



OPEN

ACP-MHCNN: an accurate multi-headed deep-convolutional neural network to predict anticancer peptides

Sajid Ahmed^{1,6}, Rafsanjani Muhammod^{1,6}, Zahid Hossain Khan^{1,6}, Sheikh Adilina¹, Alok Sharma^{2,3}, Swakkar Shatabda^{1✉} & Abdollah Dehzangi^{4,5✉}

Although advancing the therapeutic alternatives for treating deadly cancers has gained much attention globally, still the primary methods such as chemotherapy have significant downsides and low specificity. Most recently, Anticancer peptides (ACPs) have emerged as a potential alternative to therapeutic alternatives with much fewer negative side-effects. However, the identification of ACPs through wet-lab experiments is expensive and time-consuming. Hence, computational methods have emerged as viable alternatives. During the past few years, several computational ACP identification techniques using hand-engineered features have been proposed to solve this problem. In this study, we propose a new multi-headed deep convolutional neural network model called ACP-MHCNN, for extracting and combining discriminative features from different information sources in an interactive way. Our model extracts sequence, physicochemical, and evolutionary based features for ACP identification using different numerical peptide representations while restraining parameter overhead. It is evident through rigorous experiments using cross-validation and independent-dataset that ACP-MHCNN outperforms other models for anticancer peptide identification by a substantial margin on our employed benchmarks. ACP-MHCNN outperforms state-of-the-art model by 6.3%, 8.6%, 3.7%, 4.0%, and 0.20 in terms of accuracy, sensitivity, specificity, precision, and MCC respectively. ACP-MHCNN and its relevant codes and datasets are publicly available at: <https://github.com/mrzResearchArena/Anticancer-Peptides-CNN>. ACP-MHCNN is also publicly available as an online predictor at: <https://anticancer.pythonanywhere.com/>.

Cancer is one of the deadliest diseases in the world. Even though there are several ways of treating some of the cancer types, still there is no certain treatment for most of the cancers. Two of the major treatment strategies for cancer are radiation therapy and chemotherapy^{1,2}. However, they are both expensive and have long term negative side effects¹. In addition, cancer cells can become resistant to the chemotherapeutic drugs¹. Therefore, there is a demand for finding new low cost and more effective treatments for cancer³. Among the newly introduced treatment methods for this deadly disease, anticancer peptides (ACP) have gained a lot of attention in the recent years as a less toxic and potentially more effective treatment for cancer^{3,4}.

ACPs are short peptides consisting of 10 to 50 amino acids which are typically derived from antimicrobial peptides⁵. ACPs perform a wide range of cytotoxic activities against cancer cells while leave benign cells intact which is the reason behind their high specificity and low side effects⁶. Additionally, ACPs have low production cost, they are easy to synthesize and modify, and they have excellent tumour penetration capabilities⁷. In the past few years, many ACP based treatment options have been tested on a wide variety of cancer cells. However, only a few of them have been cleared for further clinical trials^{8,9}. Hence, rapid identification of potential ACPs is important for cancer therapeutic advancement. However, identification of these peptides through wet-lab

¹Department of Computer Science and Engineering, United International University, Dhaka, Bangladesh. ²Laboratory for Medical Science Mathematics, RIKEN Center for Integrative Medical Sciences, Yokohama 230-0045, Japan. ³Institute for Integrated and Intelligent Systems, Griffith University, Brisbane, QLD 4111, Australia. ⁴Department of Computer Science, Rutgers University, Camden, NJ 08102, USA. ⁵Center for Computational and Integrative Biology, Rutgers University, Camden, NJ 08102, USA. ⁶These authors contributed equally: Sajid Ahmed, Rafsanjani Muhammod and Zahid Hossain Khan. ✉email: swakkar@cse.uui.ac.bd, i.dehzangi@rutgers.edu

experiments is relatively costly and time consuming¹. Therefore, there is a demand for fast and accurate computational methods to tackle this problem. Among different computational methods, machine learning has merged as a promising approach to identify ACPs efficiently and effectively.

During the past few years, a wide range of traditional Machine Learning (ML) methods have been proposed to identify ACPs. These traditional ML techniques require a set of hand-engineered features to represent protein sequences for the classification purpose. Thus, various methods for extracting effective features to represent proteins and peptides in an effective manner that contain significant discriminatory information for the classification purpose have been proposed. AntiCP was the first ML model for ACP identification that was proposed in¹. In this model, peptide sequences are formulated by amino acid composition (AAC), split AAC (using N-terminal and C-terminal residues), dipeptide composition (DPC), and binary profiles features (BPF)¹. Afterwards, these features are passed as input to a Support Vector Machine (SVM) classifier for separating the ACPs from the non-ACPs.

Shortly after that, Hajisharifi et al., proposed two methods for ACP identification using SVM¹⁰. In the first method, SVM was employed for separating ACPs from non-ACPs. They used pseudo-amino acid composition (PseAAC) method on different combinations of 6 physicochemical properties of the amino acids to extract features. In the second method, the binary classification was performed using SVM with a local alignment-based kernel method designed for feature extraction from peptide sequence¹⁰. Later on, Chen et al. proposed iACP, where gapped dipeptide compositions (g-gap DPC) were used for feature extraction from peptide sequences, and SVM with radial basis function (RBF) kernel was used for the classification purpose³.

More recently, Manavalan et al., proposed MLACP to tackle this problem. To build this model, AAC, DPC, atomic composition (ATC) of the sequences, and physicochemical properties of the residues were used for feature extraction while, SVM and Random Forest (RF) classifiers were used for ACP identification¹¹. At the same time, Akbar et al., proposed iACP-GAEns, which used g-gap DPC, reduced amino acid alphabet composition (RAAAC), and PseAAC based on hydrophobicity and hydrophilicity of the amino acids (Am-PseAAC) for feature extraction. They also proposed an ensemble of different classifiers that combined SVM, RF, Probabilistic Neural Network (PNN), Generalized Regression Neural Network (GRNN), and K-nearest Neighbour (KNN) classification models for ACP identification¹².

Later on, Xu et al., proposed a hybrid sequence-based model, where the peptides were converted to feature vectors through g-gap DPC to tackle this problem. They also used SVM and RF as their employed classifiers¹³. At the same time Kabir et al., proposed TargetACP, where the peptides were represented using split AAC, correlation factors extracted from PSSM profiles (PsePSSM), and composite protein sequence representation (CPSR). They also used SVM, RF and KNN classifiers as their employed models¹⁴.

Most recently, Schaduangrat et al. proposed ACPred, where different combinations of AAC, DPC, PseAAC, Am-PseAAC, and physicochemical properties were used for peptide representation. They also used SVM and RF classifiers for the ACP identification prediction⁴. At the same time, Wei et al., proposed ACPred-FL, where AAC, g-gap DPC, BPF, amino acid-specific physicochemical property-based bit vectors and composition-transition-distribution (CTD) methods were used for feature extraction. Similarly, they used SVM based ensemble model as their employed classifier¹⁵.

During the revision stage of this manuscript, Charoenkwan et al. proposed a sequence-based method iACP-FSCM with an emphasis on model interpretability, where 11 local and global amino acid composition-based features were utilized with a weighted-sum-based prediction mechanism¹⁶. Furthermore, Agrawal et al. proposed a sequence-based method AntiCP 2.0 along with two ACP identification datasets¹⁷. AntiCP 2.0 has been shown to outperform all the existing ACP identification methods with state-of-the-art accuracy. In a recent review article, Basith et al.¹⁸ (Sir, please fix the citation order) presented a concise summary of 16 ML methods developed so far for ACP identification.

Using traditional ML models (SVM, RF, KNN, etc.), the systems' performances depend on the underlying manual feature extraction mechanisms. However, formulating problem-specific optimal feature representation for these sequences is not a trivial task and requires significant iterations of trial and error. In recent years, deep learning (DL) methods attracted tremendous attention to tackle challenging problems related to biological sequences because in many cases, unlike traditional ML algorithms, they do not require manual feature extraction to represent the input data^{15–25}. Several DL methods, such as Convolutional Neural Network (CNN)^{20,26}, Recurrent Neural Network (RNN)²⁰, word embedding^{27,28}, and autoencoder^{29–31} have been successfully employed for feature extraction and classification for DNA, RNA, and protein sequences. Methods such as CNN and RNN exploit spatial locality and ordering information of the residues for ensuring that the extracted features retain a significant amount of discriminatory information from biological sequences.

