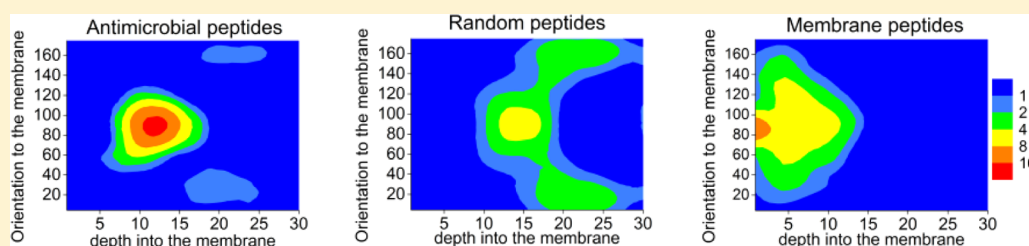# Prediction of Linear Cationic Antimicrobial Peptides Based on Characteristics Responsible for Their Interaction with the Membranes

Boris Vishnepolsky* and Malak Pirtskhalava*

I. Beritashvili Center of Experimental Biomedicine, Laboratory of Bioinformatics, Tbilisi 0160, Georgia

**ABSTRACT:** Most available antimicrobial peptides (AMP) prediction methods use common approach for different classes of AMP. Contrary to available approaches, we suggest that a strategy of prediction should be based on the fact that there are several kinds of AMP that vary in mechanisms of action, structure, mode of interaction with membrane, etc. According to our suggestion for each kind of AMP, a particular approach has to be developed in order to get high efficacy. Consequently, in this paper, a particular but the biggest class of AMP, linear cationic antimicrobial peptides (LCAP), has been considered and a newly developed simple method of LCAP prediction described. The aim of this study is the development of a simple method of discrimination of AMP from non-AMP, the efficiency of which will be determined by efficiencies of selected descriptors only and comparison the results of the discrimination procedure with the results obtained by more complicated discriminative methods. As descriptors the physicochemical characteristics responsible for capability of the peptide to interact with an anionic membrane were considered. The following characteristics such as hydrophobicity, amphiphaticity, location of the peptide in relation to membrane, charge density, propensies to disordered structure and aggregation were studied. On the basis of these characteristics, a new simple algorithm of prediction is developed and evaluation of efficacies of the characteristics as descriptors performed. The results show that three descriptors, hydrophobic moment, charge density and location of the peptide along the membranes, can be used as discriminators of LCAPs. For the training set, our method gives the same level of accuracy as more complicated machine learning approaches offered as CAMP database service tools. For the test set accuracy obtained by our method gives even higher value than the one obtained by CAMP prediction tools. The AMP prediction tool based on the considered method is available at http://www.biomedicine.org.ge/dbaasp/.

## INTRODUCTION

Antimicrobial peptides (AMP) are small peptides of low length, which interact with the bacterial cells and kill them. The great interest in these proteins is explained by their possible use for clinical purposes as a substitute for conventional antibiotics when resistance takes place.[1] Most of AMP act directly on the bacterial membrane, consequently it is difficult for bacteria to develop immunity against antimicrobial peptides.[2] Recently, there has been a large number of both theoretical and experimental studies that were focused on the properties of AMP, their mechanism of action and the design of novel peptides (see, for example, ref 3). Of particular interest are in silico methods of research of AMP that allow the capability to both predict the antimicrobial activity of the peptides based on their sequence and to serve as the first step to design new antimicrobial peptides. Methods for predicting AMP are based on some general properties that distinguish AMP from similar peptides that do not have antimicrobial activity.

Available prediction methods are generally based on discriminative analysis and essentially machine learning methods.[4−12] These methods, as a positive training set, have used a full set of antimicrobial peptide sequences, not taking into account variation in mechanisms of action, structure, mode of interaction with membrane and other differences. Contrary to available approaches, we think that strategy of prediction should be based on the fact that there are at least four kinds of AMPs for which four independent algorithms of prediction have to be developed in order to get high efficacy. For these four types of AMPs, we can consider: linear cationic antimicrobial peptides (LCAP), cationic peptides stabilizing structure by interchain covalent bond (CCP), peptides rich in proline and arginine (PRP) and anionic antimicrobial peptides (AAP).

Cationic antimicrobial peptides (CAP) of LCAP type, in addition to positive charge and amphiphilicity, possess simple mechanism of structure stabilization in membrane, hydrogen bonding only.[13,14] Absence of any other stabilization factors gives possibilities to determine the forces governing peptide—lipid or peptide—peptide interactions and predict structure of the peptides in water and membrane environment on a base of only sequence information. Consequently, quantitative characteristics for prediction would be easily revealed on the basis of sequence only. Structures of CAP of CCP type due to interchain bonds are more stable and structurally complicated both in water and in membrane environments. But despite the fact that the forces governing CCP membrane or CCP—CCP interactions are identical to the case of LCAP, complicated 3D structure and lack of information about 3D structure require a principally different approach for the development of CCP prediction algorithm. It is known that peptides of PRP type are penetrating. In other words, they do not destabilize membranes and as a rule, have a target inside cell.[15,16] It is clear that the development of the algorithm for the prediction of CAP of PRP type requires a peculiar approach. For AMP of AAP type, the mode of action principally differs from CAP and the development of the algorithm for prediction AAP indeed requires its own approach. In this work is considered CAP of LCAP type only. According to the available databases,[17] this is the biggest class of antimicrobial peptides.

Prediction accuracy is largely determined by the set of descriptors that can be used in prediction. Most current methods use a large number of characteristics for AMP prediction, using their optimization by machine learning methods, such as artificial neuron networks (ANN) and support vector machines (SVM).[4−12] Meanwhile, the influence of the individual characteristics on the AMP prediction is studied much less extensively. In this paper, we describe the influence of the characteristics that may be responsible for the prediction of LCAP on the basis of their basic function—interaction with the bacterial membrane.

There are a large number of proteins that interact with the membrane also and so resemble AMP in this regard. For instance, the so-called transmembrane proteins are generally inserted into the membrane but without destroying it. It is clear that a selection pressure on sequence random variation directs evolution of peptides with particular function (for instance, transmembrane protein fragments (TMP), LCAP, etc.). So, in order to determine what characteristics efficiently distinguish LCAP from other peptides (other membrane-interactive or nonfunctional (random)), we think that it is reasonable to make comparative analysis of sequences of the three sets of peptides: LCAP, TMP and randomly selected fragments from the soluble proteins (RFP). This work concerns just the comparative analysis of LCAP, TMP and RFP sequences.

Consequently, an attempt to reveal that characteristics that can discriminate antimicrobial peptides from both soluble nonmembrane proteins and transmembrane proteins (or fragments of membrane proteins) has been done. Taking into account the structure of the bacterial membrane, which is an anionic lipid bilayer, amphiphatic in nature, it can be assumed that, for discriminators, the following characteristics are convenient: (1) hydrophobicity, (2) amphiphaticity, (3) charge density, (4) propensity to the aggregation and (5) propensity to disordering. We think that just the values of these characteristics are responsible for: (a) capability of the peptide

to interact with an anionic membrane and (b) the results of interaction (mechanisms of action).

Quantitative estimation of all the characteristics except amphiphaticity requires information on amino acid sequences of the peptides only. Amphiphaticity in addition needs three-dimensional structure information. The exact three-dimensional structure of most linier antimicrobial peptides is unknown. But in the case of linier peptides, based on the theory of Wimly and White[13,14] and the fact that all transmembrane domains of membrane proteins consist mainly of regular secondary structure elements (α-helices or β-sheets saturated with hydrogen bonds), we can assume that the membrane environment will impel the peptide to regular conformation. So, we are motivated estimate in regular structure approximation and evaluate the hydrophobic moment of LCAP in order to see whether the hydrophobic moment can be a good discriminator and which regular structure is more suitable for effective discrimination.

There are various statistical approaches for the prediction of AMP that take into account a number of different characteristics.[4,5,11,18−23] In this paper, our goal is (a) to develop the simplest method of discrimination (based on threshold value only) of AMP from non-AMP, efficiency of which will be determined by efficiencies of selected descriptors only and (b) to compare the results of the discrimination procedure with the results obtained by more refined and complicated discriminative methods such as SVM, ANN, etc.

