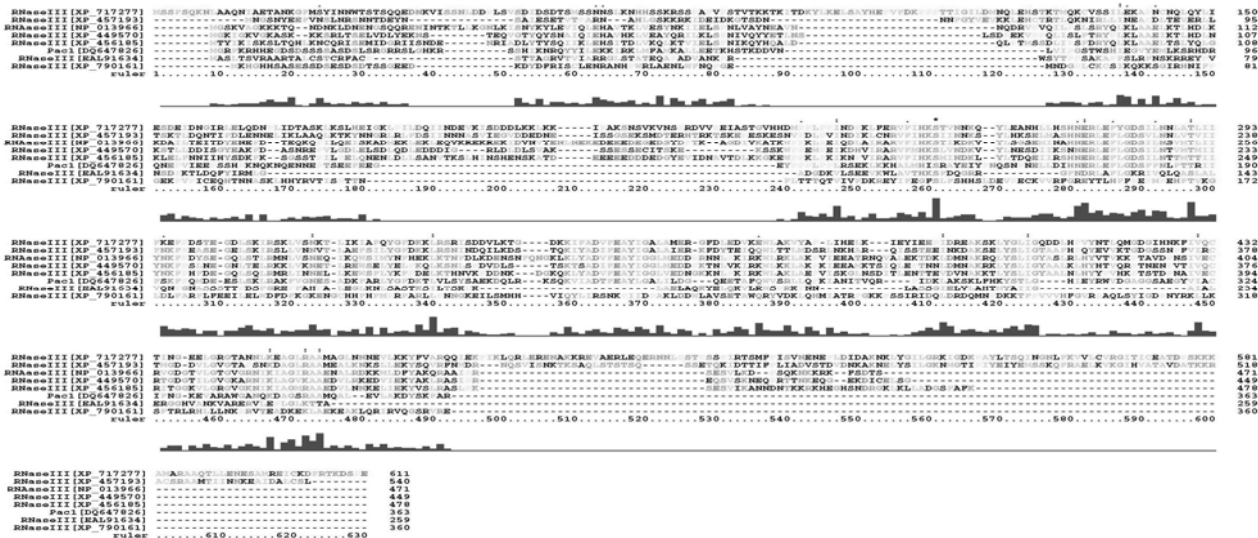
**Figure 6.** BLASTp analysis for rPac1 protein sequence [DQ647826](#). Note that the scale of scoring is progressive in darkness. Sequence names are not depicted.

We also performed an alignment between the previously selected sequences in the figure 5 and our rPac1 product using the Clustal W program, version 1.81 (see **Figure 7** and **Figure 8**). Alignment results coincide with those obtained in previous studies reported by other authors. The rPac1 showed low amino acid identity percentages in comparison to dsRNase sequences from other eukaryote organisms, even for those belonging to yeast-related species. Short and less frequent regions match along the protein sequences, especially toward the N-terminal region (see **Figure 8**). The comparison with prokaryote sequences showed a matching region toward the protein's C-terminal part, from the 170 up to 260 amino acid position. This region corresponds with the RNase III C-terminal domain (RIBOc), which is conserved in eukaryotic, bacterial and archeal RNase III and is associated with the catalytic activity. There is a significant N-terminal region in the Pac1 product that does not appear in the RNase III prokaryote family – a finding consistent with other reports (see **Figure 7**)<sup>29</sup>.



**Figure 7.** Clustal X sequence alignment involving RNase III like enzymes, each sequence is represented by its accession to GenBank Database Protein. Sequences used in the alignment were represented previously in Cartesian 2D system (**Fig. 5**). We use sequences from bacteria and rPac1 from *S. pombe*. [ZP\_01171567] *Bacillus* sp. NRRL B 14911, [AA760616] *Bacillus thuringiensis*, [BAD75477] *Geobacillus kaustophilus* HTA426, [ZP\_01275092] *Lactobacillus reuteri*, [Q822G1] *Enterococcus faecalis*, [NP\_687738] *Streptococcus agalactiae* 2603VR, [AAA79829] *E.coli*, [YP\_040620] *Staphylococcus aureus* MRSA252, [DQ647826] *S.pombe* strain 428-4-1





**Figure 8.** Clustal X sequence alignment involving RNase III like enzymes, each sequence is represented by its accession to GenBank Database Protein. Sequences used in the alignment were represented previously in Cartesian 2D system (**Figure 5**). We use sequences from some representative eukaryotes and rPac1 from *S. pombe*. [XP\_717277] *Candida albicans* SC5314, [XP\_457193] *Debaryomyces hansenii* CBS767, [NP\_013966] *Saccharomyces cerevisiae*, [XP\_449570] *Candida glabrata* CBS138, [XP\_456185] *Kluyveromyces lactis*, [DQ647826] *S.pombe* strain 428-4-1, [EAL91634] *Aspergillus fumigatus* Af293, [XP\_790161] *Strongylocentrotus purpuratus*.

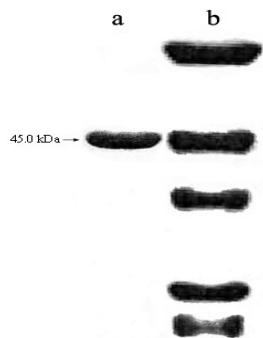
The results found in this study confirm that our model does not replace classical method for protein function annotation like BLAST or InterProt service, but becomes an interesting alternative tool – especially due to its alignment-independence and simplicity. It is also important to highlight that our methodology can be considered as a good classifier, despite its simplicity, as it gives rise to a linear equation with two variables at most. Consequently, it is a useful method to perform a quick virtual screening of a representative protein database since the protein query submission to classical sequence classifiers is generally performed on a one by one basis. Thus, once the whole database has been screened and proteins having the desired function are recognized, it would be advisable to assess results obtained using our approach by other methodologies. The search for approaches that complement or improve on classical alignment tools like BLAST with information from gene ontology, RNA secondary structure prediction, partial ordering or other sources constitutes a goal of major importance<sup>174-178</sup>.

In order to compare the MMM-QSAR approach reported here with other methodologies based on MM, training and negative (non-RNases sequences) sets were scored with a classic HMMs. Classification driven by an HMM built on the original training set resulted in an accuracy of 98.75% for the positives sequences (training set) and 96.24% for the negative sequences (see **Table 1**). Our query sequence **DQ647826** was also successfully predicted with the maximum score by the HMM.