However, none of the studies related to ML-based ACP identification explored automated feature extraction using DL methods until recently, when ACP-DL was proposed in³². Although Timmons et al. proposed a deep neural network architecture ENNAACT for ACP identification³³, it still operates on manually extracted features (AAC, DPC, g-gap DPC among others). To the best-of-our-knowledge ACP-DL is the only DL-based automated feature extraction method proposed for this problem, so far. ACP-DL uses bidirectional long-short-term-memory (LSTM) recurrent layers for extracting features from peptide sequences followed by a fully-connected layer with a sigmoid neuron for classification. ACP-DL extracts features from two one-hot vector-based peptide representation techniques (binary profile and k-mer sparse matrix) that only depict the presence of a specific amino acid or a group of amino acids along different positions of the sequences. As a result, physicochemical properties or evolutionary substitution information of the residues, which contain significant information regarding anti-cancer activities of peptide sequences are not utilized in ACP-DL's feature representation process^{4,12,14,15}. As a result, although the predictive performance of ACP-DL is quite impressive, there is still room for improvement.

Although recurrent layers are reliable for converting biological sequences into fixed-size features vectors²⁰, convolutional layers have also demonstrated promising performance addressing similar problems. In fact, CNN

have been demonstrated as an effective technique for feature extraction and classification for DNA, RNA, peptides, and protein sequences in a wide range of studies^{33–41}. However, CNN has never been used for ACP classification task.

In this study, we hypothesize that a new representation technique that depict the residues' evolutionary relationship and their physicochemical characteristics can embellish the feature extraction process for ACP identification since this type of information contains signals necessary for elucidating the structure and function of peptides. With this viewpoint, we are proposing a method called ACP-MHCNN, which consists of three jointly trained groups of stacked CNNs for interactive feature extraction from three distinct information sources for ACP identification. Our results demonstrate that ACP-MHCNN outperforms the current state-of-the-art methods on several well-established ACP identification datasets with a substantial margin. On ACP-500/ACP-164 benchmark dataset, ACP-MHCNN outperforms ACP-DL by 6.3%, 8.6%, 3.7%, 4.0%, and 0.20 in terms of accuracy, sensitivity, specificity, precision, and Matthews correlation coefficient (MCC), respectively. Our model and all its relevant codes and datasets are publicly available at: <https://github.com/mrzResearchArena/Anticancer-Peptides-CNN>, ACP-MHCNN is also publicly available as an online predictor at: <https://anticancer.pythonanywhere.com>.

Materials and methods

In this section, we represent the benchmarks that are used in this study. We also present our sequence representation as well as the proposed feature extraction and classification models.

Benchmark datasets. In this study, we use three independent benchmarks to study the effectiveness and generality of our proposed method. These benchmarks are namely, ACP-740, ACP-240, and the combination of ACP-500 and ACP-164.

ACP-740 dataset was introduced in³². For constructing ACP-740, initially, 388 experimentally validated ACPs (positive samples) were collected, among which 138 were from³ and 250 were from²⁹. Correspondingly, 456 antimicrobial peptides (AMP) without anticancer activity (negative samples) were initially collected, among which 206 were from³ and 250 were from²⁹. Subsequently, using CD-HIT, 12 positive samples and 92 negative samples were removed to ensure that none of the peptide sequence pairs have more than 90% similarity as it was done in previous studies³², which resulted in a dataset with 740 samples (376 positives + 364 negatives). The ACP-240 dataset, which was also introduced in³², consists of 240 samples where 129 experimentally validated ACPs are the positive samples, and 111 AMPs without anticancer activity are the negative samples. To avoid performance over-estimation due to homology bias, using the same procedure as ACP 740, redundancy reduction was performed with a 90% threshold to construct ACP-240.

On the other hand, ACP-500 and ACP-164, were constructed in¹⁵, where ACP-500 is used for training and validation, while ACP-164 is used as an independent test dataset. For constructing these two datasets, initially, 3212 positive samples were collected, among which 138 were from³, 225 were from¹, and 2849 were from⁴². The initial 2250 negative samples were collected from¹. After performing redundancy reduction using CD-HIT with a 90% similarity threshold, 332 positive samples and 1023 negative samples remained. From these remaining non-redundant sequences, 250 positive samples and 250 negative samples were randomly selected for constructing ACP-500, whereas ACP-164 contains the remaining 82 positive samples along with 82 randomly selected negative samples.

Numerical representation for peptide sequences. Although ACP-MHCNN does not require manual feature extraction, it is crucial to encode the sequences in numerical formats since the initial feature extraction layer of any DL architecture performs mathematical operations on the input for extracting class-discriminative activations. Such information is then passed as input to nodes in the subsequent layers. In this study, we exploit three peptide representation methods that are described in the following three sections. Since it has been shown in^{15,32} that considering k amino acids from the N-terminus of a peptide is sufficient for capturing its anticancer activity, we have represented each sequence using its k N-terminus residues. In our experiments, we have set $k = 15$. For sequences having length less than 15, post-padding has been applied as it is explained in details in⁴³.

Binary profile feature (BPF) representation. In our first representation method, each of the 20 amino acids (A, R, N, D, C, Q, E, G, H, I, L, K, M, F, P, S, T, W, Y, and V) is represented using a binary one-hot vector of length 20. For example, A is represented as [1, 0, ..., 0], R is represented as [0, 1, ..., 0], V is represented as [0, 0, ..., 1], and so on. This representation encodes each sequence into a $k \times 20$ matrix. Manually extracted short-range sequence patterns such as AAC, DPC, split AAC and long-range sequence patterns such as g-gap DPC have been successfully employed with traditional ML models for ACP identification^{1,3,10–15}. We hypothesize that our model's feature detection mechanism can capture both short-range and long-range sequence patterns that distinguish the peptides with anticancer activity from BPF representation.

Physicochemical-based (AAIs) representation. Basak et al., used a numerical representation for proteins for identifying length 5 conserved peptides through molecular evolutionary analysis⁴⁴. The underlying numerical representation method proposed in⁴⁵ utilized an alphabet reduction strategy where the amino acids are divided into non-overlapping groups based on their side chain chemical property. The findings from these two studies have implied that amino acid physicochemical properties can facilitate the identification of evolutionarily conserved motifs, which are in turn important for maintaining the appropriate structure or function of the molecules. When these conserved motifs go through changes in the primary structure level, the amino acid residues are usually replaced with the ones with similar physicochemical properties. This phenomenon highlights the signifi-

cant impact of exploring physicochemical properties for motif identification with respect to similarity among the substitute amino acids. Since our model identifies peptides with specific functions, discovering these motifs can strengthen our model.

Moreover, hand-engineered features based on amino acid physicochemical properties have been shown to improve ACP identification in a series of studies that have used traditional machine learning models^{4,10–12,15}. We hypothesize that our feature extraction mechanism can identify similar features from a peptide representation based on the amino acids' physicochemical properties. With these assumptions, our physicochemical property-based representation replaces each of the residues in a peptide sequence with a 31-dimensional vector (composed of 0/1 elements) that depict various physicochemical properties. As a result, each of the sequences is encoded into a $k \times 31$ matrix.

For each amino acid, a unique 31-dimensional vector is formed through the concatenation of a 10-bit vector and a 21-bit vector. Elements of the 10-bit vector depict the membership of a specific amino acid in 10 overlapping groups based on its physicochemical properties as it was explained in¹⁵. Elements of the 21-bit vector are determined based on membership of a specific amino acid in the $7 \times 3 = 21$ groups formed by dividing them into 3 groups based on 7 physicochemical properties namely, polarity, normalized Van der Waals volume, hydrophobicity, secondary structures, solvent accessibility, charge, and polarizability as it was done in¹⁵.

Evolutionary information-based (BLO62) representation. BLOSUM is a symmetric 20×20 matrix constructed by Henikoff et al., in⁴⁶, where each entry is proportional to the probability of substitution of a given amino acid with another amino acid in a protein (substitution probability in evolutionarily related proteins). Each entry in this matrix can be represented using the following equation:

$$M(i, j) = \frac{1}{\lambda} \log \frac{p_{ij}}{f_i f_j} \quad (1)$$

where, p_{ij} is the probability of amino acids 'i' and 'j' being aligned in homologous sequence alignments, f_i is the probability that amino acid 'i' appears in any protein sequence, f_j is the probability that amino acid 'j' appears in any protein sequence, and λ is the scaling factor for rounding off the entries in the matrix to convenient integer values.

The observed substitution frequency for every possible amino acid pair (including identity pairs) is calculated from a large number of trusted pairwise alignments of homologous sequences as it is explained in⁴⁶. If an entry $M(i, j)$ is positive, the number of observed substitutions between amino acids i and j is more than random expectation. Thus, these substitutions are conservative (these substitutions occur more frequently than other random substitutions in homologous sequences). Therefore, each of the 20 rows of this matrix is a vector containing 20 elements that depict a specific amino acid's evolutionary relationship with other amino acids. Here, we use BLOSUM matrix for retrieving a 20-dimensional vector for each of the 20 amino acids and use these vectors for encoding each peptide sequence into a $k \times 20$ matrix. We hypothesize that our feature extraction architecture can automatically extract discriminative evolutionary features for ACP identification from this sequence representation. Among different BLOSUM matrix variations, we have used BLOSUM62 as the most popular one in this study.