## ■ METHODS

**Benchmarks.** *Training Sets.* For the analysis of the characteristics, the following benchmarks were selected: set for LCAP, set for randomly selected fragments from the soluble proteins and set of fragments from transmembrane proteins. The LCAP set was selected from APD2 database[17] and consists of 1083 peptides (positive set). To estimate the discriminative efficiency of characteristics, a set of nonantimicrobial peptides has been required. Because there is a small number of peptides with experimentally verified no antimicrobial activity,[5] we have used a voluminous set of random sequences; in other words, a set of sequences with a great variety of functions. So, the last set can be considered as a nonfunctional set on average, as well as nonantimicrobial (negative set). The set of random sequences was selected from an UniProt using the filters: non-AMP, non-membrane and non-secretory proteins. Three such sets were used. The first set (RFP10000) was used for optimizing parameters for various descriptors and consists of randomly chosen fragments in the amount of 10 000 for each length of peptides from 4 to 50 amino acids. The other sets were used for the estimation of the descriptors by receiver operating characteristic (ROC) curves. 500 (for RFP500) and 10 (for RFP10) randomly selected fragments from globular proteins with lengths corresponding to each peptide from LCAP set have been included into these sets. The last set was used for comparing our results with other available prediction tools.

For membrane proteins, a full set of transmembrane (helices) fragments of more than 11 residues from database of transmembrane proteins PDB-TM[24,25] was chosen. This set contains 1691 sequences (TMP set).

*Test Sets.* Two test sets were used for the evaluation of AMP descriptors. The first test set, compiled on the basis of CAMP[11] predicted data set, contained 1153 sequences identified as antimicrobial based on the evidence of similarity or annotations in NCBI as "antimicrobial regions", without experimental

evidence. After eliminating sequences: containing nonstandard amino acids; disulfide bonds; having full negative charges; with the length of more than 50 amino acids and rich in Pro and Arg, only 98 sequences were left (TPS1). TPS1 will serve as an independent positive test data set. An additional test set was obtained from DBAASP database (http://www.biomedicine.org.ge/dbaasp) (TPS2). Only experimentally validated peptides with AMP activity have been included in this set. After peptides that were found in the training LCAP set were excluded, the above-mentioned conditions proposed for TPS1 and peptides with more than 80% homology, the TPS2 set contains 174 peptides. As mentioned above, we could not use any additional independent sample as an independent negative test set. So, we have used RFP10 as a negative data set for the evaluation of the accuracy for the selected descriptors.

**Optimization of the Parameters Defining the Characteristics of AMP.** There is evidence, especially for disulfide-bounded AMP, that despite their short length, they are unions of functional (structural) blocs.[26] So, we can propose that linear peptides are arranged in bloc principle also and not all the considered peptide, but part of it can participate in the interaction with the membrane. Accordingly, for each peptide, the descriptors were calculated for all fragments (windows) of a certain length and peptides are characterized by the particular fragment selected on certain criteria.

The values of different descriptors, in most cases, depend on various parameters, such as length of the fragment for which considered characteristic for the peptide is computed, hydrophobicity scale (see below), etc. It is necessary to choose optimal parameters for descriptors on the base of the LCAP set. Optimization of the descriptors was made by the requirement of increasing the ratio (percent) of the peptides for which the probability of appearance of their sequence as a result of random normal process is less than $P$. The value of $P$ was determined by $z$-score. That is, for each peptide's particular descriptor, its own $z$-score is defined as $z^p_d$ (where d is hydrophobicity, hydrophobic moment and other descriptors, p corresponds to certain peptide defined by its sequence). Main criterion of optimality (MOC) of descriptors was the maximality of the number of peptides from the LCAP set having $z$-score $z^p_d > 2$. The exception was a location of the peptide in relation to membrane, for which optimization has been made differently (see below).

**Hydrophobicity.** The AMP overall hydrophobicity, defined as the sum of transfer (from water into the hydrophobic environment) energy of the residue (hydrophobicity), can be used as an AMP characteristic. In the literature, there is a large number of papers[27−32] that define transfer energies of the amino acids (hydrophobicity scales). The values of the transfer energies in these scales depend on the method of determination and differ from scale to scale. Therefore, the hydrophobicity scale can be used as an optimization parameter for assessing the suitability of the hydrophobicity as AMP characteristics. The following hydrophobicity scales were considered: KD,[27] WW,[28] UHC,[29] Hes,[30] EG[31] and MF.[32]

For each peptide, hydrophobicity was calculated for all fragments of a certain length and the peptide characteristic was defined by the fragment of the highest hydrophobicity. Therefore, peptide fragment length can be the other optimization parameter. The optimal length and hydrophobicity scale were chosen by MOC. Fragment length was varied within the range of 4−50 residues. Moreover, if the peptide length was less than the length of the considered

fragment, hydrophobicity was computed for the full peptide. A similar method for optimizing the fragment length was used for the other descriptors.

**Amphipathicity.** One of the main features of antimicrobial peptides is their amphipathicity.[33] The separation of hydrophobic and hydrophilic regions in these peptides can be realized in one of the two ways: due to the internal 3D structure and by the linear separation that is due to the uneven distribution of hydrophobic and hydrophilic residues along the peptide chain. Accordingly, two characteristics were used for the evaluation of amphipathicity: hydrophobic moment[34] and linear hydrophobic moment (see below).

*Hydrophobic Moment.* Hydrophobic moment was estimated by Eisenberg:[34]

$$\mu = ([\sum_{n=1}^{N} h_n \cdot \sin(\vartheta \cdot n)]^2 + [\sum_{n=1}^{N} h_n \cdot \cos(\vartheta \cdot n)]^2)^{1/2}$$

where $\mu$ is hydrophobic moment of the peptide, contained $N$ amino acids, $h_n$ is the numerical hydrophobicity of the $n$th residue, and $\vartheta$ is turn of the residue along the helix axis.

According to the formula the existence of regular conformation is assumed. As mentioned above LCAP in membrane environment is likely to have regular secondary structure. So, $\vartheta$ is used as a parameter that determines hydrophobic moment. The last parameter was used as optimization parameter and varied from 60 to 180°. Hydrophobicity scale and fragment length were also used as optimization parameters. Optimization of the parameters was carried out by MOC.

*Linear Hydrophobic Moment.* As mentioned above, separation of the hydrophobic and hydrophilic parts may also be carried out due to an uneven distribution of hydrophobic and hydrophilic residues along the peptide chain. To estimate the separation along the chain, we have introduced the characteristic "linear hydrophobic moment", which is defined as follows:

$$M = D(\sum^{+} h_k - \sum^{-} h_k)$$

where

$$D = |\sum^{+} h_k \cdot k / \sum^{+} h_k - \sum^{-} h_k \cdot k / \sum^{-} h_k|$$

Here $D$ is the distance between the centers of hydrophobic and hydrophilic parts of the considered fragment of length $N$; $k = 1$, $N$, $h_k^+$ and $h_k^-$ are the transfer energies of the $k$-th residue from water to the hydrophobic environment under the conditions that $h_k^+ > 0$ corresponds to hydrophobic residue and $h_k^- < 0$ corresponds to hydrophilic residue.

Hydrophobicity scale and fragment length were used as optimization parameters. Optimization of the parameters was carried out by MOC.

**Charge Density.** Cationic antimicrobial peptides at neutral pH have a positive charge due to the large percentage of Lys and Arg, which facilitates them to interact with the negatively charged membrane. So, it is natural to assume that the charge of the peptide can be considered as a characteristic of LCAP. Because electrostatic interaction is long-term, we think that the net charge of the whole peptides determines the results of interaction with membrane. So the charge was calculated for the entire peptide. Charge density determined as full charge

divided by the peptide molecular weight was used as the AMP descriptor.

Initially, for the charge descriptor, full net charge normalized on the peptide length was used, but after suggestion from one of the reviewers, we have found out that charge density determined as full charge divided by the peptide molecular weight gives better discrimination AMP from non-AMP and so we have used charge density as the AMP descriptor.

**Location of the Peptide in Relation to Membrane (LPM).** Mechanism of action of AMP largely depends on their energetically most favorable location within the membrane bilayer. Taking into account the fact that the majority of the LCAP peptides has an α-helical conformation in the membrane environment (see above and the Results section), LPM was described by the penetration depth ($d$), i.e., distance of the geometrical center of peptide helix from membrane surface and angle ($\theta$) between peptide helix axis and perpendicular to the membrane surface. It would be interesting to explore the possibility of using $d$ and $\theta$ as discriminators to distinguish antimicrobial from nonantimicrobial peptides. In contrast to the previous LCAP characteristics, location of the LCAP within the bilayer is an integrated feature that will largely depend on the other previously considered characteristics (hydrophobicity, amphipathicity, charge density). To calculate $d$ and $\theta$, the hydrophobic potential designed by Senes et al.,[35] which represents the energy difference between the residue in water and within the bilayer at a given depth, is used. All calculations of LPM were performed for different fragment (window) lengths and the peptide characteristic was defined by the fragment with minimal energy.