**4.2.3 Experimental evidence for RNase III activity.** Recombinant Pac1 protein from *Schizosaccharomyces pombe* strain 428-4-1 was purified in order to measure its double-stranded RNase activity *in vitro*. The corresponding product size (45.5 kDa) coincided with the reported size for the native protein (see **Figure 9**). Double-stranded activity was measured *in vitro* by following the protocol described above. The unit definition for all RNase III types is the amount of enzyme able to solubilize 1 nmol of acid-precipitable radioactivity per hour.<sup>17</sup> Pac1 activity showed values comparable to other results ( $5 \times 10^5$  U/mg) obtained for a recombinant Pac1 product from *Schizosaccharomyces pombe* by Rotondo and



Frendewey<sup>29</sup>. Results derived from the enzymatic activity assay are shown in **Table 2** for each experiment; the mean value was  $6.96 \times 10^5$  U/mg.



**Figure 9.** Electrophoresis of rPac1 protein, 45 kDa rPac1 was purified and loaded in 12.5% PAGE-SDS and stained with Coomassie brilliant. (a) Band corresponding to rPac1 purified (b) Molecular weight marker; 66.2 kDa, 45.7 kDa, 31 kDa, 21.5 kDa, 14.4 kDa (Unstained SDS-PAGE Standards Broad Range, BioRad)

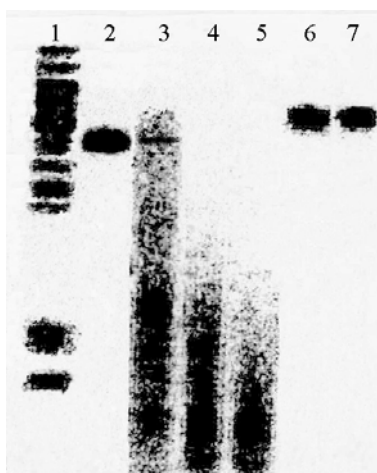
**Table 2.** Enzymatic assay of double-stranded RNase recombinant Pac1 **DQ647826** extracted from *Schizosaccharomyces pombe* strain 428-4-1

Conc. rPac 1	1nM	10 nM	100 nM
EUV <sup>a</sup>	$6.2 \times 10^5$	$7.4 \times 10^5$	$7.2 \times 10^5$
	$6.0 \times 10^5$	$6.8 \times 10^5$	$7.3 \times 10^5$
	$6.6 \times 10^5$	$6.9 \times 10^5$	$7.9 \times 10^5$
Mean	$6.4 \times 10^5$	$7.0 \times 10^5$	$7.5 \times 10^5$

<sup>a</sup> Enzymatic unit value for rPac 1 (U/mg)

The kinetic enzymatic reaction of rPac1 by monitoring dsRNA integrity (lanes 2–5) is illustrated in **Figure 10**. This particular results for the **DQ647826** sequence were not carry out to validate the MMM-QSAR model but to shown how to use it for predicting RNase III-like protein function annotation. We recall that the validation of the MMM-QSAR model was assessed with the external prediction series as recommended for any QSAR (see previous sections).<sup>179</sup>





**Figure 10.** Autoradiography of rPac1 enzymatic assay. To visualize the cleavage activity of dsRNA substrate generated by T3/T7 "in vitro" transcription, aliquots of enzymatic assay were taken at 2, 5 and 10 minutes, loaded in 12.5% PAGE/7M Urea followed by autoradiography. Lane 1 is pBR322 digested by Msp1, Lane 2 is the intact dsRNA substrate, Lane 3 to 5 are the results of rPac1 enzymatic activity at 2, 5 and 10 minutes of reaction at 30°C. Lane 6 and 7 are the T3 and T7 ssRNA obtained by "in vitro" transcription too which are not degraded by RNase activity of Pac1.

## 5. Conclusions

The work described here introduces a new approach to predict RNase type III function from protein sequences irrespective of sequence alignment. The methodology uses the MMMs associated with a 2D sequence representation as the input for an LDA classifier. This MMM-QSAR classifier successfully discriminates between RNase-like sequences and a control group. The Pac1 gene product was chosen as a representative example of a sequence with low amino acid identity compared to other enzymes with similar activity. The present methodology achieves high classification scores similar than bioinformatics tools based on sequence alignment (BLASTp) and comparable results to other predicting protein function annotation methods like InterProt and HMMs. The predictions made by the present model coincide with outcomes from experimental isolation, expression, and enzymatic activity measurement of a novel *pac1*<sup>+</sup> gene sequence **DQ647826** isolated from a new isolate *Schizosaccharomyces pombe* strain 428-4-1. The work opens up new possibilities for the use of the experience accumulated in small molecules QSAR in the field kind of alignment-independent sequence function annotation.

## Supplementary Material

Detailed information on the proteins used in the study is supplied in the online **Supplementary material** and this includes organism, accession number, protein definition, values of the stochastic spectral moments and scores (**Tables ISM and IISM**). This information is available free of charge via the Internet at <http://pubs.acs.org>.

## References

1. Robertson, H., *Escherichia coli* ribonuclease III cleavage sites. *Cell* **1982**, 30, 669–672.
2. Chelladurai, B., Li, H., Zhang, K. and Nicholson, A. W., Mutational analysis of a ribonuclease III processing signal *Biochemistry* **1993**, 32, 7549–7558.
3. Abeyrathne, P. D.; Nazar, R. N., Parallels in rRNA processing: conserved features in the processing of the internal transcribed spacer 1 in the pre-rRNA from *Schizosaccharomyces pombe*. *Biochemistry* **2005**, 44, (51), 16977–87.
4. Dunn, J. J.; Studier, F. W., T7 early RNAs and *Escherichia coli* ribosomal RNAs are cut from large precursor RNAs in vivo by ribonuclease III. *Proc Natl Acad Sci U S A* **1973**, 70, (12), 3296–3300.