Multi-headed convolutional neural network architecture. CNN is a specialized neural network where each neuron in a given layer is connected to a group of neighbouring nodes in the previous layer. These layers drastically reduce parameter overhead and extract translation-invariant meaningful features by exploiting spatial locality structure in data through local connectivity and weight sharing⁴⁷. A convolutional layer usually consists of several kernels where each kernel detects some specific local pattern in different input locations⁴⁷. Since hand-engineered feature extraction methods such as AAC, DPC, g-gap DPC, PseAAC, and PsePSSM utilize ordering of neighbouring residues and their correlation information with respect to evolutionary and physicochemical properties for feature generation from peptide sequences, using convolutional kernels for automatically approximating similar features is a rational choice. Moreover, well-defined ordering among the residues in peptide primary structure, the residues' inherent local neighbourhood structures, and the presence of similar patterns (sequence motifs) at different locations across a peptide make these sequences perfect candidates for feature extraction through convolutional kernels.

The feature extraction mechanism in our proposed model consists of groups of stacked convolutional layers. Each convolutional layer group extracts features from a different representation of the peptide sequence. Since we have use three representation methods that serve as sources of discriminative information, our model contains three parallel layer groups. Each of these groups extract short-range and long-range patterns from a unique sequence representation using two stacked convolutional layers with varying number of kernels. The number of kernels in the layers and size of these filters are hyperparameters tuned through cross-validation⁴⁸.

The output feature maps of the second convolutional layer of each of the three groups are flattened, and the three resulting vectors are concatenated. The unified vector from this concatenation is passed through a dense layer with ReLU (Rectified Linear Unit) activation function for recombining the features extracted from different sequence representations⁴⁹. It is to be mentioned that each element of the input vector for this dense recombination layer is calculated from a single information source (BPF or physicochemical or evolutionary representation) during forward-propagation. In contrast, elements of this layer's output vector can be aggregated from multiple information sources. Hence, this layer enables seamless interaction between different convolutional groups that extract patterns from different representations and facilitates joint feature learning from multiple information sources during back-propagation⁵⁰. These complex non-linear features are then passed as inputs to a dense layer

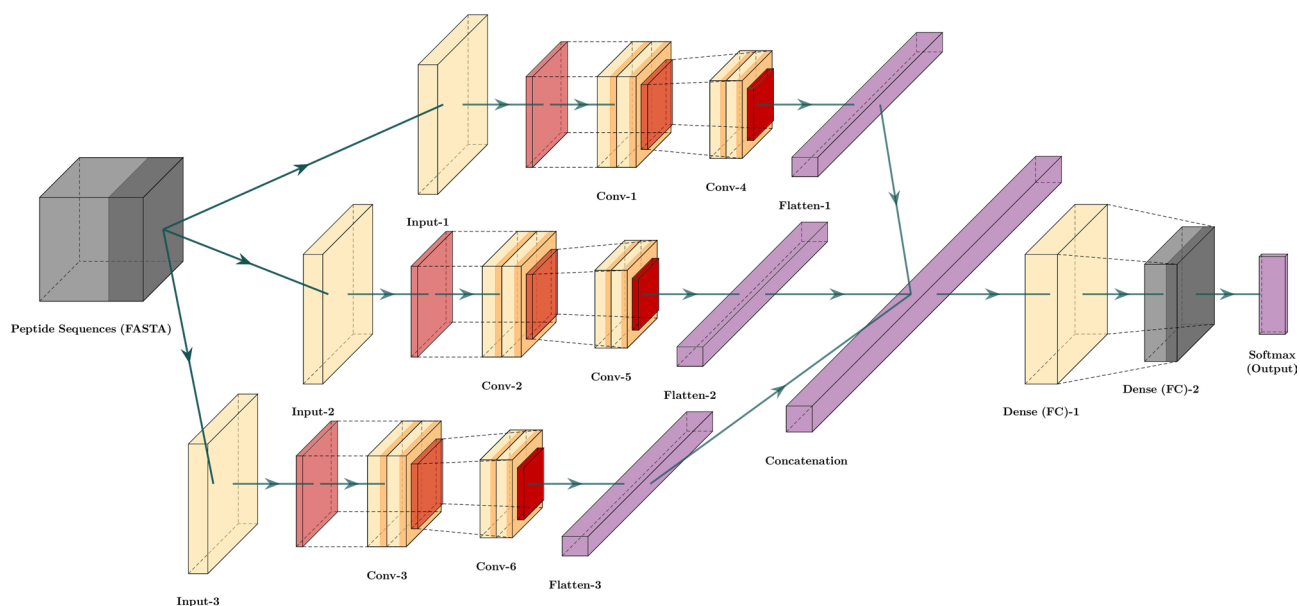


Figure 1. The general architecture of ACP-MHCNN. We extract BPF, physicochemical, and evolutionary-based features. We then feed the extracted features to a multi-headed deep convolutional neural network (MHCNN) to predict Anti-Cancer peptides.

with SoftMax activation function⁵¹, which draws a linear decision boundary on the derived feature space for separating the anticancer peptides from peptides without anticancer activity. Figure 1 represents the architecture of our proposed model for joint feature extraction from multiple information sources.

Since the training data is limited for this task, there is a possibility of overfitting when training a deep-CNN model. To avoid overfitting, we adopt both L2 regularization and dropout methods in the feature extraction step to build out model⁵². L2 and dropout have been shown to be effective methods to address overfitting issue when the number of training samples are limited⁵². To be specific, the feature extraction occurs in all layers of the three parallel convolutional groups and the dense recombination layer after concatenation. Therefore, here high dropout rates (≥ 0.5) are employed after each of these layers during the training phase to mitigate overfitting. These dropout rates are determined through cross-validation. Note that, the three convolutional layer groups that extract features from three distinct sequence representations are jointly trained alongside the dense recombination layer for minimizing cross entropy loss function⁵³. Therefore, our model can simultaneously interact with the three information sources for detecting complex and ambiguous patterns. Optimal values for our model's weights and biases are learned by employing Adam optimizer⁵⁰ with a learning rate determined through cross-validation.

ACP-DL, the only deep learning-based architecture proposed to date for anticancer peptide identification, employed stacked bidirectional LSTM layers for feature extraction which is an intuitive choice given a recurrent model's capability of capturing global sequence-order information³². However, the recurrent connections and the gates also introduce a large number of parameters that need to be tuned. This can lead to overfitting since the number of training instances is limited. Moreover, since only 15 N-terminus amino acids have been considered for feature extraction, LSTM's long-range sequence-order-effect detection capabilities remain underutilized while the parameter overhead remains³². In this study, we do not add any recurrent layer on top of the output feature maps from the final convolutional layers to avoid this issue.

Furthermore, it is to be noted that the kernels in the final layer of each convolutional group have an effective receptive field of length 6 due to hierarchical relationship between the stacked layers (length 4 kernels to length 3 kernels)⁴⁷. This effective receptive field should provide sufficient coverage for extracting both short-range and long-range patterns from sub-sequences of length 15. In addition, since we extract features from short sub-sequences, reducing the temporal dimension of the intermediate feature maps (outputs of the first and second convolutional layers of each group) is not required for learning higher order features. Hence, we do not add any pooling layers between the feature extraction layers within the convolutional groups⁴⁷. The absence of pooling layers also reduces potential loss of sequence order information that can be exploited by the kernels in the final convolutional layers in the groups for detecting long-range patterns⁴⁷.

To analyse the contribution of features extracted from each of the information sources, we carry out experiments using all possible combinations of the three representations. This results in seven models (${}^3C_1 + {}^3C_2 + {}^3C_3$) with 1, 2 or 3 convolutional groups. All these combinations are summarized in Table 1. The performance for our architecture using these seven combinations is reported in the following section.

For ACP-740 and ACP-240, our model's hyperparameters are tuned on ACP-740 through cross-validation, and the same model configuration is used for ACP-240. For ACP-500 and ACP-164, hyperparameter tuning is performed on ACP-500 through cross-validation. ACP-240 and ACP-164 have been kept untouched during hyperparameter tuning to avoid performance overestimation. Table 2 shows detailed hyperparameter configurations for different ACP identification datasets used in this study.

Combination number	Feature encoding technique	Number of convolutional layer groups
C1	BPF	1
C2	Physicochemical Properties	1
C3	Evolutionary Information	1
C4	BPF & Physicochemical Properties	2
C5	BPF & Evolutionary Information	2
C6	Physicochemical Properties & Evolutionary Information	2
C7	BPF & Physicochemical Properties & Evolutionary Information	3

Table 1. Summary of seven combinations of the three sequence representations explored in this study. On the first column of the table, we present the name of the combination, on the second column we present the name of the representations used to build the given combination, and in the third column we present the number of convolutional groups for the given combination.