Another (different from MOC) approach was used for the optimization of $d$ and $\theta$. The approach is based on receiver operating characteristic (ROC) curve analysis. ROC curve, which represents the dependence of sensitivity ($S_n$) ($y$-axis) versus $1 -$ specificity ($S_p$) ($x$-axis) was used for quantifying differences of LCAP from membrane proteins and soluble protein fragments. RFP500 and TMP were used as the negative sets. For each LPM, the area under the ROC curve AUC was calculated, defined as AUC_R, relative to the RFP500 negative set and as AUC_T, relative to the TMP negative set. As mentioned above, peptide fragment length varied and so it was used as an optimization parameter. The maximization (AUC_R + AUC_T) value was used to optimize LPM ($d$ and $\theta$). During optimization, $d$ and $\theta$ vary from 0 to 30 Å and 0−180°, respectively and were used as optimization parameters for LPM. Other variables $\delta(d_k)$ and $\delta(\theta_k)$ (for each $k$th $d$ and $\theta$) were used for plotting ROC curve also. For each of the values, $d_k$ and $\theta_k$, $\delta(d_k)$ and $\delta(\theta_k)$ varied and for $i$th their, values $\delta(d_k)_i$ and $\delta(\theta_k)_i$, intervals $d_k \pm \delta(d_k)_i$, $\theta_k \pm \delta(\theta_k)_i$, i.e., $i$th area on the ($d, \theta$) plane is determined. The number of peptides from the positive data sets with energetically most favorable depth and orientation lying within the interval $d_k \pm \delta(d_k)_i$, $\theta_k \pm \delta(\theta_k)_i$ determines sensitivity $S_{ni}$ and the number of peptides from the negative data sets with energetically most favorable depth and orientation lying within the same interval ($d_k \pm \delta(d_k)_i$, $\theta_k \pm \delta(\theta_k)_i$) determines specificity $S_{pi}$. $S_{ni}$ and $S_{pi}$ give $i$th point of the ROC curve. For each length, $d_k$ and $\theta_k$ ROC curves and consequently (AUC_R + AUC_T) values were calculated and maxima among the calculated values correspond to optimums of length, $d_k$ and $\theta_k$.

**Disordering.** It is reasonable to consider such short cationic peptides as LCAP disordered in water environments. Indeed, there is experimentally proved data for many LCAP showing disordered structure in water environments.[36] It is interesting to mention, whether the disordering connected with the short length only, or other causes for structure destabilization exist (for example, total positive charge). Uversky[37] investigated disordered protein and concluded that disordered protein can be predicted on the basis of the estimation of hydrophobic/charge (h/r) ratio. As the LCAPs are characterized by very peculiar balance between hydrophobic and positively charged residues, we think that it will be interesting to estimate if the h/r ratio can be the cause of the disordered structure of LCAP in water environment according to the Uversky's rule. So, we assume that it is interesting to estimate efficiency of the Uversky's relations as discriminator.

According to Uversky's formula, the degree of disordering of globular protein under physiological conditions is defined by the relation

$$S = 2.785\langle H\rangle - 1.151 - \langle R\rangle$$

where $\langle H\rangle$ is the average hydrophobicity of the protein and $\langle R\rangle$ its charge.

Negative values of $S$ correspond to the protein to be disorded.

**Aggregation Propensity.** We have used two descriptors for aggregation propensity; aggregation in solution (in vitro aggregation) and aggregation in bacteria membrane (in vivo aggregation). In vitro aggregation propensity evaluation was made by employing the TANGO software.[38] Tango counts the partition function of the conformational phase space assuming that every segment on the protein populates one state: random coil, β-turn, α-helix, α-helix aggregation and β-sheet aggregation. Therefore, TANGO software can predict aggregation in solution, considering only structural parameters defined by the peptide sequence.

In vivo aggregation was propensity calculated using AGGRESCAN, an algorithm based on an amino acid aggregation-propensity scale derived from in vivo experiments and on the assumption that short and specific sequence stretches modulate protein aggregation. The algorithm can actually predict the aggregation propensity of peptides in the presence of cell material.[39]

**Evaluation of the Efficiency of Characteristics.** Receiver operating characteristic (ROC) curves were used to evaluate the effectiveness of various characteristics. Each point, $i$, of the ROC curve corresponds to values of sensitivity and specificity ($S_{ni}$ and $S_{pi}$), which are calculated for the variable $z_i$ (where $z_i$ changes from $\min(z^p_d)$ to $\max(z^p_d)$ with step 0.1). The number of peptides from the positive data set (LCAP) with $z$-score $z^p_d > z_i$ (for hydrophobic moment, charge density and linear hydrophobic moment) and $z$-score $z^p_d < z_i$ for (hydrophobicity and disordering) determines $S_{ni}$ and the number of peptides from the negative data sets (RFP500 for ROC_R and TMP for ROC_T) with $z$-score $z^p_d > z_i$ (for hydrophobic moment, charge density and linear hydrophobic moment) and $z$-score $z^p_d < z_i$ (for hydrophobicity and disordering) determines $S_{pi}$. Quantitative evaluations of the effectiveness of the characteristics are made on the basis of area under the ROC curve.

**Evaluation of the Prediction Quality.** A threshold for each characteristic was evaluated and the prediction of the existence of antimicrobial activity of the peptide was done on the basis of it. The threshold is determined by a point on the ROC curve closest to the point (0,1). Sensitivity, specificity and accuracy for the thresholds have been evaluated.

The following equations were used for the prediction quality:

**Table 1. Optimal Parameters for Different Descriptors**[a]

|  | hydrophobicity scale | fragment length | angle $\vartheta$ (deg) | MOC[b] No. | MOC[b] % | $d$ | $\theta$ | AUC_R | AUC_T |
|---|---|---|---|---|---|---|---|---|---|
| hydrophobic moment | MF | 24 | 96 | 679 | 62.70 | | | | |
| hydrophobicity | KD | 21 | | 258 | 22.30 | | | | |
| linear hydrophobic moment | EG | 31 | | 119 | 10.00 | | | | |
| LPM | | 17 | | | | 12.9 | 81.0 | 0.76 | 0.78 |

[a]For charge density, disordering, propensity to aggregation in vitro and in vivo parameters optimization was not carried out. [b]A number (No.) and percent (%) of the peptides from the LCAP set, which satisfy MOC criterion.

$$S_n = TP/(TP + FN)$$

$$S_p = TN/(TN + FP)$$

$$BAC = (S_n + S_p)/2$$

$$AC = (TP + TN)/(TP + FN + TN + FP)$$

where $S_n$ is the sensitivity, $S_p$ the specificity, BAC the balanced accuracy and AC the accuracy.

The calculation of balanced accuracy is used for the evaluation of the prediction quality because the negative sets contain more peptides than the positive ones and the balanced accuracy reflects equal influence of positive and negative sets irrespective of the number of contained peptides in them.

## ■ RESULTS AND DISCUSSION

**Optimization of the Descriptors.** The following descriptors were considered: hydrophobic moment, charge density, location of the peptide in relation to membrane (LPM), linear hydrophobic moment, disordering and propensities to in vitro and in vivo aggregation.

The optimization of the most descriptors (except LPM) has been made by MOC criterion (see the Methods section). The corresponding data are given in Table 1. For normal (random) distribution, the probability that $z$-score $> 2$ is equal to 0.02. So, on the basis of the obtained results, we can say that for all optimized descriptors probabilities that $z$-score $> 2$ are higher than expected from fully random processes. It means we can assume some kind of selection pressure on sequence random variation.

For hydrophobic moment, for example, the optimal value of the turn ($\vartheta$) of residue ($\vartheta$ varied from $60°$ to $180°$) in regular structure approximation is $96°$, which shows that the optimal secondary structure for discrimination LCAP from Non-AMP is an $\alpha$-helix. This result can be expected.

For LPM, another criterion of optimality (different from MOC) was used that was based on $d$ and $\theta$ distributions in the considerable set (see the Methods section). Assuming the peptide is $\alpha$-helical (see above) for each peptide, we can calculate the energetically most favorable location ($d$ and $\theta$) of the peptide fragments of the particular length in the membrane. Figure 1a shows a plot of relative density of the orientation and the depth of energetically most favorable fragments of length 17 on the basis of the three considered sets (LCAP, RFP500 and TMP).