5. Young, R. A.; Steitz, J. A., Complementary sequences 1700 nucleotides apart from a ribonuclease III cleavage site in *Escherichia coli* ribosomal precursor RNA. *Proc. Natl. Acad. Sci. USA* **1978**, *75*, 3593–3597.
6. Court, D., *RNA processing and degradation by RNase III*. In: *Control of mRNA stability*. Brawerman, G., Belasco, J. ed.; Academic Press: New York, 1993; p 70–116.
7. Koraimann, G., Schroller, C., Graus, H., Angerer, D., Teferle, K. and; Hogenauer, G., Expression of gene 19 of the conjugative plasmid R1 is controlled by RNase III. *Mol. Microbiol.* **1993**, *9*, 717–727.
8. Gerdes, K., Nielsen, A., Thorsted, P. and Wagner, E. G, Mechanism of killer gene activation. Antisense RNA-dependent RNase III cleavage ensures rapid turn-over of the stable hok, srnB and pndA effector messenger RNAs. *J. Mol. Biol.* **1992**, *226*, 637–649.
9. Blomberg, P., Wagner, E. G. H. and Nordström, K., Control of replication of plasmid R1: The duplex between the antisense RNA, CopA, and its target, CopT, is processed specifically in vivo and in vitro by RNase III. *EMBO J* **1990**, *9*, 2331–2340.
10. Gregory, R. I.; Chendrimada, T. P.; Shiekhattar, R., MicroRNA biogenesis: isolation and characterization of the microprocessor complex. *Methods Mol Biol* **2006**, *342*, 33-47.
11. Carmell, M. A.; Hannon, G. J., RNase III enzymes and the initiation of gene silencing. *Nat Struct Mol Biol* **2004**, *11*, (3), 214-8.
12. Meister, G.; Landthaler, M.; Peters, L.; Chen, P. Y.; Urlaub, H.; Luhrmann, R.; Tuschl, T., Identification of novel argonaute-associated proteins. *Curr Biol* **2005**, *15*, (23), 2149-55.
13. Han, J.; Lee, Y.; Yeom, K. H.; Kim, Y. K.; Jin, H.; Kim, V. N., The Drosha-DGCR8 complex in primary microRNA processing. *Genes Dev* **2004**, *18*, (24), 3016-27.
14. Yang, W. J.; Yang, D. D.; Na, S.; Sandusky, G. E.; Zhang, Q.; Zhao, G., Dicer is required for embryonic angiogenesis during mouse development. *J Biol Chem* **2005**, *280*, (10), 9330-5.
15. Pekarik, V., Design of shRNAs for RNAi-A lesson from pre-miRNA processing: possible clinical applications. *Brain Res Bull* **2005**, *68*, (1-2), 115-20.
16. Robertson, H. D., Webster, R. E. and Zinder, N. D, Purification and properties of ribonuclease III from *Escherichia coli* *J. Biol. Chem* **1968**, *243*, 82-91.
17. Dunn J, J., *Ribonulcease III*. In: *The Enzymes*. Boyer, P.D ed.; Academic Press: New York, 1982; Vol. 15 (Part B), p 485–499.
18. Nicholson, A. W., Structure, reactivity and biology of double-stranded RNA. *Prog. Nucleic Acid Res. Mol. Biol.* **1996**, *52*, 1–65.
19. Grummt, I., Hall, S. H. and Crouch, R. J., Localisation of an endonuclease specific for double-stranded RNA within the nucleolus and its implication in processing ribosomal transcripts. *Eur. J. Biochem* **1979**, *94*, 437–443.
20. Ohtsuki, K., Groner, Y. and Hurwitz, J, Isolation and purification of double-stranded ribonuclease from calf thymus. *J. Biol. Chem.* **1977**, *252*, 483–491.
21. Lipardi, C.; Baek, H. J.; Wei, Q.; Paterson, B. M., Analysis of short interfering RNA function in RNA interference by using *Drosophila* embryo extracts and schneider cells. *Methods Enzymol* **2005**, *392*, 351-71.
22. Xu, H.-P., Riggs, M., Rodgers, L. and Wigler, M., A gene from *S. Pombe* with homology to *E. Coli* RNase III blocks conjugation and sporulation when overexpressed in wild type cells. *Nucleic Acids Res* **1990**, *18*, 5304.
23. Iino, Y., Sugimoto, A. and Yamamoto, M., *S. pombe pac1<sup>+</sup>*, whose over-expression inhibits sexual development, encodes a ribonuclease III-like RNase *EMBO J* **1991**, *10*, 221–226.
24. Qian, Z.; Xuan, B.; Hong, J.; Hao, Z.; Wang, L.; Huang, W., Expression and purification of the carboxyl terminus domain of *Schizosaccharomyces pombe* dicer in *Escherichia coli*. *Protein Pept Lett* **2005**, *12*, (4), 311-4.
25. Xu, H. P.; Riggs, M.; Rodgers, L.; Wigler, M., A gene from *S. pombe* with homology to *E. coli* RNase III blocks conjugation and sporulation when overexpressed in wild type cells. *Nucleic Acids Res* **1990**, *18*, (17), 5304.

26. Potashkin J, F. D., A mutation in a single gene of *Schizosaccharomyces pombe* affects the expression of several snRNAs and causes defects in RNA processing. *EMBO J* **1990**, 9, 525–534.
27. Lamontagne, B.; Elela, S. A., Evaluation of the RNA determinants for bacterial and yeast RNase III binding and cleavage. *J Biol Chem* **2004**, 279, (3), 2231-41.
28. Saunders, L. R.; Barber, G. N., The dsRNA binding protein family: critical roles, diverse cellular functions. *The FASEB Journal* **2003**, 17, 961-983.
29. Rotondo, G.; Frenthewey, D., Purification and characterization of the Pac1 ribonuclease of *Schizosaccharomyces pombe*. *Nucleic Acids Res* **1996**, 24, (12), 2377-86.
30. Rotondo, G.; Huang, J. Y.; Frenthewey, D., Substrate structure requirements of the Pac1 ribonuclease from *Schizosaccharmyces pombe*. *RNA* **1997**, 3, (10), 1182-93.
31. Dobson, P. D.; Doig, A. J., Predicting enzyme class from protein structure without alignments. *J Mol Biol* **2005**, 345, (1), 187-99.
32. Dobson, P. D.; Cai, Y. D.; Stapley, B. J.; Doig, A. J., Prediction of protein function in the absence of significant sequence similarity. *Curr Med Chem* **2004**, 11, (16), 2135-42.
33. Han, L.; Cui, J.; Lin, H.; Ji, Z.; Cao, Z.; Li, Y.; Chen, Y., Recent progresses in the application of machine learning approach for predicting protein functional class independent of sequence similarity. *Proteomics* **2006**, 6, (14), 4023-37.
34. Lin, H. H.; Han, L. Y.; Zhang, H. L.; Zheng, C. J.; Xie, B.; Chen, Y. Z., Prediction of the functional class of lipid binding proteins from sequence-derived properties irrespective of sequence similarity. *J Lipid Res* **2006**, 47, (4), 824-31.
35. Cui, J.; Han, L. Y.; Cai, C. Z.; Zheng, C. J.; Ji, Z. L.; Chen, Y. Z., Prediction of functional class of novel bacterial proteins without the use of sequence similarity by a statistical learning method. *J Mol Microbiol Biotechnol* **2005**, 9, (2), 86-100.
36. Han, L. Y.; Cai, C. Z.; Ji, Z. L.; Cao, Z. W.; Cui, J.; Chen, Y. Z., Predicting functional family of novel enzymes irrespective of sequence similarity: a statistical learning approach. *Nucleic Acids Res* **2004**, 32, (21), 6437-44.
37. Han, L. Y.; Cai, C. Z.; Ji, Z. L.; Chen, Y. Z., Prediction of functional class of novel viral proteins by a statistical learning method irrespective of sequence similarity. *Virology* **2005**, 331, (1), 136-43.
38. Han, L. Y.; Cai, C. Z.; Lo, S. L.; Chung, M. C.; Chen, Y. Z., Prediction of RNA-binding proteins from primary sequence by a support vector machine approach. *Rna* **2004**, 10, (3), 355-68.
39. Chen, J.; Liu, H.; Yang, J.; Chou, K. C., Prediction of linear B-cell epitopes using amino acid pair antigenicity scale. *Amino Acids* **2007**, 33, (3), 423-8.
40. Guo, Y. Z.; Li, M.; Lu, M.; Wen, Z.; Wang, K.; Li, G.; Wu, J., Classifying G protein-coupled receptors and nuclear receptors based on protein power spectrum from fast Fourier transform. *Amino Acids* **2006**, 30, 397-402.
41. Chou, K. C.; Elrod, D. W., Prediction of enzyme family classes. *J Proteome Res* **2003**, 2, (2), 183-90.
42. Elrod, D. W.; Chou, K. C., A study on the correlation of G-protein-coupled receptor types with amino acid composition. *Protein Eng* **2002**, 15, (9), 713-5.
43. Sun, X. D.; Huang, R. B., Prediction of protein structural classes using support vector machines. *Amino Acids* **2006**, 30, 469-475.
44. Wen, Z.; Li, M.; Li, Y.; Guo, Y.; Wang, K., Delaunay triangulation with partial least squares projection to latent structures: a model for G-protein coupled receptors classification and fast structure recognition. *Amino Acids* **2006**, 32, 277-283.
45. Zhang, S. W.; Pan, Q.; Zhang, H. C.; Shao, Z. C.; Shi, J. Y., Prediction protein homo-oligomer types by pseudo amino acid composition: Approached with an improved feature extraction and naive Bayes feature fusion. *Amino Acids* **2006**, 30, 461-468.
46. Chou, K. C., Using amphiphilic pseudo amino acid composition to predict enzyme subfamily classes. *Bioinformatics* **2005**, 21, (1), 10-9.
47. Chou, K. C., Prediction of protein cellular attributes using pseudo-amino acid composition. *Proteins* **2001**, 43, (3), 246-55.