ACP-740 and ACP-240			ACP-500 and ACP-164		
Convolutional group-1			Convolutional group-1		
Conv-1			Conv-		
filter = 10	kernel = 4	drop = 0.8	filter = 16	kernel = 3	drop = 0.7
Conv-2			Conv-2		
filter = 8	kernel = 3	drop = 0.7	filter = 8	kernel = 3	drop = 0.5
Convolutional group-2			Convolutional Group-2		
Conv-1			Conv-1		
filter = 10	kernel = 4	drop = 0.8	filter = 16	kernel = 3	drop = 0.7
Conv-2			Conv-2		
filter = 8	kernel = 3	drop = 0.7	filter = 8	kernel = 3	drop = 0.5
Convolutional Group-3			Convolutional Group-3		
Conv-1			Conv-1		
filter = 10	kernel = 4	drop = 0.8	filter = 16	kernel = 3	drop = 0.7
Conv-2			Conv-2		
filter = 8	kernel = 3	drop = 0.7	filter = 8	kernel = 3	drop = 0.5
Dense recombination			Dense recombination		
Dense-1			Dense-1		
units = 8	drop = 0.7		units = 16	drop = 0.6	
			Dense-2		
			units = 8	drop = 0.5	

Table 2. Hyperparameter configurations employed for different ACP datasets. In this table, 'Conv' = a convolutional layer, 'Dense' = a fully connected layer, 'filter' = number of filters in a convolutional layer, 'kernel' = size of filters in a convolutional layer, 'drop' = dropout rate, and 'units' = number of neurons in a fully connected layer.

Results and discussion

In this section, we present how we carry out the performance evaluation of our proposed model, our achieved results, and then discuss them.

Evaluation metrics. The evaluation metrics that have been used for measuring the performance of our classification method are Accuracy, Sensitivity, Specificity, Precision, and Matthews correlation coefficient (MCC). These metrics are described through the following equations:

$$Accuracy = \frac{tp + tn}{tp + tn + fp + fn} * 100 \quad (2)$$

$$Sensitivity = \frac{tp}{tp + fn} * 100 \quad (3)$$

$$Specificity = \frac{tn}{tn + fp} * 100 \quad (4)$$

Combination	Accuracy (STD)	Sensitivity (STD)	Specificity (STD)	Precision (STD)	MCC (STD)
C1	76.0 (2.9)	78.9 (7.8)	73.0 (8.1)	75.0 (6.2)	0.52 (0.02)
C2	73.1 (4.8)	74.7 (13.5)	71.3 (11.6)	72.8 (11.6)	0.46 (0.11)
C3	81.1 (3.1)	81.3 (3.7)	80.7 (3.7)	81.3 (4.1)	0.62 (0.05)
C4	76.9 (2.9)	75.7 (7.5)	78.4 (2.9)	78.2 (2.5)	0.54 (0.05)
C5	84.0 (3.7)	87.6 (8.3)	80.3 (4.2)	82.0 (3.7)	0.68 (0.07)
C6	81.8 (3.2)	82.9 (3.3)	81.1 (5.2)	81.8 (4.2)	0.64 (0.07)
C7	86.0 (1.6)	88.9 (3.2)	83.1 (4.4)	84.4 (3.9)	0.72 (0.03)

Table 3. Results achieved using fivefold cross validation for ACP-740 dataset for different input feature groups. The STD is also presented in the brackets for each measurement. Bold items indicate the best values found by the methods.

Combination	Accuracy (STD)	Sensitivity (STD)	Specificity (STD)	Precision (STD)	MCC (STD)
C1	73.5 (3.1)	82.7 (9.9)	63.6 (8.8)	72.9 (9.4)	0.47 (0.06)
C2	71.2 (4.5)	82.3 (11.0)	59.6 (14.9)	70.6 (4.6)	0.43 (0.07)
C3	79.1 (2.1)	84.6 (6.0)	72.7 (6.0)	78.6 (5.9)	0.58 (0.08)
C4	75.1 (4.4)	84.6 (4.4)	63.6 (7.1)	73.3 (6.4)	0.50 (0.08)
C5	79.9 (2.3)	85.4 (5.9)	73.6 (15.8)	79.3 (1.1)	0.60 (0.08)
C6	81.5 (1.9)	83.2 (8.6)	79.6 (9.3)	82.8 (5.8)	0.63 (0.08)
C7	83.0 (1.1)	90.1 (5.1)	75.6 (3.5)	81.1 (3.9)	0.67 (0.04)

Table 4. Results achieved using fivefold cross validation for ACP-240 dataset for different input feature groups. The STD is also presented in the brackets for each measurement. Bold items indicate the best values found by the methods.

Combination	Accuracy	Sensitivity	Specificity	Precision	MCC
C1	83.8	85.4	81.6	82.3	0.67
C2	74.2	77.9	70.6	72.6	0.49
C3	89.0	91.4	86.6	87.2	0.78
C4	85.6	88.7	82.6	83.6	0.71
C5	90.0	93.7	86.3	87.3	0.80
C6	88.4	89.4	86.7	87.1	0.76
C7	91.0	97.6	84.2	86.0	0.82

Table 5. Results achieved using independent test for ACP-500/164 dataset. Model trained on ACP-500 and tested on ACP-164. Bold items indicate the best values found by the methods.

$$Precision = \frac{tp}{tp + fp} * 100 \quad (5)$$

$$MCC = \frac{(tp * tn) - (fp * fn)}{\sqrt{(tp + fp)(tp + fn)(tn + fp)(tn + fn)}} \quad (6)$$

where, tp is the number of correctly predicted positive instances, tn is the number of correctly predicted negative instances, fp is the number of incorrectly predicted negative instances, and fn is the number of incorrectly predicted positive instances. The range of values for Accuracy, Sensitivity, Specificity, and Precision is 0 to 100 percent. 100% represents an ideal classifier (totally accurate) and 0% represents the worst possible model (totally inaccurate). In addition, MCC has a range from -1 to $+1$. A value of 0 in MCC represent a random classifier with no correlation, $+1$ represent perfect positive correlation and -1 represents perfect negative correlation.

Contribution analysis for different sequence representations. For each of the representation combinations summarized in Table 1, we have performed experiments on ACP-740 and ACP-240 using fivefold-cross validation, and the corresponding results are reported in Table 3 and 4, respectively. For ACP-500 and ACP-164, we train and tune the models on ACP-500 and test them on ACP-164. The corresponding results are reported in Table 5.

As shown in Table 3, for the ACP-740 dataset, among the single-representation combinations (C1, C2, and C3), the representation depicting evolutionary information of the amino acid residues (C3) performs better compared to BPF and physicochemical-based representations (C1 and C2) on all six performance measures. As shown in Tables 4 and 5, similar results are observed for single representation models for ACP-240 and ACP-164. These results indicate that when it comes to feature extraction from a single peptide representation, evolutionary information contributes the most for separating the ACPs from the non-ACPs compared to BPF and physicochemical-based representation.

Among the two-representation combinations (C4, C5, and C6), C5 (BPF + evolutionary), and C6 (physicochemical property + evolutionary information) performs better than C4 (BPF + physicochemical property) which further underscores the importance of the features extracted from evolutionary information for ACP identification. Moreover, C5 and C6 (two-representation combinations containing evolutionary information) perform better than C3 (the best performing single-representation combination containing evolutionary information only). This aspect of the results manifests that our proposed joint pattern extraction strategy from multiple representations through parallel-convolutional-groups can effectively embellish the features learned from a strong primary representation (evolutionary information in this case) through potential ambiguity resolution using other secondary representations (BPF and physicochemical property-based information in this case).

This hypothesis has been further corroborated by the performance of the all-representation combination (C7) on all datasets. As shown in Tables 3, 4, and 5, the model trained on C7 consisting of three parallel convolutional groups that extract features from all three representations performs better than the other combinations (C1 to C6). Therefore, we use this all-representation combination model to train ACP-MHCNN and compare its performance with state-of-the-art methods in the following subsection. To provide more insight into our achieved results, we present receiver operating characteristic (ROC) curves for our achieved results. The ROC curve for ACP-740 (using fivefold cross validation), ACP-240 (using fivefold cross validation), and ACP-164 (using ACP-500 as the training dataset) are shown in Figs. 2, 3, and 4, respectively. The results for ACP-MHCNN when it is trained on ACP-740 dataset and tested on ACP-240 and ACP-164 datasets are provided in Table S1.

As shown in these figures, we constantly achieve very high Area Under the Curve (AUC) value. We achieve 0.90, 0.88, and 0.93 for ACP-740, ACP-240, and ACP-164, respectively. The consistent AUC achieved on these three benchmarks using different evaluation methods demonstrates the generality of our proposed model. In addition, achieving 0.93 in AUC on ACP-164 which is an independent test set demonstrates the potential of ACP-MHCNN on identifying ACP for new unseen samples.

We perform additional experiments to study the performance of our proposed method when full sequences are utilized instead of partial sequences. For these experiments, the longest sequence in each dataset was kept untouched and rest of the sequences were post-padded accordingly for matching the longest sequence's length⁴². These results are reported in Tables 6, 7, and 8, respectively.