From Figure 1a, it is clear that the orientation distribution of peptides in different sets varies from each other. For most LCAP, the more energetically favorable depth is within $8-15$ Å, which corresponds to the boundary between the interface site and the hydrophobic core of the membrane (Figure 1a). It also shows that most of the peptides are located at a relatively small angle to the membrane surface ($\theta \sim 90°$). These results are consistent with experimental data, according to which most of
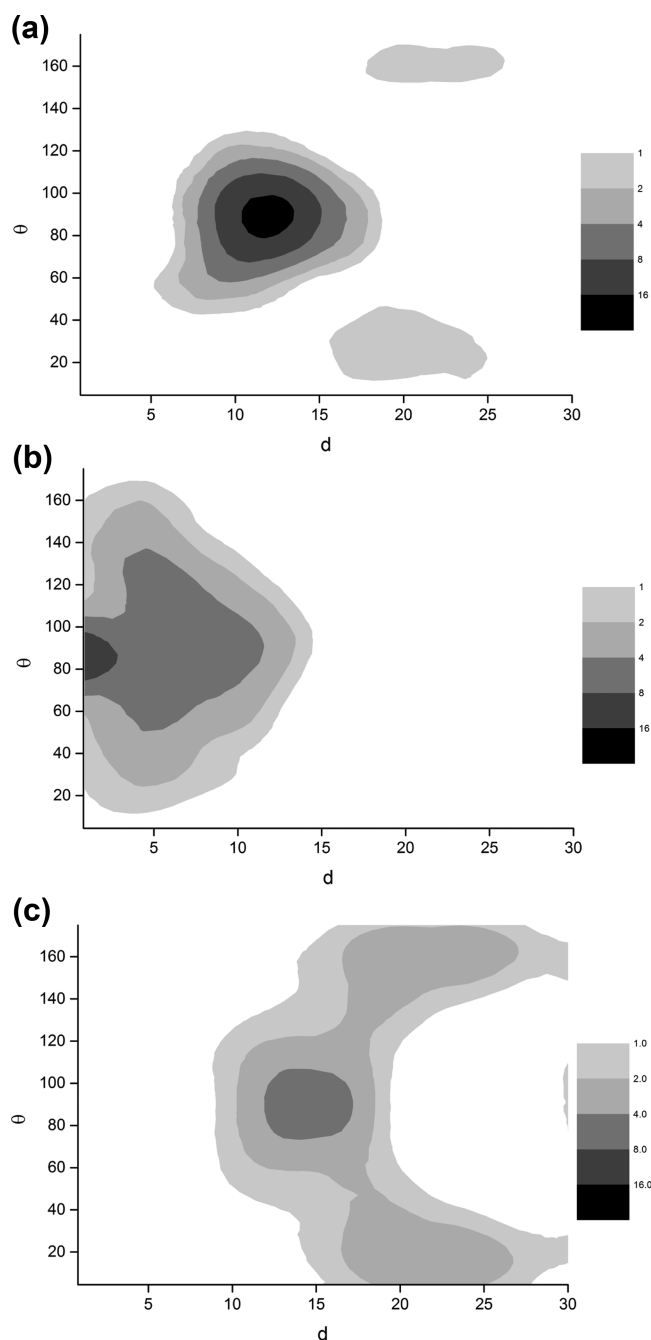


**Figure 1.** Plots of relative density of the orientation ($\theta$) and the depth ($d$) of energetically most favorable fragments of length 17 on the basis of the (A) LCAP, (B) TMP and (C) RFP500 set. The values of the density are given relative to the densities of uniform distribution.

the CAP penetrate into the membrane at a shallow depth parallel to the membrane surface.[40] Maximum density on the
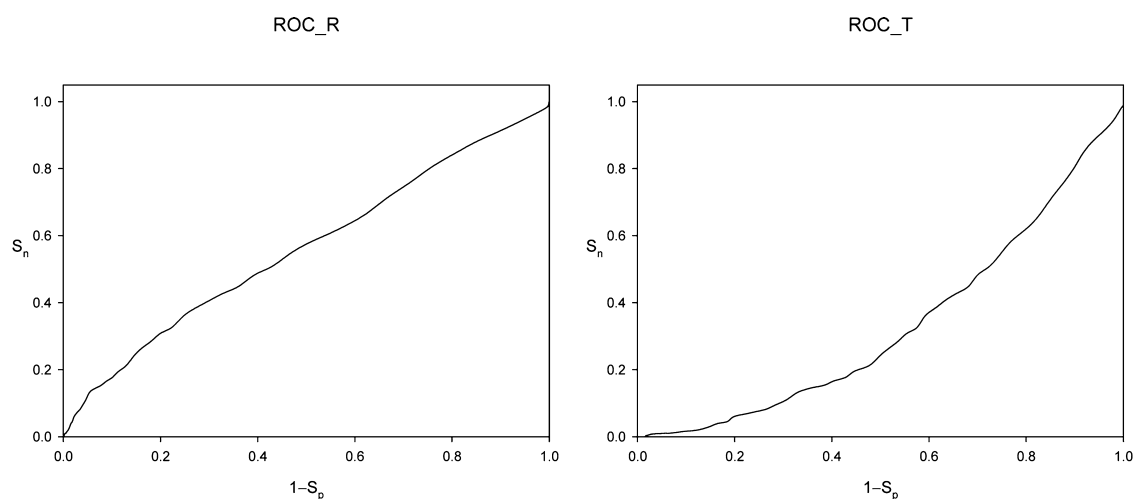
ROC_R  ROC_T



**Figure 2.** ROC curves for evaluation prediction quality of linear hydrophobic moment for training sets: ROC_R corresponds to positive LCAP and negative RFP500 sets and ROC_T corresponds to positive LCAP and negative TMP sets.
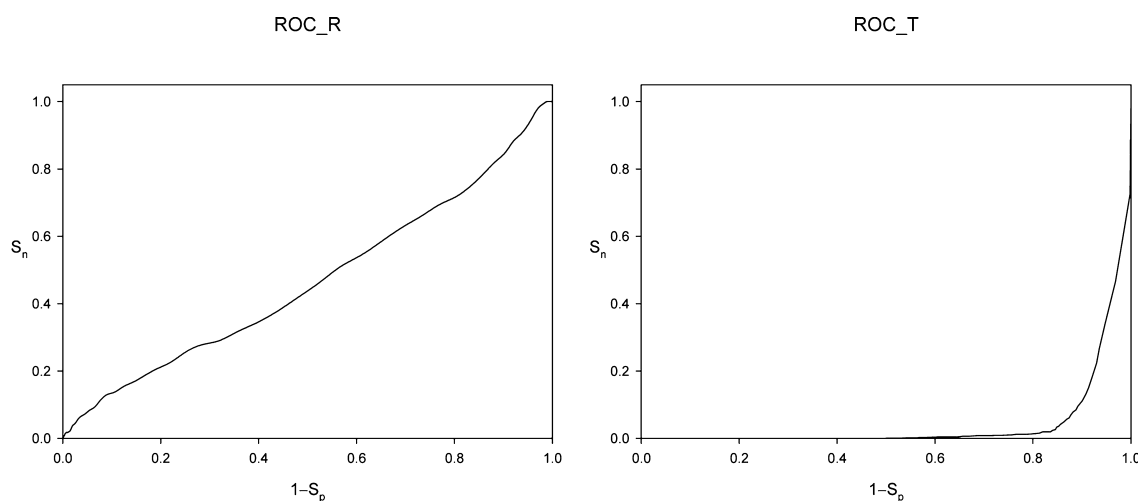
ROC_R  ROC_T



**Figure 3.** ROC curves for evaluation prediction quality of propensity to aggregation in vitro for training sets: ROC_R corresponds to positive LCAP and negative RFP500 sets and ROC_T corresponds to positive LCAP and negative TMP sets.
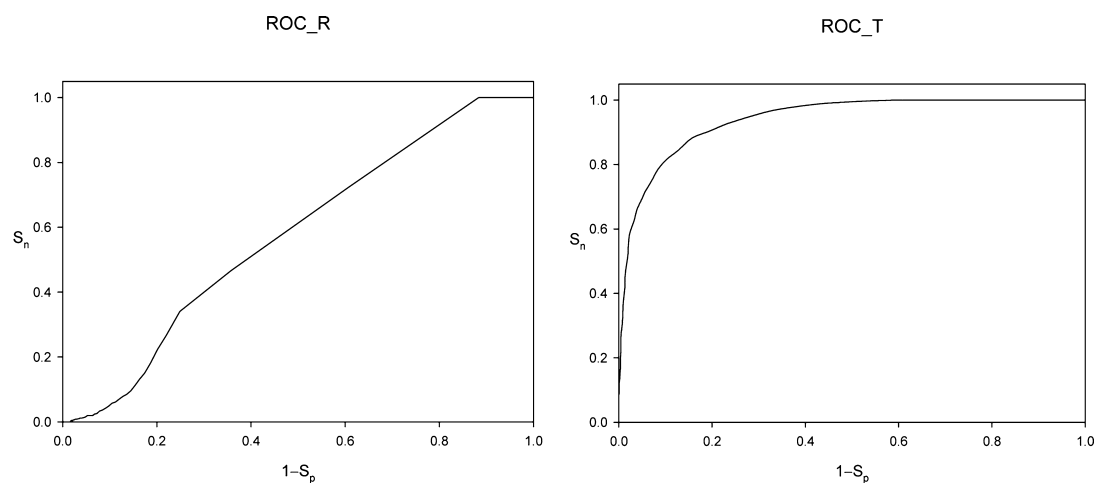
ROC_R  ROC_T



**Figure 4.** ROC curves for evaluation prediction quality of disordering for training sets: ROC_R corresponds to positive LCAP and negative RFP500 sets and ROC_T corresponds to positive LCAP and negative TMP sets.