48. Chou, K. C.; Shen, H. B., Hum-PLoc: a novel ensemble classifier for predicting human protein subcellular localization. *Biochem Biophys Res Commun* **2006**, 347, (1), 150-7.
49. Xiao, X.; Shao, S.; Ding, Y.; Huang, Z.; Huang, Y.; Chou, K. C., Using complexity measure factor to predict protein subcellular location. *Amino Acids* **2005**, 28, (1), 57-61.
50. Xiao, X.; Shao, S.; Ding, Y.; Huang, Z.; Chen, X.; Chou, K. C., Using cellular automata to generate image representation for biological sequences. *Amino Acids* **2005**, 28, (1), 29-35.
51. Gao, Y.; Shao, S.; Xiao, X.; Ding, Y.; Huang, Y.; Huang, Z.; Chou, K. C., Using pseudo amino acid composition to predict protein subcellular location: Approached with Lyapunov index, Bessel function, and Chebyshev filter. *Amino Acids* **2005**, 28, (4), 373-6.
52. Wang, S. Q.; Yang, J.; Chou, K. C., Using stacked generalization to predict membrane protein types based on pseudo-amino acid composition. *J Theor Biol* **2006**, 242, (4), 941-6.
53. Shen, H. B.; Chou, K. C., Ensemble classifier for protein fold pattern recognition. *Bioinformatics* **2006**, 22, (14), 1717-22.
54. Chou, K. C.; Shen, H. B., Predicting eukaryotic protein subcellular location by fusing optimized evidence-theoretic K-Nearest Neighbor classifiers. *J Proteome Res* **2006**, 5, (8), 1888-97.
55. Chou, K. C.; Shen, H. B., Large-scale predictions of gram-negative bacterial protein subcellular locations. *J Proteome Res* **2006**, 5, (12), 3420-8.
56. Shen, H. B.; Chou, K. C., Virus-PLoc: a fusion classifier for predicting the subcellular localization of viral proteins within host and virus-infected cells. *Biopolymers* **2007**, 85, (3), 233-40.
57. Shen, H. B.; Chou, K. C., Hum-mPLoc: an ensemble classifier for large-scale human protein subcellular location prediction by incorporating samples with multiple sites. *Biochem Biophys Res Commun* **2007**, 355, (4), 1006-11.
58. Xiao, X.; Shao, S. H.; Huang, Z. D.; Chou, K. C., Using pseudo amino acid composition to predict protein structural classes: approached with complexity measure factor. *J Comput Chem* **2006**, 27, (4), 478-82.
59. Xiao, X.; Shao, S.; Ding, Y.; Huang, Z.; Chou, K. C., Using cellular automata images and pseudo amino acid composition to predict protein subcellular location. *Amino Acids* **2005**.
60. Chou, K. C.; Shen, H. B., MemType-2L: a web server for predicting membrane proteins and their types by incorporating evolution information through Pse-PSSM. *Biochem Biophys Res Commun* **2007**, 360, (2), 339-45.
61. Chou, K. C.; Shen, H. B., Euk-mPLoc: a fusion classifier for large-scale eukaryotic protein subcellular location prediction by incorporating multiple sites. *J Proteome Res* **2007**, 6, (5), 1728-34.
62. Pu, X.; Guo, J.; Leung, H.; Lin, Y., Prediction of membrane protein types from sequences and position-specific scoring matrices. *J Theor Biol* **2007**, 247 (2007) 247 (2007) 259-265, 259-265.
63. Zhou, X. B.; Chen, C.; Li, Z. C.; Zou, X. Y., Using Chou's amphiphilic pseudo-amino acid composition and support vector machine for prediction of enzyme subfamily classes. *J Theor Biol* **2007**, 248, 546-551.
64. Pan, Y. X.; Zhang, Z. Z.; Guo, Z. M.; Feng, G. Y.; Huang, Z. D.; He, L., Application of pseudo amino acid composition for predicting protein subcellular location: stochastic signal processing approach. *J Protein Chem* **2003**, 22, 395-402.
65. Chen, C.; Tian, Y. X.; Zou, X. Y.; Cai, P. X.; Mo, J. Y., Using pseudo-amino acid composition and support vector machine to predict protein structural class. *J Theor Biol* **2006**, 243, 444-448.
66. Chen, C.; Zhou, X.; Tian, Y.; Zou, X.; Cai, P., Predicting protein structural class with pseudo-amino acid composition and support vector machine fusion network. *Anal Biochem* **2006**, 357, 116-121.
67. Du, P.; Li, Y., Prediction of protein submitochondria locations by hybridizing pseudo-amino acid composition with various physicochemical features of segmented sequence. *BMC Bioinformatics* **2006**, 7, 518.
68. Mondal, S.; Bhavna, R.; Babu, R. M.; Ramakumar, S., Pseudo amino acid composition and multi-class support vector machines approach for conotoxin superfamily classification. *J Theor Biol* **2006**, 243, 252-260.