By comparing Tables 6 (ACP-740 full sequence), Table 7 (ACP-240 full sequence), and Table 8 (ACP-500/164 full sequence) with Tables 3 (ACP-740 partial sequence), Table 4 (ACP-240 partial sequence), and Table 5 (ACP-500/164 partial sequence), respectively, it can be observed that using full sequences degrade our model's performance for most of the representation combinations. Moreover, for all three datasets, the performance of the model with the all-representation combination (C7) degrades significantly (for ACP-240, C7 performs much worse compared to C3) when full sequences are used. These observations suggest that using k N-terminus sequence performs better than complete sequences for ACP identification task using the current version of our model.

One of the potential causes behind performance degradation using full sequence is that the sufficient effective receptive field assumptions for long-range pattern extraction discussed in “Multi-headed convolutional neural network architecture” no longer holds when long sequences are used. These results have corroborated our decision of considering only k N-terminus residues for feature extraction.

We also compared ACP-MHCNN with some of the widely used classical Machine Learning classifiers in similar studies such as Support Vector Machine (SVM), Random Forest RF, Extra Tree (ET), eXtreme Gradient Boosting (XGB), k-Nearest Neighbours (KNN), Decision Tree (DT), Naive Bayes (NB), and Adaptive Boosting (AB)^{54–56}. To do this, we convert BPF, Physicochemical Properties, and Evolutionary Information to vector from matrix and use to train these classifiers. The result for this comparison on ACP-740, ACP-240, and ACP-500/164 are shown in Table 9. As shown in this Table, ACP-MHCNN significantly outperform these classifiers. The main reason is the ability of ACP-MHCNN to automatically extract related features from the input matrix compared to traditional ML models which require further steps to extract relevant information. Such comparison demonstrates the importance of automated feature extraction to enhance the prediction performance.

Comparison with state-of-the-art methods. In this section, we compare ACP-MHCNN with ACP-DL as the state of the art and also the only DL based ACP identification model proposed to date³². Yi et al., tested their proposed ACP-DL on ACP-740 and ACP-240 datasets using 5-fold cross-validation. We use the same evaluation strategies and metrics for a fair comparison while estimating our ACP-MHCNN's performance on ACP-740 and ACP-240 datasets. To investigate the generality of ACP-MHCNN even further, we compare it with ACP-DL on ACP500/ACP164 dataset as well. In this experiment, ACP-500 is used for training and tuning the model, and ACP-164 is used as the independent dataset. During all these experiments, ACP-DL is trained using the implementation details available in the accompanying GitHub repository (<https://github.com/haichengyi/ACP-DL>). It is to be noted that, during our experiments, ACP-DL obtained accuracies of 80% and 81.3% on ACP-740 and ACP-240, respectively.

Comparison between ACP-MHCNN and ACP-DL on all the datasets is shown in Table 10. As shown in this table, ACP-MHCNN outperforms ACP-DL on all datasets for every evaluation metric. To be precise, on ACP-740,

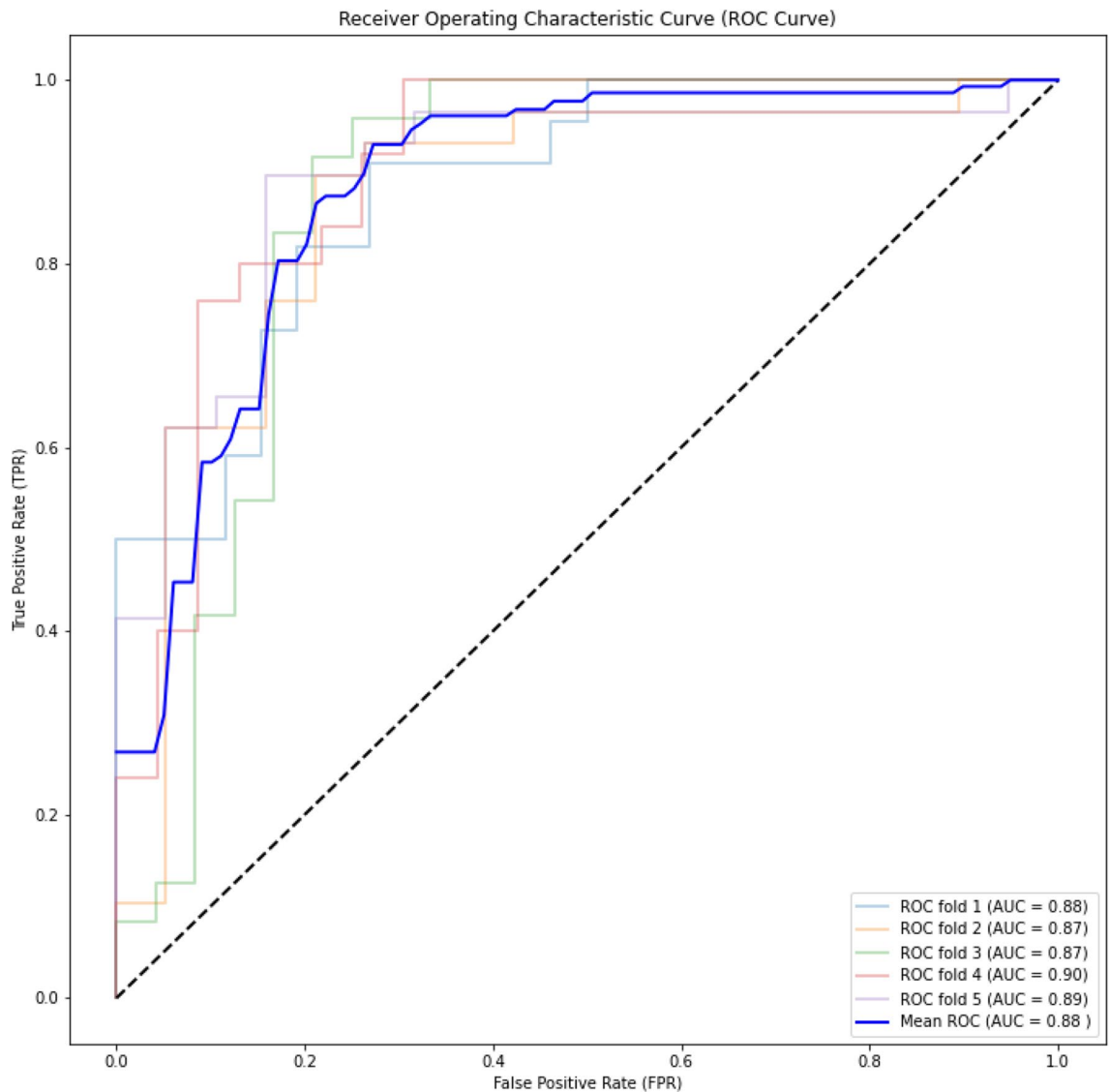


Figure 2. ROC curve for ACP-740 dataset for the fivefold cross-validation on the experiment. As shown in these figures, we constantly achieve very high Area Under the Curve (AUC) value.

ACP-MHCNN scores 6.0%, 7.5%, 4.5%, 4.7%, and 0.12 more than ACP-DL in terms of accuracy, sensitivity, specificity, precision, and MCC, respectively. Similarly, on ACP-240 ACP-MHCNN scores 1.8%, 6.0%, 4.4% and 0.02 more than ACP-DL in terms of accuracy, specificity, and MCC, respectively.

ACP-MHCNN also significantly outperforms ACP-DL on the ACP-500/ACP-164 dataset that was used to investigate the generalizability of our model. On ACP-500/ACP-164 ACP-MHCNN outperforms ACP-DL by 6.3%, 8.6%, 3.7%, 4.0%, and 0.20 in terms of accuracy, sensitivity, specificity, precision, and MCC respectively. ACP-MHCNN and its relevant codes as well as the datasets used in this study are all publicly available at: <https://github.com/mrzResearchArena/Anticancer-Peptides-CNN>. ACP-MHCNN is also publicly available as an online predictor at: <https://anticancer.pythonanywhere.com>.

Additionally, we have trained and tested ACP-MHCNN on two datasets proposed by Agrawal et al. in the recently published method AntiCP 2.0¹⁷. The two datasets are main and alternate and contain their respective training and external validation partitions. ACP-MHCNN has substantially outperformed ACP-DL on both datasets. These results are shown in Table 11.