$(d, \theta)$ plot for the LCAP set is higher than for the other peptide sets (RFP500 and TMP). From Figure 1b, it can be seen that it is energetically more favorable for the membrane proteins to penetrate more deeply into the membrane ($d$ = 2−10 Å). At the same time, peptides from the data set of random protein fragments are located closer to the membrane surface ($d$ = 12−
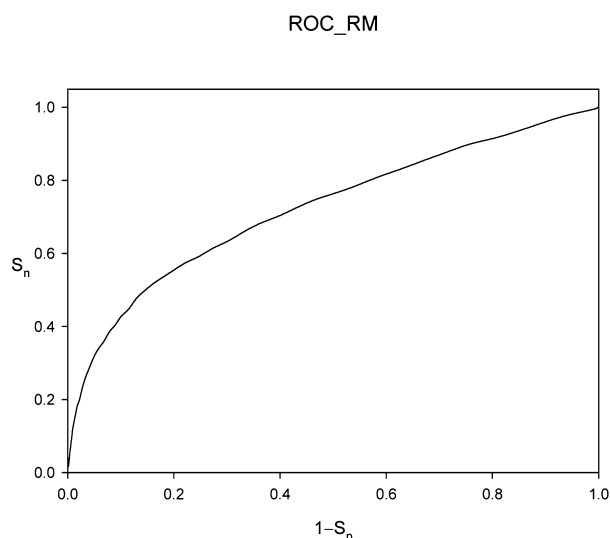
ROC_RM



**Figure 5.** ROC curve for evaluation prediction quality of the transemembrane helixes of linear hydrophobic moment (ROC_RM) for TMP and RFP500 sets.

30 Å; see Figure 1c). Peptides from the later (last) set are distributed on the $(d, \theta)$ plot less densely than from the other peptide sets.

On the basis of these data, we have decided to use the receiver operating characteristic (ROC) curve to quantify differences of LCAP from membrane proteins and soluble proteins fragments (see the Methods section). Calculations have shown that the optimal values of the penetration depth and angle ($d$ and $\theta$) are 12.9 Å and 81° at a fragment length of 17 amino acids. The optimal AUC_R and AUC_T values (see the Methods section) for the two data sets (RFP500 and TMP) at the same time are 0.76 and 0.78, respectively. Therefore, we concluded that the location of the peptide in relation to membrane can be used as a descriptor to distinguish linear cationic antimicrobial peptides from other peptides.

**Evaluation of the Efficiency of LCAP Prediction.** Receiver operating characteristic (ROC) curves were used to evaluate the effectiveness of various characteristics for LCAP prediction. ROC curves, plotted for each characteristic, are shown in Figures 2−10. Quantitative evaluation of effectiveness of the characteristics is made on the basis of the following quantities: (a) area under the ROC curve (defined as AUC_R relative to the RFP500 negative set and defined as AUC_T relative to the TMP negative set) and (b) a threshold for each characteristic by which prediction of peptide antimicrobiality will be done. A threshold is determined by a point on the ROC curve closest to the point (0,1). Sensitivity, specificity and balanced accuracy for the thresholds have been evaluated.

Varying $z_i$ from $\min(z^p_d)$ to $\max(z^p_d)$ and based on the assumption that the values of the descriptors must be higher for LCAP than for non-AMP (as in the case of hydrophobic moment, charge density, linear hydrophobic moment, propensities to aggregation in vitro and in vivo), condition of $z^p_d > z_i$ was used to calculate sensitivity and specificity ($S_{ni}$ and $S_{pi}$) that is $i$th point of the ROC curve. When we assumed that the values of the descriptors must be less for LCAP than for non-AMP (as in the case of hydrophobicity and disordering), the condition $z^p_d < z_i$ was used to calculate sensitivity and specificity ($S_{ni}$ and $S_{pi}$) that is $i$th point of the ROC curve. If the assumption is true that the value of AUC for each descriptor will be higher than 0.5. It can be noted that the higher the value of AUC, the better the descriptor discriminates AMP from non-AMP. The value of AUC for good descriptors must be no less than 0.7. But as we can see, our results show that the values of AUC_R for linear moment and in vitro aggregation are close to 0.5 (see Table 2 and Figures 2 and 3) and for disordered even less than 0.5 (see Table 2 and Figure 4). It means that the last characteristics cannot distinguish antimicrobial from non-antimicrobial peptides. The low value of AUC_R = 0.56 for the linear moment suggests that for the most antimicrobial peptides, there is no significant linear separation of hydrophobic and hydrophilic residues along the peptide chain. On the other hand, the ROC curve plotted for linear moment of TMP set relative to the negative set RFP500 (ROC_RM (Figure 5)) gives the value 0.73 for the area under the ROC curve (AUC_RM = 0.73). These differences between the values of AUC_R and AUC_RM can be explained by the fact that in contrast to antimicrobial peptides, in the transmembrane peptides, linear separation of hydrophobic and hydrophilic group of residues occurs. Such separation was revealed by other authors[41,42] also, who supposed that amphyphilic residues are
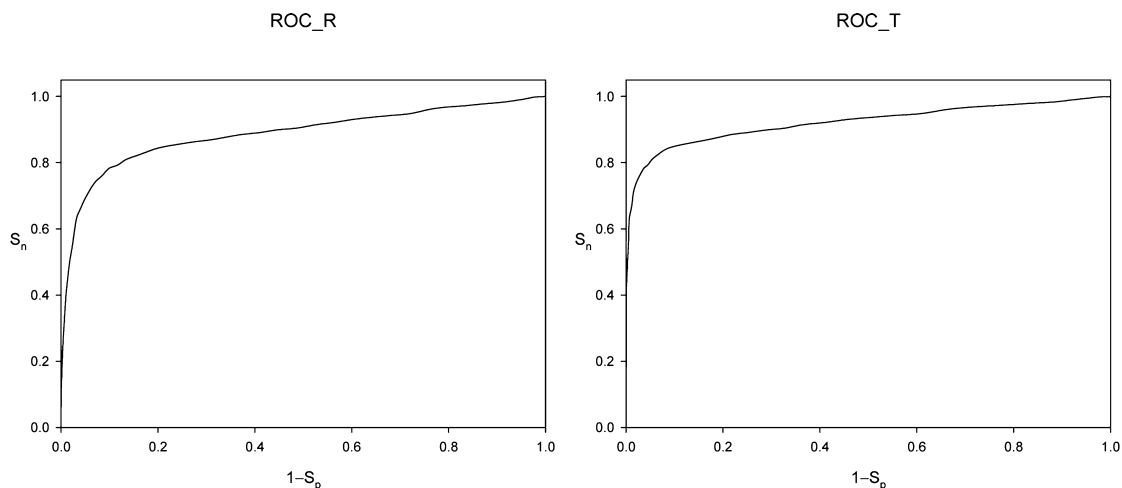
ROC_R



ROC_T



**Figure 6.** ROC curves for evaluation prediction quality of hydrophobic moment for training sets: ROC_R corresponds to positive LCAP and negative RFP500 sets and ROC_T corresponds to positive LCAP and negative TMP sets.
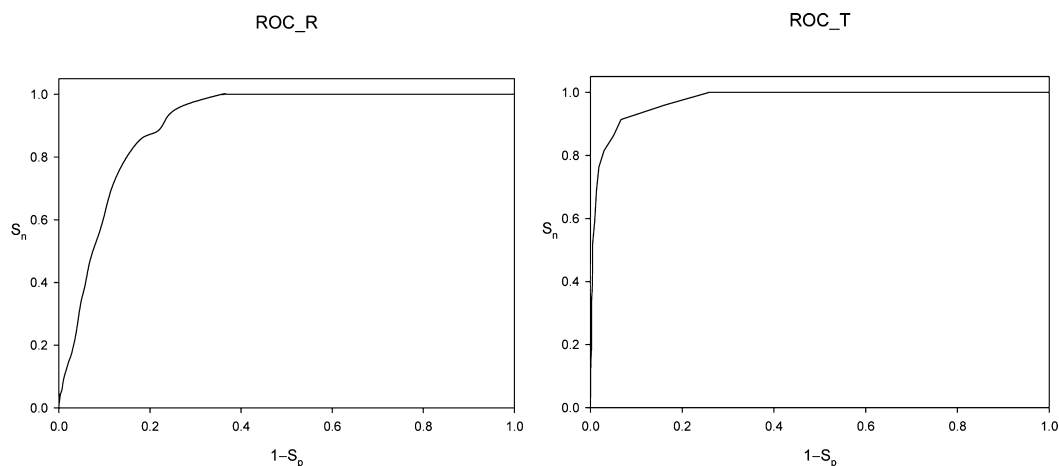
**Figure 7.** ROC curves for the evaluation of prediction quality of charge density for training sets: ROC_R corresponds to positive LCAP and negative RFP500 sets and ROC_T corresponds to positive LCAP and negative TMP sets.
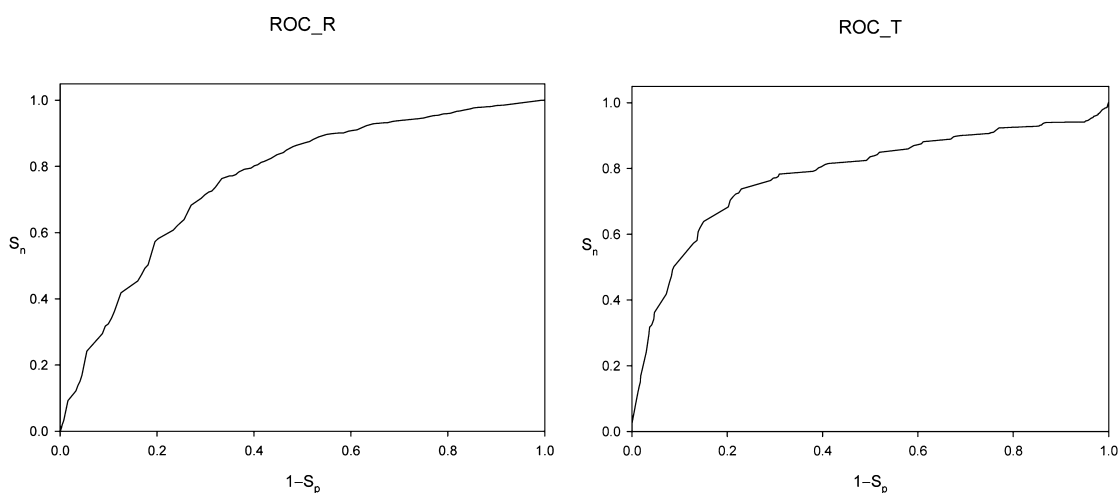