69. Chen, Y. L.; Li, Q. Z., Prediction of apoptosis protein subcellular location using improved hybrid approach and pseudo amino acid composition. *J Theor Biol* **2007**, 248 377–381.
70. Chen, Y. L.; Li, Q. Z., Prediction of the subcellular location of apoptosis proteins. *J Theor Biol* **2007**, 245 775-783.
71. Lin, H.; Li, Q. Z., Predicting conotoxin superfamily and family by using pseudo amino acid composition and modified Mahalanobis discriminant. *Biochem Biophys Res Commun* **2007**, 354, 548-551.
72. Lin, H.; Li, Q. Z., Using Pseudo Amino Acid Composition to Predict Protein Structural Class: Approached by Incorporating 400 Dipeptide Components. *J Comput Chem* **2007**, 28, 1463-1466.
73. Shi, J. Y.; Zhang, S. W.; Pan, Q.; Cheng, Y.-M.; Xie, J., Prediction of protein subcellular localization by support vector machines using multi-scale energy and pseudo amino acid composition. *Amino Acids* **2007**, 33, 69-74.
74. Kurgan, L. A.; Stach, W.; Ruan, J., Novel scales based on hydrophobicity indices for secondary protein structure. *J Theor Biol* **2007**, 248, 354–366.
75. Du, Q.; Mezey, P. G.; Chou, K. C., Heuristic molecular lipophilicity potential (HMLP): a 2D-QSAR study to LADH of molecular family pyrazole and derivatives. *J Comput Chem* **2005**, 26, (5), 461-70.
76. Du, Q. S.; Huang, R. B.; Wei, Y. T.; Du, L. Q.; Chou, K. C., Multiple field three dimensional quantitative structure-activity relationship (MF-3D-QSAR). *J Comput Chem* **2007**.
77. Caballero, J.; Fernandez, L.; Abreu, J. I.; Fernandez, M., Amino Acid Sequence Autocorrelation vectors and ensembles of Bayesian-Regularized Genetic Neural Networks for prediction of conformational stability of human lysozyme mutants. *J Chem Inf Model* **2006**, 46, (3), 1255-68.
78. Liao, B.; Ding, K., Graphical approach to analyzing DNA sequences. *J Comput Chem* **2005**, 26, (14), 1519-23.
79. Liao, B.; Wang, T. M., Analysis of similarity/dissimilarity of DNA sequences based on nonoverlapping triplets of nucleotide bases. *J Chem Inf Comput Sci* **2004**, 44, (5), 1666-70.
80. Liao, B.; Wang, T. M., New 2D graphical representation of DNA sequences. *J Comput Chem* **2004**, 25, (11), 1364-8.
81. Liao, B.; Xiang, X.; Zhu, W., Coronavirus phylogeny based on 2D graphical representation of DNA sequence. *J Comput Chem* **2006**, 27, (11), 1196-1202.
82. Yu-Hua, Y.; Liao, B.; Tian-Ming, W., A 2D graphical representation of RNA secondary structures and the analysis of similarity/dissimilarity based on it. *Journal of molecular structure:THEOCHEM* **2005**, 755, 131-136.
83. Liao, B.; Wang, T., A 3D Graphical representation of RNA secondary structure. *J Biomol Struct Dynamics* **2004**, 21, 827-832.
84. Liao, B.; Ding, K.; Wang, T., On A Six-Dimensional Representation of RNA Secondary Structures. *J Biomol Struct Dynamics* **2005**, 22, 455-464.
85. Liao, B.; Wang, T.; Ding, K., On A Seven-Dimensional Representation of RNA Secondary Structures. *Molecular Simulation* **2005**, 31, (14 ), 1063-1071.
86. Liao, B.; Luo, J.; Li, R.; Zhu, W., RNA Secondary structure 2D graphical representation without degeneracy. *International Journal of Quantum Chemistry* **2006** 106 (8), 1749-1755.
87. Zhu, W.; Liao, B.; Ding, K., A condensed 3D Graphical representation of RNA secondary structures. *Journal of Molecular Structure: THEOCHEM* **2005** 757 193-198.
88. Song, J.; Tang, H., A new 2-D graphical representation of DNA sequences and their numerical characterization. *J Biochem Biophys Methods* **2005**, 63, (3), 228-39.
89. Randic, M.; Zupan, J., Highly compact 2D graphical representation of DNA sequences. *SAR QSAR Environ Res* **2004**, 15, (3), 191-205.
90. Randic, M.; Balaban, A. T., On a four-dimensional representation of DNA primary sequences. *J Chem Inf Comput Sci* **2003**, 43, (2), 532-9.
91. Liu, Y.; Guo, X.; Xu, J.; Pan, L.; Wang, S., Some notes on 2-D graphical representation of DNA sequence. *J Chem Inf Comput Sci* **2002**, 42, (3), 529-33.