Table 11 clearly shows ACP-MHCNN outperforms ACP-DL by a large margin. We also compare ACP-MHCNN with several existing ACP identification methods on both main and alternate datasets used in¹⁷, and the results are shown in Table 12. This comparison shows that ACPred-LAF¹⁶, iACP-FSCM⁵⁷, and AntiCP-2.0¹⁷ slightly outperforms ACP-MHCNN, and all outperform other existing methods by significant margin on these two specific datasets. It is worth noting that, since AntiCP-2.0 and all of the existing methods reported in Table 12 are traditional machine learning models while ACP-MHCNN is composed of several convolutional layers with much larger effective hypotheses space, the sizes of the training partitions of main and alternate datasets are the

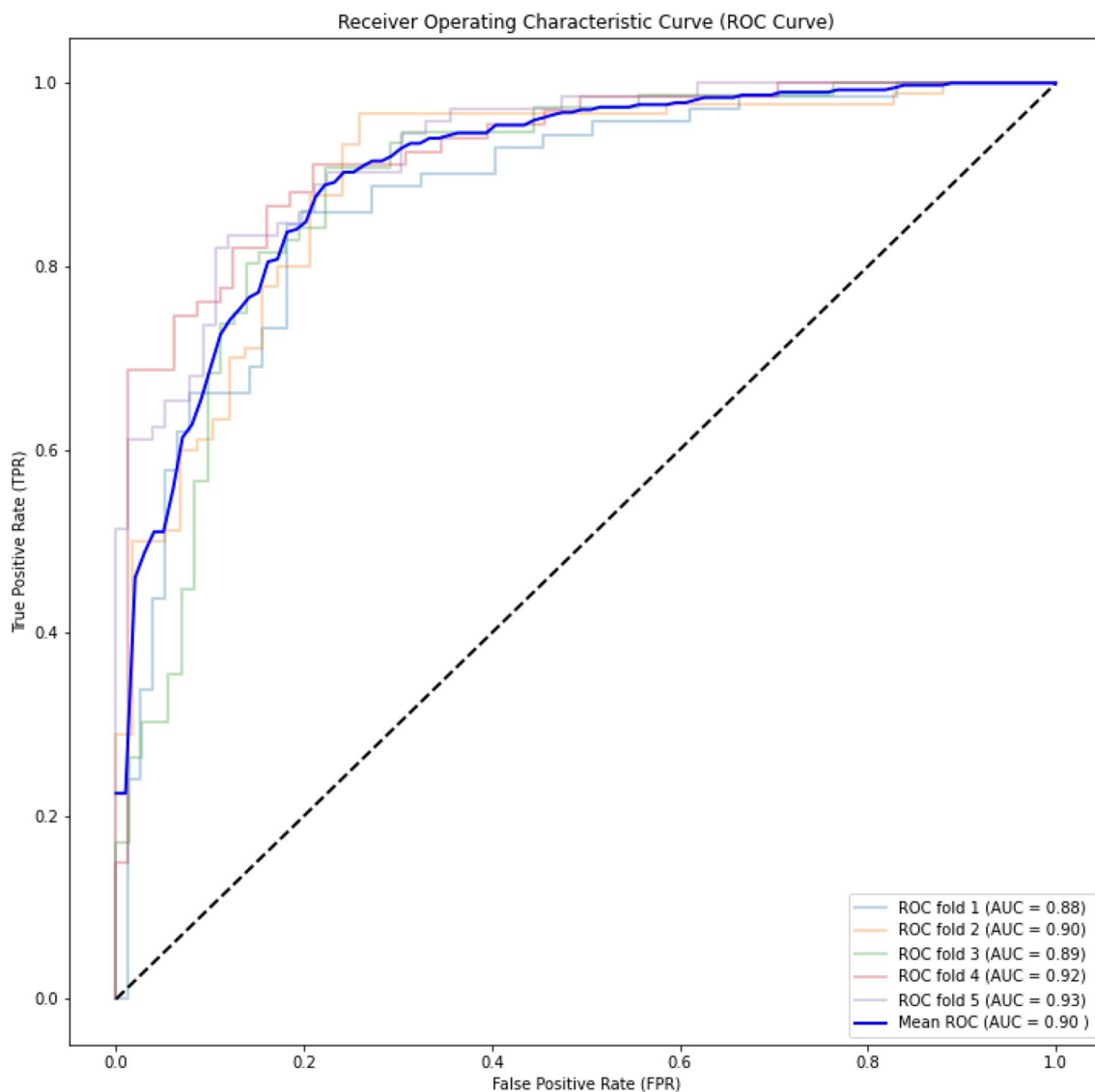


Figure 3. ROC curve for ACP-240 dataset for the fivefold cross-validation on the experiment. Similar to the results reported for ACP-740 dataset, we constantly achieve very high Area Under the Curve (AUC) value.

bottleneck for ACP-MHCNN when it comes to generalization capability. In future work, we need to mitigate this limitation through some data augmentation scheme or self-supervised pre-training or both.

Conclusion

In this study, we propose a new deep neural network architecture called ACP-MHCNN consisting of parallel convolutional groups which jointly learn and combine features from three different peptide representation methods for accurate identification of ACPs. The architecture extracts sequence-based features from residue-order information (using BPF representation), physicochemical property-based features from 31 bit-vector representation of the residues (elements of these vectors depict various physicochemical properties of the amino acids), and evolutionary features from BLOSUM62 matrix-based representation of the peptides.

Although hand-engineered features extracted from these information sources have been successfully employed for ACP identification, automatic feature extraction has hardly been explored for this problem. Before this study, ACP-DL was the only method that has used deep feature extraction for ACP identification³². It has used recurrent layers for extracting features from two residue-order-based peptide representations and leaves significant room for improvement. In the current study, we attempt to address the limitations of ACP-DL by improving the sequence representation and feature extraction methods. For sequence representation, we consider the residues' evolutionary and physicochemical characteristics alongside their ordering so that the downstream feature extraction layers can embed the sequences in spaces with additional discriminative information. For feature extraction, we jointly train three parallel convolutional layer groups so that the combined feature vector contains discriminative patterns extracted from three sources. Our method's performance could improve further by incorporating some carefully chosen manually extracted features that have been applied successfully in

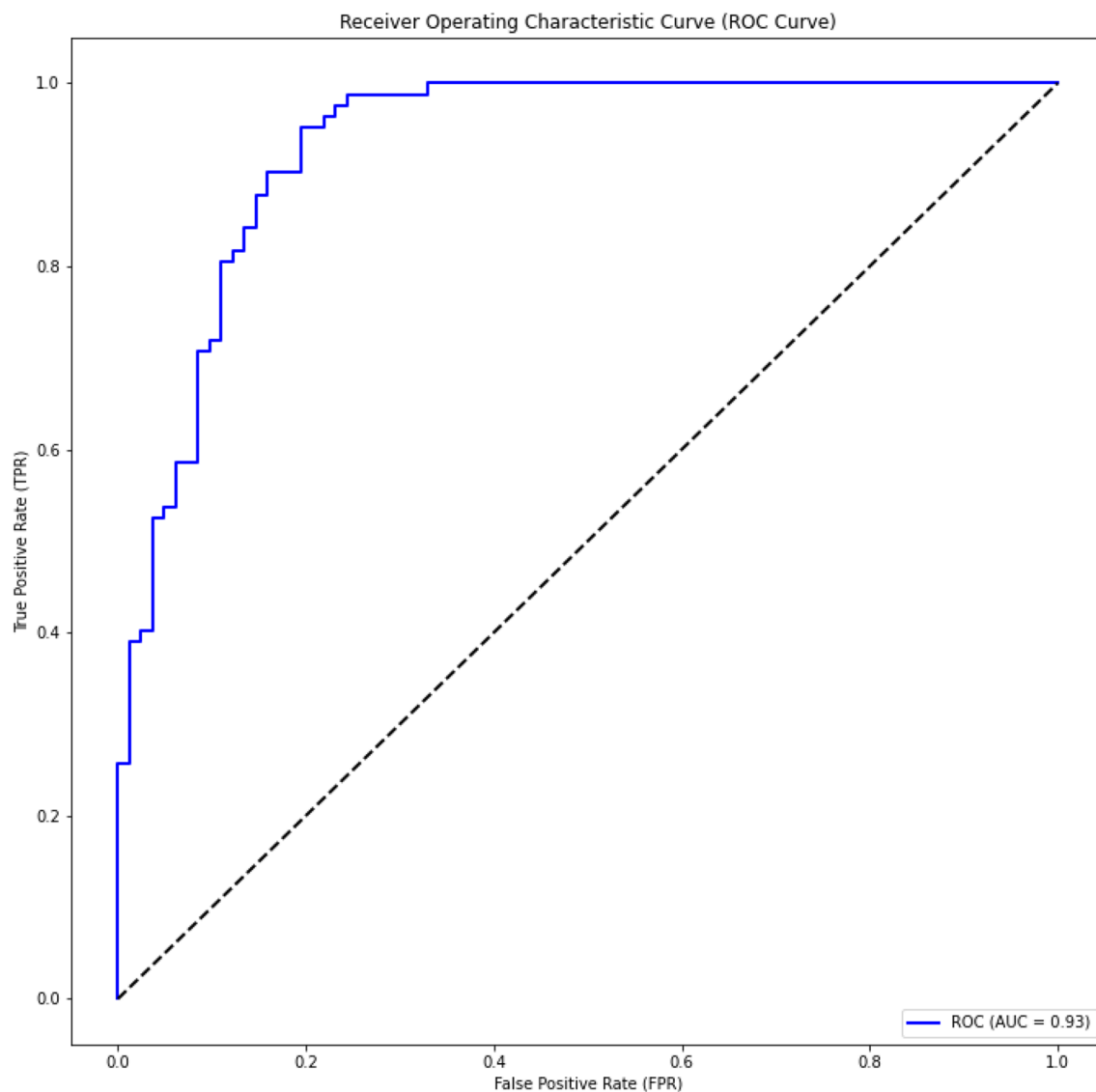


Figure 4. ROC curve for ACP-500/164. Here we used ACP-500 as a training dataset and ACP-164 as a testing dataset on the experiment.