**Figure 8.** ROC curves for evaluation prediction quality of location of the peptide along the membrane (LPM) for training sets: ROC_R corresponds to positive LCAP and negative RFP500 sets and ROC_T corresponds to positive LCAP and negative TMP sets.
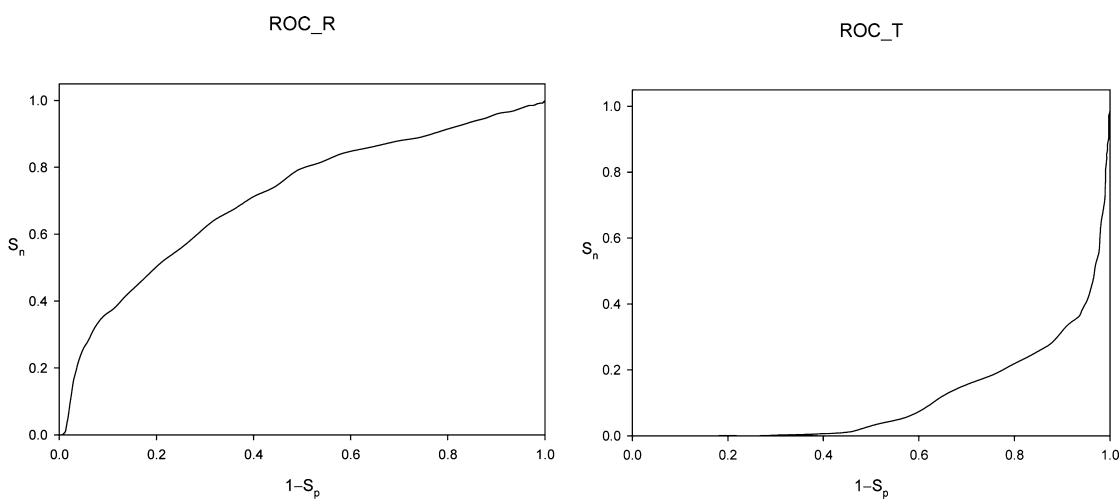


**Figure 9.** ROC curves for evaluation prediction quality of hydrophobicity for training sets: ROC_R corresponds to positive LCAP and negative RFP500 sets and ROC_T corresponds to positive LCAP and negative TMP sets.

concentrated at the ends of the transmembrane helix while hydrophobic residues are located in the middle.

AUC_R value for disordered is 0.47, which is less than 0.5, and it can be said that according to the proposed by Uversky criteria,[37] antimicrobial peptides are more ordering than
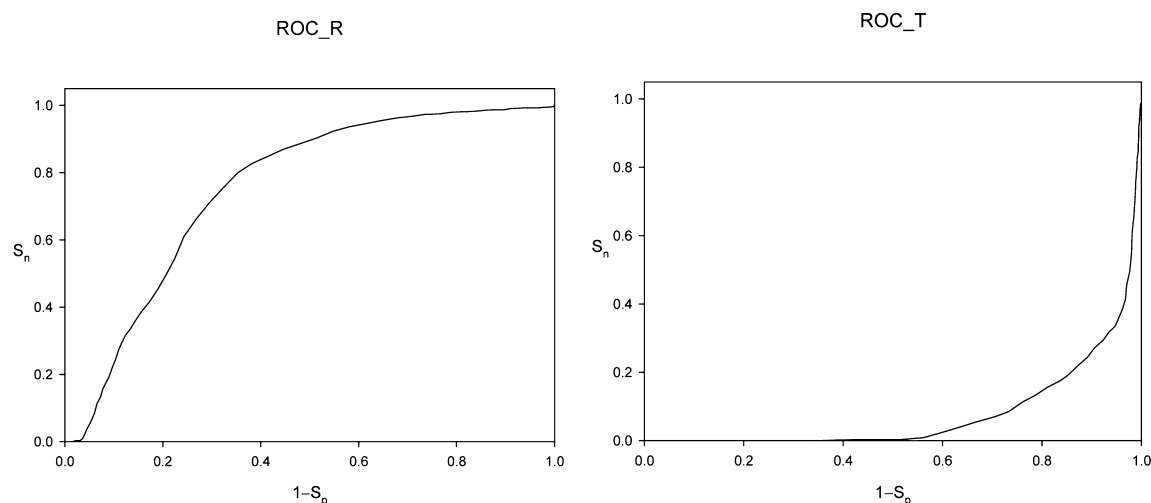
**Figure 10.** ROC curves for evaluation prediction quality of propensity to aggregation in vitro for training sets: ROC_R corresponds to positive LCAP and negative RFP500 sets and ROC_T corresponds to positive LCAP and negative TMP sets.

**Table 2. Comparison of the Different Descriptors for Training Set (LCAP and RFP500)**

|  | AUC_R*100 | AUC_T *100 | $R_{min}$[a]*100 | $S_n$*100 | $S_p$*100 | BAC*100 |
|---|---|---|---|---|---|---|
| hydrophobic moment | 88.63 | 92.01 | 23.32 | 80.79 | 86.77 | 83.78 |
| charge | 90.29 | 97.61 | 22.99 | 86.24 | 81.58 | 83.91 |
| LPM | 76.12 | 78.32 | 37.38 | 76.36 | 68.12 | 72.24 |
| hydrophobicity | 71.18 | 11.09 | 47.96 | 64.27 | 69.15 | 66.14 |
| linear hydrophobic moment | 56.09 | 33.37 | 68.65 | 40.63 | 65.54 | 53.08 |
| disordering | 46.59 | 93.75 | 74.96 | 50.97 | 43.30 | 47.13 |
| in vitro aggregation[b] | 57.41 | 4.34 | 64.23 | 46.63 | 64.26 | 55.45 |
| in vivo aggregation[c] | 75.38 | 7.87 | 40.57 | 75.81 | 67.43 | 71.62 |

[a]Distance from the point (0,1) to the point on the ROC curve closest to the point (0,1). [b]TANGO AGG index. [c]AGGERSCAN Na4vSS index.

**Table 3. Prediction Quality of Combined Use of Three Descriptors for Training Set (LCAP and RFP500)**

|  | AUC_R*100 | $S_n$*100 | $S_p$*100 | BAC*100 |
|---|---|---|---|---|
| hydrophobic moment + charge | 91.23 | 84.21 | 91.69 | 87.95 |
| hydrophobic moment + charge + LPM | 91.38 | 84.03 | 92.50 | 88.26 |

**Table 4. Prediction Quality for Different Methods for Sets LCAP and RFP10**

|  | $S_n$*100 | $S_p$*100 | BAC*100 | AC*100 |
|---|---|---|---|---|
| hydrophobic moment | 80.79 | 89.54 | 85.17 | 88.74 |
| charge | 86.24 | 74.78 | 80.51 | 75.83 |
| hydrophobic moment + charge | 84.21 | 92.36 | 88.23 | 89.61 |
| hydrophobic moment + charge + LPM | 84.03 | 93.00 | 88.52 | 90.20 |
| SVM | 93.44 | 87.57 | 90.51 | 88.11 |
| RF | 95.57 | 86.51 | 91.04 | 87.33 |
| ANN | 90.21 | 86.20 | 88.21 | 86.56 |
| DA | 92.89 | 86.72 | 89.81 | 87.28 |

fragments of random proteins. This may be due to the fact that the Uversky criterion, which determines the degree of the disorder of globular proteins, is not suitable for the evaluation of the disorder of small peptides.