92. Randić, M.; Vračko, M., On the similarity of DNA primary sequences. *J Chem Inf Comput Sci* **2000**, 40, (3), 599-606.
93. Agüero-Chapin, G.; Gonzalez-Diaz, H.; Molina, R.; Varona-Santos, J.; Uriarte, E.; Gonzalez-Diaz, Y., Novel 2D maps and coupling numbers for protein sequences. The first QSAR study of polygalacturonases; isolation and prediction of a novel sequence from *Psidium guajava* L. *FEBS lett* **2006**, 580, 723-730.
94. Randić, M.; Vračko, M.; Nandy, A.; Basak, S. C., On 3-D Graphical Representation of DNA Primary Sequences and Their Numerical Characterization. *J Chem Inf Comput Sci* **2000**, 40, 1235-1244.
95. Randić, M.; Vračko, M.; Nandy, A.; Basak, S. C., On 3-D graphical representation of DNA primary sequences and their numerical characterization. *J Chem Inf Comput Sci* **2000**, 40, (5), 1235-44.
96. Nandy, A., Two-dimensional graphical representation of DNA sequences and intron-exon discrimination in intron-rich sequences. *Comput Appl Biosci* **1996**, 12, (1), 55-62.
97. Nandy, A., Recent investigations into global characteristics of long DNA sequences. *Indian J Biochem Biophys* **1994**, 31, (3), 149-55.
98. Nandy, A.; Basak, S. C., Simple numerical descriptor for quantifying effect of toxic substances on DNA sequences. *J Chem Inf Comput Sci* **2000**, 40, (4), 915-9.
99. Zupan, J.; Randić, M., Algorithm for coding DNA sequences into "spectrum-like" and "zigzag" representations. *J Chem Inf Model* **2005**, 45, (2), 309-13.
100. Woodcock, S.; Mornon, J. P.; Henrissat, B., Detection of secondary structure elements in proteins by hydrophobic cluster analysis. *Protein Eng* **1992**, 5, (7), 629-35.
101. Randić, M., 2-D graphical representation of proteins based on virtual genetic code. *SAR QSAR Environ Res* **2004**, 15, (3), 147-57.
102. Gonzalez-Diaz, H.; Uriarte, E.; Ramos de Armas, R., Predicting stability of Arc repressor mutants with protein stochastic moments. *Bioorg Med Chem* **2005**, 13, (2), 323-31.
103. Vilar, S.; Santana, L.; Uriarte, E., Probabilistic neural network model for the in silico evaluation of anti-HIV activity and mechanism of action. *J Med Chem* **2006**, 49, (3), 1118-1124.
104. Vilar, S.; Estrada, E.; Uriarte, E.; Santana, L.; Gutierrez, Y., In silico studies toward the discovery of new anti-HIV nucleoside compounds through the use of TOPS-MODE and 2D/3D connectivity indices. 2. Purine derivatives. *J Chem Inf Model* **2005**, 45, (2), 502-14.
105. Gonzalez, M. P.; del Carmen Teran Moldes, M., A TOPS-MODE approach to predict affinity for A1 adenosine receptors. 2-(Arylamino)adenosine analogues. *Bioorg Med Chem* **2004**, 12, (11), 2985-93.
106. Gonzalez, M. P.; Helguera, A. M.; Cabrera, M. A., Quantitative structure-activity relationship to predict toxicological properties of benzene derivative compounds. *Bioorg Med Chem* **2005**, 13, (5), 1775-81.
107. Gonzalez, M. P.; Moldes del Carmen Teran, M., A TOPS-MODE approach to predict adenosine kinase inhibition. *Bioorg Med Chem Lett* **2004**, 14, (12), 3077-9.
108. Gonzalez, M. P.; Teran, C.; Teijeira, M., A topological function based on spectral moments for predicting affinity toward A3 adenosine receptors. *Bioorg Med Chem Lett* **2006**, 16, (5), 1291-6.
109. Estrada, E., On the topological sub-structural molecular design (TOSS-MODE) in QSPR/QSAR and drug design research. *SAR QSAR Environ Res* **2000**, 11, (1), 55-73.
110. Estrada, E., Characterization of the folding degree of proteins. *Bioinformatics* **2002**, 18, 697-704.
111. Estrada, E.; Gonzalez, H., What are the limits of applicability for graph theoretic descriptors in QSPR/QSAR? Modeling dipole moments of aromatic compounds with TOPS-MODE descriptors. *J Chem Inf Comput Sci* **2003**, 43, (1), 75-84.
112. Estrada, E.; Uriarte, E.; Vilar, S., Effect of Protein Backbone Folding on the Stability of Protein-Ligand Complexes. *Journal of Proteome Research* **2006**, 5, 105-111.
113. Du, Q. S.; Wei, Y. T.; Pang, Z. W.; Chou, K. C.; Huang, R. B., Predicting the affinity of epitope-peptides with class I MHC molecule HLA-A\*0201: an application of amino acid-based peptide prediction. *Protein Eng Des Sel* **2007**, 20, (9), 417-23.
114. Du, Q. S.; Huang, R. B.; Wei, Y. T.; Wang, C. H.; Chou, K. C., Peptide reagent design based on physical and chemical properties of amino acid residues. *J Comput Chem* **2007**, 28, (12), 2043-50.

115. Meneses-Marcel, A.; Marrero-Ponce, Y.; Machado-Tugores, Y.; Montero-Torres, A.; Pereira, D. M.; Escario, J. A.; Nogal-Ruiz, J. J.; Ochoa, C.; Aran, V. J.; Martinez-Fernandez, A. R.; Garcia Sanchez, R. N., A linear discrimination analysis based virtual screening of trichomonacidal lead-like compounds: outcomes of in silico studies supported by experimental results. *Bioorg Med Chem Lett* **2005**, 15, (17), 3838-43.
116. Marrero-Ponce, Y.; Castillo-Garit, J. A.; Olazabal, E.; Serrano, H. S.; Morales, A.; Castanedo, N.; Ibarra-Velarde, F.; Huesca-Guillen, A.; Sanchez, A. M.; Torrens, F.; Castro, E. A., Atom, atom-type and total molecular linear indices as a promising approach for bioorganic and medicinal chemistry: theoretical and experimental assessment of a novel method for virtual screening and rational design of new lead anthelmintic. *Bioorg Med Chem* **2005**, 13, (4), 1005-20.
117. Marrero-Ponce, Y.; Diaz, H. G.; Zaldivar, V. R.; Torrens, F.; Castro, E. A., 3D-chiral quadratic indices of the 'molecular pseudograph's atom adjacency matrix' and their application to central chirality codification: classification of ACE inhibitors and prediction of sigma-receptor antagonist activities. *Bioorg Med Chem* **2004**, 12, (20), 5331-42.
118. Garcia-Garcia, A.; Galvez, J.; de Julian-Ortiz, J. V.; Garcia-Domenech, R.; Munoz, C.; Guna, R.; Borrás, R., New agents active against Mycobacterium avium complex selected by molecular topology: a virtual screening method. *J Antimicrob Chemother* **2004**, 53, (1), 65-73.
119. Gozalbes, R.; Galvez, J.; Garcia-Domenech, R.; Derouin, F., Molecular search of new active drugs against Toxoplasma gondii. *SAR QSAR Environ Res.* **1999**, 10, (1), 47-60.
120. Ponce, Y. M.; Diaz, H. G.; Zaldivar, V. R.; Torrens, F.; Castro, E. A., 3D-chiral quadratic indices of the 'molecular pseudograph's atom adjacency matrix' and their application to central chirality codification: classification of ACE inhibitors and prediction of sigma-receptor antagonist activities. *Bioorg Med Chem* **2004**, 12, (20), 5331-42.
121. Santana, L.; Uriarte, E.; González-Díaz, H.; Zagotto, G.; Soto-Otero, R.; Mendez-Alvarez, E., A QSAR model for in silico screening of MAO-A inhibitors. Prediction, synthesis, and biological assay of novel coumarins. *J Med Chem* **2006**, 49, (3), 1149-56.
122. Zdobnov, E. M.; R., A., InterProScan - an integration platform for the signature-recognition methods in InterPro. *Bioinformatics* **2001**, 17, (9), 847-848.
123. Gonzalez-Diaz, H.; Vina, D.; Santana, L.; de Clercq, E.; Uriarte, E., Stochastic entropy QSAR for the in silico discovery of anticancer compounds: Prediction, synthesis, and in vitro assay of new purine carbanucleosides. *Bioorg Med Chem* **2005**.
124. González-Díaz, H., Molina, R.R., Uriarte, E, Stochastic molecular descriptors for polymers. 1. Modelling the properties of icosahedral viruses with 3D-Markovian negentropies. *Polymer* **2003**, (45), 3845-3853.
125. Gonzalez-Diaz, H.; Agüero, G.; Cabrera, M. A.; Molina, R.; Santana, L.; Uriarte, E.; Delogu, G.; Castanedo, N., Unified Markov thermodynamics based on stochastic forms to classify drugs considering molecular structure, partition system, and biological species: distribution of the antimicrobial G1 on rat tissues. *Bioorg Med Chem Lett* **2005**, 15, (3), 551-7.
126. Gonzalez-Diaz, H.; de Armas, R. R.; Molina, R., Vibrational Markovian modelling of footprints after the interaction of antibiotics with the packaging region of HIV type 1. *Bull Math Biol* **2003**, 65, (6), 991-1002.
127. Dobson, P. D.; Doig, A. J., Distinguishing Enzyme Structures from Non-enzymes Without Alignments. *J. Mol. Biol.* **2003**, 330, 771-783.
128. Yuan, Z., Prediction of protein subcellular locations using Markov chain models. *FEBS Lett* **1999**, 451, (1), 23-6.
129. Gonzalez-Diaz, H.; Cruz-Monteagudo, M.; Vina, D.; Santana, L.; Uriarte, E.; De Clercq, E., QSAR for anti-RNA-virus activity, synthesis, and assay of anti-RSV carbonucleosides given a unified representation of spectral moments, quadratic, and topologic indices. *Bioorg Med Chem Lett* **2005**, 15, (6), 1651-7.