Combination	Accuracy (STD)	Sensitivity (STD)	Specificity (STD)	Precision (STD)	MCC (STD)
C1	78.2 (1.5)	82.5 (8.2)	74.1 (8.8)	77.2 (6.0)	0.57 (0.03)
C2	71.1 (5.6)	69.9 (16.9)	72.5 (13.7)	73.9 (13.7)	0.44 (0.11)
C3	81.0 (3.3)	81.4 (4.1)	81.7 (3.7)	82.0 (4.5)	0.63 (0.07)
C4	77.1 (3.0)	74.1 (8.0)	80.8 (3.1)	79.9 (2.6)	0.55 (0.06)
C5	82.9 (4.1)	86.7 (9.2)	78.8 (4.7)	80.9 (3.7)	0.66 (0.09)
C6	81.3 (3.8)	81.6 (3.8)	81.2 (5.7)	81.9 (4.3)	0.63 (0.08)
C7	83.2 (1.7)	80.4 (4.5)	84.8 (5.4)	84.9 (4.3)	0.65 (0.03)

Table 6. Results achieved using fivefold cross validation for ACP-740 dataset (Complete sequences utilized instead of 15 N-terminus amino acids). The STD is also presented in the brackets for each measurement.

different ACP identification methods through a fourth parallel track with fully connected layers. Additionally, since the BPF representation is sparse, our feature extraction method could benefit from adding an embedding layer at the beginning of the BPF track. Once more experimental training data is available, we would be able to incorporate more parameters in our model without the risk of overfitting and explore these directions. Additionally, we would like to employ embedding techniques used in natural language processing (NLP) tasks,

Combination	Accuracy (STD)	Sensitivity (STD)	Specificity (STD)	Precision (STD)	MCC (STD)
C1	75.4 (4.3)	81.6 (8.2)	71.8 (9.2)	76.5 (10.9)	0.53 (0.07)
C2	62.6 (4.8)	77.6 (16.9)	44.5 (15.0)	63.1 (5.6)	0.25 (0.08)
C3	82.1 (4.0)	86.4 (4.1)	78.6 (6.7)	82.3 (6.0)	0.65 (0.09)
C4	79.0 (5.5)	81.9 (8.0)	75.3 (8.2)	79.1 (7.1)	0.57 (0.08)
C5	78.2 (2.8)	84.5 (9.2)	67.0 (13.2)	77.1 (2.2)	0.53 (0.08)
C6	77.0 (4.4)	81.8 (3.8)	70.8 (9.2)	77.2 (5.1)	0.54 (0.09)
C7	78.1 (2.8)	85.6 (4.5)	68.9 (4.5)	76.2 (4.5)	0.56 (0.05)

Table 7. Results achieved using fivefold cross validation for ACP-240 dataset (Complete sequences utilized instead of 15 N-terminus amino acids). The STD is also presented in the brackets for each measurement.

Combination	Accuracy	Sensitivity	Specificity	Precision	MCC
C1	82.3	86.6	78.1	79.8	0.65
C2	84.0	84.2	82.9	83.1	0.67
C3	84.1	87.8	80.5	81.8	0.68
C4	88.1	90.2	86.6	87.1	0.77
C5	85.0	87.8	81.7	82.7	0.70
C6	86.3	81.7	90.2	89.3	0.72
C7	87.2	87.8	85.4	85.7	0.73

Table 8. Results achieved using independent test for ACP-500/164 dataset (Complete sequences utilized instead of 15 N-terminus amino acids). Model trained on ACP-500 and tested on ACP-164.

Classifier	ACP-740 dataset				ACP-240 dataset				ACP-500/164 dataset			
	Acc	Sen	Spe	MCC	Acc	Sen	Spe	MCC	Acc	Sen	Spe	MCC
SVM	80.4	77.6	83.2	0.61	68.7	65.1	72.9	0.38	78.0	74.3	81.7	0.56
RF	81.2	79.2	84.8	0.64	71.0	72.0	74.7	0.48	84.1	82.9	85.3	0.68
ET	81.5	78.4	85.9	0.65	72.7	72.8	80.1	0.53	81.0	79.2	82.9	0.62
XGB	81.6	82.4	81.8	0.64	74.2	82.1	74.7	0.57	85.3	86.5	84.1	0.71
KNN	79.3	64.3	75.5	0.40	70.6	91.4	15.3	0.11	68.9	51.2	86.5	0.40
DT	78.4	76.8	70.8	0.48	70.9	75.1	68.4	0.44	78.6	71.9	85.3	0.58
NB	78.2	80.0	73.6	0.54	70.6	75.1	62.1	0.38	71.9	74.3	69.5	0.44
AB	78.1	77.3	78.5	0.56	71.3	79.0	72.0	0.52	79.8	79.2	80.4	0.60
ACP-MHCNN	86.0	88.9	83.1	0.72	83.0	90.1	75.6	0.67	91.0	97.6	84.2	0.82

Table 9. The results achieved for ACP-MHCNN compared to traditional ML models on ACP-740, ACP-240, and ACP-500/164 using fivefold cross validation. Bold items indicate the best values found by the methods.

Dataset	Model	Accuracy	Sensitivity	Specificity	Precision	MCC
ACP-740	ACP-DL	80.0	81.4	78.6	79.7	0.60
ACP-740	ACP-MHCNN	86.0	88.9	83.1	84.4	0.72
ACP-240	ACP-DL	81.3	92.0	69.6	76.7	0.64
ACP-240	ACP-MHCNN	83.0	90.1	75.6	81.1	0.67
ACP-500/ACP-164	ACP-DL	84.7	89.0	80.5	82.0	0.62
ACP-500/ACP-164	ACP-MHCNN	91.0	97.6	84.2	86.0	0.82

Table 10. Comparing the results achieved for ACP-MHCNN to ACP-DL as the state-of-the-art anticancer peptide predictor. Bold items indicate the best values found by the methods.

Dataset	Model	Accuracy	Sensitivity	Specificity	Precision	MCC
AntiCP-2.0 (Main validation)	ACP-DL	66.0	58.1	74.4	69.4	0.33
AntiCP-2.0 (Main validation)	ACP-MHCNN	73.0	78.5	67.4	70.6	0.46
AntiCP-2.0 (Alternate Validation)	ACP-DL	83.0	82.9	82.9	82.9	0.66
AntiCP-2.0 (Alternate Validation)	ACP-MHCNN	90.0	86.6	94.3	93.8	0.81

Table 11. Comparing the results achieved for ACP-MHCNN to ACP-DL as the state-of-the-art anticancer peptide predictor. Bold items indicate the best values found by the methods.

Methods	Main datasets				Alternative datasets			
	Acc	Sen	Spe	MCC	Acc	Sen	Spe	MCC
ACP-MHCNN	73.00	78.50	67.40	0.46	90.00	86.60	86.60	0.81
AntiCP-2.0	75.43	77.46	73.41	0.51	92.01	92.27	92.27	0.84
AntiCP	50.58	100.00	1.16	0.07	89.95	89.69	89.69	0.80
ACPred	53.47	85.55	21.39	0.09	85.31	87.11	87.11	0.71
ACPred-FL	44.80	67.05	22.54	-0.12	43.80	60.21	60.21	-0.15
ACPred-Fuse	68.90	69.19	68.60	0.38	78.87	64.43	64.43	0.60
PEPred-Suite	53.49	33.14	73.84	0.08	57.47	40.21	40.21	0.16
iACP	55.10	77.91	32.16	0.11	77.58	78.35	78.35	0.55
iACP-FSCM	82.50	72.60	90.30	0.65	88.90	87.60	90.20	0.78
ACPred-LAF	85.75	84.24	87.20	0.72	96.41	96.26	96.52	0.93

Table 12. Comparing the results achieved for ACP-MHCNN to the state-of-the-art anticancer peptide predictors on the main and alternative datasets used in^{16,17,60}. Bold items indicate the best values found by the methods.

such as Word2Vec⁵⁸ and FastText⁵⁹ for k-mer feature extraction. Since these embeddings are local and preserve sequence-order information, sequence representations consisting of these embeddings can be readily added as parallel branches to our model. Furthermore, inspired by the success of self-supervised pre-training on NLP tasks, several pre-trained models for protein sequences have recently been made publicly available. Among them, UDSPProt⁶⁰, a LSTM sequence model trained on unlabeled Swiss-Prot protein sequences in a self-supervised autoregressive manner has shown remarkable performance on protein-level classification tasks after fine tuning. Another convolutional transformation and attention-based model ProteinBERT⁶¹, pre-trained on sequence-correction and GO annotation prediction tasks, has shown impressive performance on protein-level regression tasks after fine tuning. We want to explore the possibility of combining ACP-MHCNN for fine tuning these pre-trained models for ACP identification in future work.