In the case of hydrophobicity, hydrophobic moment, LPM and charge density and propensity to aggregation in vivo AUC_R > 0.7 (see Table 2 and Figures 6−10), which indicate that these characteristics can be used to distinguish

antimicrobial from soluble nonantimicrobial peptides. For the hydrophobic moment and charge density, AUC_T > AUC_R. On the basis of this, we can suggest that if the value of the last characteristics can discriminate a peptide (potential LCAP) from the nonmembrane peptides, it should rather discriminate the peptide from the transmembrane peptides also. So, for these characteristics, only a single threshold defined from ROC_R can be used, because for this threshold, sensitivities are the same for ROC_R and ROC_T, but specificity and thus accuracy obtained from ROC_T is larger than from ROC_R. Consequently, LCAP peptides can be discriminated from the membrane peptides with accuracy obtained from the threshold defined from ROC_R only.

LPM was already optimized in such a way that the greater difference from the RFP500 and TMP sets was reached (see above). AUC_R for this descriptor is 0.76, so it can be used as a LCAP characteristic (see Table 2).

Though, for the hydrophobicity, AUC_R = 0.71, but AUC_T = 0.11 < 0.5 (see Table 2), it means that the LCAP has a lower average hydrophobicity than the transmembrane helices, but greater than random fragments from the soluble proteins. Therefore, we cannot use single threshold to discriminate nonantimicrobial and antimicrobial peptides.

Analogous results were obtained for the propensity to in vivo aggregation (AUC_R = 0.75, but AUC_T = 0.08 < 0.5). The question of AMP aggregation is difficult and unclear. There is speculation that AMP greatly differ in the predisposition to aggregation.[43] Our results confirm this speculation because various AMP peptides from the considered benchmarks great

**Table 5. Prediction Quality for Different Methods for Test Sets[a]**

|  | TPS1 set | | | TPS2 set | | |
|---|---|---|---|---|---|---|
|  | $S_n$*100 | BAC*100 | AC*100 | $S_n$*100 | BAC | AC*100 |
| hydrophobic moment | 80.61 | 85.08 | 88.82 | 86.78 | 88.18 | 89.49 |
| charge | 72.45 | 73.62 | 74.76 | 81.03 | 77.91 | 74.88 |
| hydrophobic moment + charge | 81.63 | 87.00 | 92.27 | 89.66 | 91.01 | 92.32 |
| hydrophobic moment + charge + LPM | 81.63 | 87.32 | 92.93 | 89.66 | 91.33 | 92.98 |
| SVM | 84.69 | 86.13 | 87.55 | 91.95 | 89.76 | 87.64 |
| RF | 81.63 | 84.07 | 86.47 | 93.68 | 90.10 | 86.62 |
| ANN | 83.67 | 84.93 | 86.17 | 89.66 | 87.93 | 86.25 |
| DA | 83.67 | 85.20 | 86.69 | 91.38 | 89.05 | 86.80 |

[a]Specificities for test sets have been calculated for the RFP10 set (see Table 4).

differ by aggregation index, especially for in vitro aggregation. A wide range of proposed mechanisms of AMP action can be explained by the fact that AMP behave differently in terms of the stability of their aggregates both in the membrane and in the aqueous environment. Our results (Table 1) show that propensity to in vitro aggregation does not discriminate AMP from non-AMP and propensity to in vivo aggregation discriminates AMP from non membrane non-AMP but does not do it from transemembrane non-AMP. So, we have not used these descriptors as discriminated LCAP characteristics.

The highest values of AUC_R correspond to hydrophobic moment and charge density. Thus, we can suggest that these characteristics are the best separators between nonantimicrobial and antimicrobial peptides.

So, three descriptors: hydrophobic moment, charge density and LPM were selected as the most effective LCAP descriptors. Given the above, it can be assumed that the combined use of these three characteristics can improve the prediction of LCAP. To combine these characteristics, we have taken into account the fact that for the separation of LCAP and non-AMP, specific set of threshold values that can be obtained from the analysis of ROC curves were used. Accordingly, by changing synchronously thresholds for different characteristics, we can simply optimize these thresholds to obtain the greatest accuracy. The corresponding balanced accuracy for the charge density alone, hydrophobic moment alone, hydrophobic moment and charge density together and for the three characteristics together in case of the training set (LCAP and RFP500) are 83.91, 83.78, 87.95 and 88.26, respectively (Tables 2 and 3) .

We have also tried to evaluate our results with other prediction methods. As we have mentioned above, several computational methods[4,5,11,18−23] have been proposed for the predicting AMPs. However, some methods[4,5,18] did not contain available web services for testing our data sets. BACTIBASE[19,20] and PhytAMP[21] methods were specifically designed for bacteriocin and plant, respectively. As for AntiBP[22] and AntiBP2 methods,[23] they were designed to identify the AMPs in a protein sequence, and hence could not be used to compare with our method. So, to make the comparison meaningful, our method was compared with CAMP method,[11] which was developed based on the random forests (RF), SVM, ANN and discriminant analysis (DA). This method can be used for the evaluation of the sensitivity, specificity and accuracy for the considered training positive set. As a set of nonantimicrobial peptides (negative set), we have used a set of 10 peptide fragments (instead of 500) for each peptide in the AMP set (RFP10). The corresponding balanced accuracy for three considered characteristics together when using this set is 88.52. For the same positive and negative sets, CAMP method gives

the following values for the balanced accuracy: random forests (RF), 91.04; SVM, 90.51; discriminant analysis (DA), 89.81; ANN, 88.21 (Table 4).

For testing purposes, two independent positive sets (TPS1 and TPS2) (see the Methods section) were used. The results of comparison for test sets are shown in Table 5.

We cannot use any additional independent test sets for non-AMP (negative set), so the RFP10 set, which was not employed for training purposes, was used as a negative test set. As we can see from Tables 4 and 5 the best prediction quality for the training and both test sets was obtained when all three descriptors were used together, although a pair of the descriptors (hydrophobic moment and charge density) gives very close results. We have also noted that for the test sets, the prediction quality (balanced accuracy) based on the hydrophobic moment and charge density gives better results than the one obtained from the all CAMP prediction algorithms.

We want to emphasize the fact that the CAMP method uses a combination of numerous characteristics and more complicated, refined and effective discriminative methods.[11] High performance of our approach can be explained by the fact that we have used only one class of AMP, cationic linier peptides. Our results confirm the assumption that prediction of AMP is preferable to make for the peculiar class separately, using a particular approach in each case.

The AMP prediction tool based on the considered method is included into the Database of Antimicrobial Activity and Structure of Peptides (DBAASP) and available at http://www.biomedicine.org.ge/dbaasp/.

## ■ AUTHOR INFORMATION

### Corresponding Authors
*B. Vishnepolsky. Phone: +995 32 2371019. E-mail: b.vishnepolsky@lifescience.org.ge.
*M. Pirtskhalava. Phone: +995 574162397. E-mail: m.pirtskhalava@lifescience.org.ge.

### Notes
The authors declare no competing financial interest.

## ■ REFERENCES

(1) Sang, Y.; Blecha, F. Antimicrobial peptides and bacteriocins: alternatives to traditional antibiotics. *Anim. Health Res. Rev.* **2008**, *9*, 227−235.

(2) Hancock, R. E. W.; Rozek, A. Role of membranes in the activities of antimicrobial cationic peptides. *FEMS Microbiol. Lett.* **2002**, *206* (2), 143−149.

(3) Fjell, C. D.; Hiss, J. A.; Hancock, R. E.; Schneider, G. Designing antimicrobial peptides: form follows function. *Nat. Rev. Drug Discovery* **2011**, *11* (1), 37−51.

(4) Torrent, M.; Andreu, D.; Nogués, V. M.; Boix, E. Connecting Peptide Physicochemical and Antimicrobial Properties by a Rational Prediction Model. *PLoS One* **2011**, *6* (2), e16968.

(5) Wang, P.; Hu, L.; Liu, G.; Jiang, N.; Chen, X.; Xu, J.; Zheng, W.; Li, L.; Tan, M.; Chen, Z.; Song, H.; Cai, Y.; Chouet, K. Prediction of Antimicrobial Peptides Based on Sequence Alignment and Feature Selection Methods. *PLoS One* **2011**, *6* (4), e18476.

(6) Frecer, V.; Ho, B.; Ding, J. L. De novo design of potent antimicrobial peptides. *Antimicrob. Agents Chemother.* **2004**, *48*, 3349−3357.