130. Ramos de Armas, R.; Gonzalez-Diaz, H.; Molina, R.; Perez Gonzalez, M.; Uriarte, E., Stochastic-based descriptors studying peptides biological properties: modeling the bitter tasting threshold of dipeptides. *Bioorg Med Chem* **2004**, 12, (18), 4815-22.
131. Gonzalez-Diaz, H.; Uriarte, E., Biopolymer stochastic moments. I. Modeling human rhinovirus cellular recognition with protein surface electrostatic moments. *Biopolymers* **2005**, 77, (5), 296-303.
132. Randić, M., Graphical representation of DNA as a 2-D map. *Chem. Phys. Lett.* **2004**, (386), 468–471.
133. Jacchieri, S. G., Mining combinatorial data in protein sequences and structures. *Molecular Diversity* **2000**, (5), 145–152.
134. Ramos de Armas, R.; Gonzalez-Diaz, H.; Molina, R.; Uriarte, E., Markovian Backbone Negentropies: Molecular descriptors for protein research. I. Predicting protein stability in Arc repressor mutants. *Proteins* **2004**, 56, (4), 715-23.
135. González-Díaz, H.; Hernández, I. *BIOMARKS 1.0*; Sta Clara: 2005.
136. STATISTICA for Windows release 6.0. Statsoft Inc. 2001.
137. Mc Farland, J. W.; Gans, D. J., *Cluster Significance Analysis. In Method and Principles in Medicinal Chemistry.* van Waterbeemd, H ed.; VCH: Weinheim, Germany, 1995; Vol. 2, p 295-307.
138. Kutner, M. H.; Nachtsheim, C. J.; Neter, J.; Li, W., Standardized Multiple Regression Model. In *Applied Linear Statistical Models*, Fifth ed.; McGraw Hill: New York, 2005; pp 271-277.
139. Rotondo, G.; Gillespie, M.; Frendewey, D., Rescue of the fission yeast snRNA synthesis mutant *snm1* by overexpression of the double-strand-specific Pac1 ribonuclease. *Mol Gen Genet* **1995**, 247, (6), 698-708.
140. Sambrook, J.; Fritsh, E. F.; Maniatis, T., *Molecular Cloning: A laboratory Manual.* second ed.; Cold Spring Harbor Laboratory Press: USA, 1989.
141. Pawar, A. P.; Dubay, K. F.; Zurdo, J.; Chiti, F.; Vendruscolo, M.; Dobson, C. M., Prediction of "aggregation-prone" and "aggregation-susceptible" regions in proteins associated with neurodegenerative diseases. *J Mol Biol* **2005**, 350, (2), 379-92.
142. Chiti, F.; Stefani, M.; Taddei, N.; Ramponi, G.; Dobson, C. M., Rationalization of the effects of mutations on peptide and protein aggregation rates. *Nature* **2003**, 424, (6950), 805-8.
143. Gonzalez-Diaz, H.; Agüero-Chapin, G.; Varona-Santos, J.; Molina, R.; de la Riva, G.; Uriarte, E., 2D RNA-QSAR: assigning ACC oxidase family membership with stochastic molecular descriptors; isolation and prediction of a sequence from *Psidium guajava* L. *Bioorg Med Chem Lett* **2005**, 15, (11), 2932-7.
144. Di Francesco, V. M., P. J.; Garnier, J, FORESST: fold recognition from secondary structure predictions of proteins. *Bioinformatics* **1999**, 15, (2), 131-40.
145. Chou, K. C., Prediction of signal peptides using scaled window. *Peptides* **2001**, 22, (12), 1973-1979.
146. Chou, K. C.; Shen, H. B., Signal-CF: a subsite-coupled and window-fusing approach for predicting signal peptides. *Biochem Biophys Res Commun* **2007**, 357, (3), 633-40.
147. Chou, K.-C., Prediction and classification of alpha-turn types. *Biopolymers* **1997**, 42, (7), 837-53
148. Chou, K. C., Prediction of beta-turns. *J Pept Res* **1997**, 49, (2), 120-44.
149. Chou, K. C., Prediction of tight turns and their types in proteins. *Anal Biochem* **2000**, 286, (1), 1-16.
150. Chou, K. C., A sequence-coupled vector-projection model for predicting the specificity of GalNAc-transferase. *Protein Sci* **1995**, 4, (7), 1365-83.
151. Chou, K. C., A vectorized sequence-coupling model for predicting HIV protease cleavage sites in proteins. *J Biol Chem* **1993**, 268, (23), 16938-48.
152. Chou, J. J., A formulation for correlating properties of peptides and its application to predicting human immunodeficiency virus protease-cleavable sites in proteins. *Biopolymers* **1993**, 33, 1405-1414.
153. Chou, J. J., Predicting cleavability of peptide sequences by HIV protease via correlation-angle approach. *J Protein Chem* **1993**, 12, 291-302.
154. Chou, K. C., Prediction of human immunodeficiency virus protease cleavage sites in proteins. *Anal Biochem* **1996**, 233, (1), 1-14.

155. Chou, K. C.; Wei, D. Q.; Zhong, W. Z., Binding mechanism of coronavirus main proteinase with ligands and its implication to drug design against SARS. *Biochem Biophys Res Commun* **2003**, 308, (1), 148-51.
156. Chou, K. C., Structural bioinformatics and its impact to biomedical science. *Curr Med Chem* **2004**, 11, (16), 2105-34.
157. Sirois, S.; Sing, T.; Chou, K. C., HIV-1 gp120 V3 loop for structure-based drug design. *Curr Protein Pept Sci* **2005**, 6, (5), 413-22.
158. Chou, K. C.; Wei, D. Q.; Du, Q. S.; Sirois, S.; Zhong, W. Z., Progress in computational approach to drug development against SARS. *Curr Med Chem* **2006**, 13, (27), 3263-70.
159. Du, Q.; Wang, S.; Jiang, Z.; Gao, W.; Li, Y.; Wei, D.; Chou, K. C., Application of bioinformatics in search for cleavable peptides of SARS-CoV M(pro) and chemical modification of octapeptides. *Med Chem* **2005**, 1, (3), 209-13.
160. Du, Q. S.; Sun, H.; Chou, K. C., Inhibitor design for SARS coronavirus main protease based on "distorted key theory". *Med Chem* **2007**, 3, (1), 1-6.
161. Wei, D. Q.; Zhang, R.; Du, Q. S.; Gao, W. N.; Li, Y.; Gao, H.; Wang, S. Q.; Zhang, X.; Li, A. X.; Sirois, S.; Chou, K. C., Anti-SARS drug screening by molecular docking. *Amino Acids* **2006**, 31, (1), 73-80.
162. Zhang, R.; Wei, D. Q.; Du, Q. S.; Chou, K. C., Molecular modeling studies of peptide drug candidates against SARS. *Med Chem* **2006**, 2, (3), 309-14.
163. Wang, S. Q.; Du, Q. S.; Zhao, K.; Li, A. X.; Wei, D. Q.; Chou, K. C., Virtual screening for finding natural inhibitor against cathepsin-L for SARS therapy. *Amino Acids* **2007**, 33, (1), 129-35.
164. Van Waterbeemd, H., *Chemometric methods in molecular design*. Wiley-VCH: New York, 1995; Vol. 2.
165. Morales Helguera A, C. P. M. A., Pérez González M, Molina Ruiz R, González Díaz H., A topological substructural approach applied to the computational prediction of rodent carcinogenicity. *Bioorganic & Medicinal Chemistry* **2005**, 13, 2477-2488.
166. Marrero-Ponce, Y.; Montero-Torres, A.; Zaldivar, C. R.; Veitia, M. I.; Perez, M. M.; Sanchez, R. N., Non-stochastic and stochastic linear indices of the 'molecular pseudograph's atom adjacency matrix': application to 'in silico' studies for the rational discovery of new antimalarial compounds. *Bioorg Med Chem* **2005**, 13, (4), 1293-304.
167. Marrero-Ponce, Y.; Medina-Marrero, R.; Torrens, F.; Martinez, Y.; Romero-Zaldivar, V.; Castro, E. A., Atom, atom-type, and total nonstochastic and stochastic quadratic fingerprints: a promising approach for modeling of antibacterial activity. *Bioorg Med Chem* **2005**, 13, (8), 2881-99.
168. Marrero-Ponce, Y.; Medina-Marrero, R.; Castillo-Garit, J. A.; Romero-Zaldivar, V.; Torrens, F.; Castro, E. A., Protein linear indices of the 'macromolecular pseudograph alpha-carbon atom adjacency matrix' in bioinformatics. Part 1: prediction of protein stability effects of a complete set of alanine substitutions in Arc repressor. *Bioorg Med Chem* **2005**, 13, (8), 3003-15.
169. Swets, J. A., Measuring the accuracy of diagnostic systems. *Science* **1988**, 240, (4857), 1285-93.
170. Dieterich, G.; Karst, U.; Wehland, J.; Jansch, L., MineBlast: a literature presentation service supporting protein annotation by data mining of BLAST results. *Bioinformatics* **2005**, 21, (16), 3450-1.
171. Balakrishnan, R.; Christie, K. R.; Costanzo, M. C.; Dolinski, K.; Dwight, S. S.; Engel, S. R.; Fisk, D. G.; Hirschman, J. E.; Hong, E. L.; Nash, R.; Oughtred, R.; Skrzypek, M.; Theesfeld, C. L.; Binkley, G.; Dong, Q.; Lane, C.; Sethuraman, A.; Weng, S.; Botstein, D.; Cherry, J. M., Fungal BLAST and Model Organism BLASTP Best Hits: new comparison resources at the Saccharomyces Genome Database (SGD). *Nucleic Acids Res* **2005**, 33, (Database issue), D374-7.
172. Zhou, Y.; Huang, G. M.; Wei, L., UniBLAST: a system to filter, cluster, and display BLAST results and assign unique gene annotation. *Bioinformatics* **2002**, 18, (9), 1268-9.
173. Zhang, J.; Madden, T. L., PowerBLAST: a new network BLAST application for interactive or automated sequence analysis and annotation. *Genome Res* **1997**, 7, (6), 649-56.
174. Song, J.; Burrage, K.; Yuan, Z.; Huber, T., Prediction of cis/trans isomerization in proteins using PSI-BLAST profiles and secondary structure information. *BMC Bioinformatics* **2006**, 7, (1), 124.



175. Rettie, A. E.; Jones, J. P., Clinical and toxicological relevance of CYP2C9: drug-drug interactions and pharmacogenetics. *Annu Rev Pharmacol Toxicol* **2005**, 45, 477-94.
176. Zehetner, G., OntoBlast function: From sequence similarities directly to potential functional annotations by ontology terms. *Nucleic Acids Res* **2003**, 31, (13), 3799-803.
177. Yang, A. S., Structure-dependent sequence alignment for remotely related proteins. *Bioinformatics* **2002**, 18, (12), 1658-65.
178. Lee, C.; Grasso, C.; Sharlow, M. F., Multiple sequence alignment using partial order graphs. *Bioinformatics* **2002**, 18, (3), 452-64.
179. Marrero-Ponce, Y.; Medina-Marrero, R.; Castro, A. E.; Ramos de Armas, R.; González-Díaz, H.; Romero-Zaldivar, V.; Torrens, F., Protein Quadratic Indices of the “Macromolecular Pseudograph’s  $\alpha$ -Carbon Atom Adjacency Matrix”. 1. Prediction of Arc Repressor Alanine-mutant’s Stability. *Molecules* **2004**, 9, 1124–1147.