(7) Jenssen, H.; Lejon, T.; Hilpert, K.; Fjell, C. D.; Cherkasov, A.; Hancock, R. E. W. Evaluating different descriptors for model design of antimicrobial peptides with enhanced activity toward *P. aeruginosa*. *Chem. Biol. Drug Des.* **2007**, *70*, 134−142.

(8) Jenssen, H.; Fjell, C. D.; Cherkasov, A.; Hancock, R. E. QSAR modeling and computer-aided design of antimicrobial peptides. *J. Pept. Sci.* **2008**, *14*, 110−114.

(9) Frecer, V. QSAR analysis of antimicrobial and haemolytic effects of cyclic cationic antimicrobial peptides derived from protegrin-1. *Bioorg. Med. Chem.* **2006**, *14*, 6065−6074.

(10) Cherkasov, A.; Jankovic, B. Application of 'inductive' QSAR descriptors for quantification of antibacterial activity of cationic polypeptides. *Molecules* **2004**, *9*, 1034−1052.

(11) Waghu, F. H.; Gopi, L.; Barai, R. S.; Ramteke, P.; Nizami, B.; Idicula-Thomas, S. CAMP: Collection of sequences and structures of antimicrobial peptides. *Nucleic Acids Res.* **2014**, *42*, D1154−D1158.

(12) Taboureau, O.; Olsen, O. H.; Nielsen, J. D.; Raventos, D.; Mygind, P. H.; Kristensen, H. H. Design of novispirin antimicrobial peptides by quantitative structure-activity relationship. *Chem. Biol. Drug Des.* **2006**, *68*, 48−57.

(13) White, S. H.; Wimley, W. C. Membrane protein folding and stability: physical principles. *Annu. Rev. Biophys. Biomol. Struct.* **1999**, *28*, 319−365.

(14) Jayasinghe, S.; Hristova, K.; White, S. H. Energetics, stability, and prediction of transmembrane helices. *J. Mol. Biol.* **2001**, *312*, 927−934.

(15) Scocchi, M.; Tossi, A.; Gennaro, R. Proline-rich antimicrobial peptides: converging to a non-lytic mechanism of action. *Cell. Mol. Life Sci.* **2011**, *68*, 2317−30.

(16) Kragol, G.; Hoffmann, R.; Chatergoon, M. A.; Lovas, S.; Cudic, M.; Bulet, P.; Condie, B. A.; Rosengren, K. J.; Montaner, L. J.; Otvos, L. Identification of crucial residues for the antibacterial activity of the proline-rich peptide, pyrrhocoricin. *Eur. J. Biochem.* **2002**, *269*, 4226−4237.

(17) Wang, Z.; Wang, G. APD: the Antimicrobial Peptide Database. *Nucleic Acids Res.* **2004**, *32*, D590−D592.

(18) Fjell, C. D.; Hancock, R. E.; Cherkasov, A. AMPer: a database and an automated discovery tool for antimicrobial peptides. *Bioinformatics* **2007**, *23*, 1148−1155.

(19) Hammami, R.; Zouhir, A.; Ben Hamida, J.; Fliss, I. BACTIBASE: a new web-accessible database for bacteriocin characterization. *BMC Microbiol.* **2007**, *7*, 89.

(20) Hammami, R.; Zouhir, A.; Le Lay, C.; Ben Hamida, J.; Fliss, I. BACTIBASE second release: a database and tool platform for bacteriocin characterization. *BMC Microbiol.* **2010**, *10*, 22.

(21) Hammami, R.; Ben Hamida, J.; Vergoten, G.; Fliss, I. PhytAMP: a database dedicated to antimicrobial plant peptides. *Nucleic Acids Res.* **2009**, *37*, D963−968.

(22) Lata, S.; Sharma, B. K.; Raghava, G. P. Analysis and prediction of antibacterial peptides. *BMC Bioinf.* **2007**, *8*, 263 (2007).

(23) Lata, S.; Mishra, N. K.; Raghava, G. P. AntiBP2: improved version of antibacterial peptide prediction. *BMC Bioinf.* **2010**, *11* (Suppl1), S19; *Bioinformatics* **2005**, *59*, 252−265.

(24) Tusnády, G. E.; Dosztányi, Z.; Simon, I. PDB_TM: selection and membrane localization of transmembrane proteins in the protein data bank. *Bioinformatics* **2004**, *20* (17), 2964−2972.

(25) Kozma, D.; Simon, I.; Tusnády, G. E. PDBTM: Protein Data Bank of transmembrane proteins after 8 years. *Nucleic Acids Res.* **2013**, *41* (D1), D524−D529.

(26) Yeaman, M. R.; Yount, N. Y. Unifying themes in host defence effector polypeptides. *Nat. Rev. Microbiol.* **2007**, *5*, 727−740.

(27) Kyte, J.; Doolittle, R. F. A simple method for displaying the hydropathic character of a protein. *J. Mol. Biol.* **1982**, *157* (1), 105−132.

(28) Wimley, W. C.; White, S. H. Experimentally determined hydrophobicity scale for proteins at membrane interfaces. *Nat. Struct. Biol.* **1996**, *3*, 842−848.

(29) Koehler, J.; Woetzel, N.; Staritzbichler, R.; Sanders, C. R.; Meiler, J. A unified hydrophobicity scale for multispan membrane proteins. *Proteins* **2009**, *76*, 13−29.

(30) Hessa, T.; Meindl-Beinker, N. M.; Bernsel, A.; Kim, H.; Sato, Y.; Lerch-Bader, M.; Nilsson, I.; White, S. H.; von Heijne, G. Molecular code for transmembrane-helix recognition by the Sec61 translocon. *Nature* **2007**, *450*, U1026−U1032.

(31) Eisenberg, D.; Weiss, R. M.; Terwilliger, T. C. The hydrophobic moment detects periodicity in protein hydrophobicity. *Proc. Natl. Acad. Sci. U. S. A.* **1984**, *81* (1), 140−144.

(32) Moon, C. P.; Fleming, K. G. Side-chain hydrophobicity scale derived from transmembrane protein folding into lipid bilayers. *Proc. Natl. Acad. Sci. U. S. A.* **2011**, *108* (25), 10174−10177.

(33) Yeaman, M. R.; Yount, N. Y. Mechanisms of antimicrobial peptide action and resistance. *Pharmacol. Rev.* **2003**, *55* (1), 27−55.

(34) Eisenberg, D.; Weiss, R. M.; Terwilliger, T. C. The helical hydrophobic moment: a measure of the amphiphilicity of a helix. *Nature* **1982**, *299* (5881), 371−374.

(35) Senes, A.; Chadi, D. C.; Law, P. B.; Walters, R. F. S.; Nanda, V.; DeGrado, W. F. Ez, a Depth-dependent Potential for Assessing the Energies of Insertion of Amino Acid Side-chains into Membranes: Derivation and Applications to Determining the Orientation of Transmembrane and Interfacial Helices. *J. Mol. Biol.* **2007**, *366*, 436−448.

(36) Shai, Y. Mode of action of membrane active antimicrobial peptides. *Biopolymers* **2002**, *66*, 236−248.

(37) Uversky, V.; Gillespie, J.; Fink, A. Why are "natively unfolded" proteins unstructured under physiological conditions? *Proteins: Struct. Funct. Gen.* **2000**, *41* (3), 415−427.

(38) Fernandez-Escamilla, A. M.; Rousseau, F.; Schymkowitz, J.; Serrano, L. Prediction of sequence-dependent and mutational effects on the aggregation of peptides and proteins. *Nat. Biotechnol.* **2004**, *22*, 1302−1306.

(39) Conchillo-Sole, O.; de Groot, N. S.; Aviles, F. X.; Vendrell, J.; Daura, X.; Ventura, S. AGGRESCAN: a server for the prediction and evaluation of "hot spots" of aggregation in polypeptides. *BMC Bioinf.* **2007**, *8*, 65.

(40) Bechinger, B. The structure, dynamics and orientation of antimicrobial peptides in membranes by multidimensional solid-state NMR spectroscopy. *Biochim. Biophys. Acta* **1999**, *1462* (1−2), 157−83.

(41) Mitaku, S.; Hirokava, T.; Tsuji, T. Amphiphilicity index of polar amino acids as an aid in the characterization of amino acid preference at membrane-water interfaces. *Bioinformatics* **2002**, *18*, 608−16.

(42) Ulmschneider, M. B.; Sansom, M. S. P.; Di Nola, A. Properties of Integral Membrane Protein Structures:Derivation of an Implicit Membrane Potential. *Proteins: Struct., Funct., Bioinf.* **2005**, *59*, 252−262.

(43) Wimley, W. C. Describing the Mechanism of Antimicrobial Peptide Action with the Interfacial Activity Model. *ACS Chem. Biol.* **2010**, *5*, 905−917.

1523

dx.doi.org/10.1021/ci4007003 | *J. Chem. Inf. Model.* 2014, 54, 1512−1523