The positive effects of these improvements are manifested in the experimental results obtained on well-established ACP identification datasets, where ACP-MHCNN has significantly outperformed ACP-DL using different evaluation measures for every dataset investigated in this study. Hence, we believe our current study's findings add significantly to the existing knowledge on computational method development for ACP identification. ACP-MHCNN, its relevant codes, and the datasets used in this study are all publicly available at: <https://github.com/mrzResearchArena/Anticancer-Peptides-CNN>, ACP-MHCNN is also publicly available as an online predictor at: <https://anticancer.pythonanywhere.com>.

Received: 21 January 2021; Accepted: 17 November 2021

Published online: 08 December 2021

References

1. Tyagi, A. *et al.* In silico models for designing and discovering novel anticancer peptides. *Sci. Rep.* **3**, 1–8 (2013).
2. Shoombuatong, W., Schaduangrat, N. & Nantasenammat, C. Unraveling the bioactivity of anticancer peptides as deduced from machine learning. *EXCLI J.* **17**, 734 (2018).
3. Chen, W., Ding, H., Feng, P., Lin, H. & Chou, K. C. iACP: A sequence based tool for identifying anticancer peptides. *Oncotarget* **7**, 16895 (2016).
4. Schaduangrat, N., Nantasenammat, C., Prachayasittikul, V. & Shoombuatong, W. Acpred: A computational tool for the prediction and analysis of anticancer peptides. *Molecules* **24**(10), 1973 (2019).
5. Mader, J. S. & Hoskin, D. W. Cationic antimicrobial peptides as novel cytotoxic agents for cancer treatment. *Expert Opin. Investig. Drugs* **15**, 933–946 (2006).
6. Huang, Y., Feng, Q., Yan, Q., Hao, X. & Chen, Y. Alpha-helical cationic anticancer peptides: A promising candidate for novel anticancer drugs. *Mini Rev. Med. Chem.* **15**, 73–81 (2015).
7. Otvos, L. Jr. Peptide-based drug design: Here and now. *Methods Mol. Biol.* **494**, 1–8 (2008).
8. Boohaker, R. J., Lee, M. W., Vishnubhotla, P., Perez, J. M. & Khaled, A. R. The use of therapeutic peptides to target and to kill cancer cells. *Curr. Med. Chem.* **19**, 3794–3804 (2012).

9. Thundimadathil, J. Cancer treatment using peptides: Current therapies and future prospects. *J. Amino Acids* **2012**, 967347 (2012).
10. Hajisharifi, Z., Piryaeie, M., Beigi, M. M., Behbahani, M. & Mohabatkari, H. Predicting anticancer peptides with chous pseudo amino acid composition and investigating their mutagenicity via ames test. *J. Theor. Biol.* **341**, 34–40 (2014).
11. Manavalan, B. *et al.* Mlaccp: Machine-learning-based prediction of anticancer peptides. *Oncotarget* **8**, 77121 (2017).
12. Akbar, S., Hayat, M., Iqbal, M. & Jan, M. A. iACP-GAEnsC: Evolutionary genetic algorithm based ensemble classification of anticancer peptides by utilizing hybrid feature space. *Artif. Intell. Med.* **79**, 62–70 (2017).
13. Lei, X., Liang, G., Wang, L. & Liao, C. A novel hybrid sequence-based model for identifying anticancer peptides. *Genes* **9**, 158 (2018).
14. Kabir, M. *et al.* Intelligent computational method for discrimination of anticancer peptides by incorporating sequential and evolutionary profiles information. *Chemom. Intell. Lab. Syst.* **182**, 158–165 (2018).
15. Wei, L., Zhou, C., Chen, H., Song, J. & Ran, Su. Acpred-fl: A sequence-based predictor using effective feature representation to improve the prediction of anti-cancer peptides. *Bioinformatics* **34**(23), 4007–4016 (2018).
16. Charoenkwan, P. *et al.* Improved prediction and characterization of anticancer activities of peptides using a novel flexible scoring card method. *Sci. Rep.* **11**(1), 1–13 (2021).
17. Agrawal, P. *et al.* AntiCP 2.0: An updated model for predicting anticancer peptides. *Brief. Bioinform.* **22**(3), 153 (2021).
18. Basith, S. *et al.* Machine intelligence in peptide therapeutics: A next-generation tool for rapid disease screening. *Med. Res. Rev.* **40**(4), 1276–1314 (2020).
19. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**(7553), 436–444 (2015).
20. Daniel, Q. & Xie, X. DanQ: A hybrid convolutional and recurrent deep neural network for quantifying the function of DNA sequences. *Nucleic Acids Res.* **44**(11), e107–e107 (2016).
21. Yang, B. *et al.* BiRen: predicting enhancers with a deep-learning-based model using the DNA sequence alone. *Bioinformatics* **33**(13), 1930–1936 (2017).
22. Alipanahi, B., Delong, A., Weirauch, M. T. & Frey, B. J. Predicting the sequence specificities of dna- and rna-binding proteins by deep learning. *Nat. Biotechnol.* **33**, 831 (2015).
23. Bosco, G. L. & Di Gangi, M. A. Deep learning architectures for dna sequence classification. *Fuzzy Logic Soft Comput.* **10147**, 162–171 (2017).
24. Busia, A. *et al.* A deep learning approach to pattern recognition for short dna sequences. *BioRxiv* **2019**, 353474 (2019).
25. Rizzo, R., Fiannaca, A., La Rosa, M. & Urso, A. A deep learning approach to dna sequence classification. *Comput. Intell. Method Bioinform. Biostat.* **9874**, 129–140 (2016).
26. Wang, L., You, Z. H., Huang, D. S. & Zhou, F. Combining high speed elm learning with a deep convolutional neural network feature encoding for predicting protein-rna interactions. *IEEE/ACM Trans. Comput. Biol. Bioinform.* **17**, 972–982 (2018).
27. Zou, Q., Xing, P., Wei, L. & Liu, B. Gene2vec: Gene subsequence embedding for prediction of mammalian n6-methyladenosine sites from mRNA. *RNA* **25**(2), 205–218 (2019).
28. You, Z.-H., Lei, Y.-K., Gui, J., Huang, D.-S. & Zhou, X. Using manifold embedding for assessing and predicting protein interactions from highthroughput experimental data. *Bioinformatics* **26**(21), 2744–2751 (2010).
29. Wei, L., Ding, Y., Ran, Su., Tang, J. & Zou, Q. Prediction of human protein subcellular localization using deep learning. *J. Parallel Distrib. Comput.* **117**, 212–217 (2018).
30. Wang, Y. *et al.* Predicting protein interactions using a deep learning method-stacked sparse autoencoder combined with a probabilistic classification vector machine. *Complexity* <https://doi.org/10.1155/2018/4216813> (2018).
31. Yi, H.-C. *et al.* A deep learning framework for robust and accurate prediction of ncrnprotein interactions using evolutionary information. *Mol. Ther.* **11**, 337–344 (2018).
32. Yi, H. C. *et al.* HAcP-dl: A deep learning long short-term memory model to predict anticancer peptides using high-efficiency feature representation. *Mol. Ther.* **17**, 1–9 (2019).
33. Timmons, P. B. & Hewage, C. M. ENNAACT is a novel tool which employs neural networks for anticancer activity classification for therapeutic peptides. *Biomed. Pharmacother.* **133**, 111051 (2021).
34. Gu, J. *et al.* Recent advances in convolutional neural networks. *Pattern Recogn* **77**, 3354–3377 (2015).
35. Yamashita, R., Nishio, M., Do, R. K. G. & Togashi, K. Convolutional neural networks: An overview and application in radiology. *Insights Imaging* **9**(4), 611–629 (2018).
36. Shin, H. *et al.* Deep convolutional neural networks for computer-aided detection: CNN architectures, dataset characteristics and transfer learning. *IEEE Trans. Med. Imaging* **35**(5), 1285–1298 (2016).
37. Amin, R. *et al.* iPromoter-BnCNN: A novel branched CNN based predictor for identifying and classifying sigma promoters. *Bioinformatics* **36**, 4869–4875 (2019).
38. Zeng, H., Edwards, M. D., Liu, G. & Gifford, D. K. Convolutional neural network architectures for predicting DNA–protein binding. *Bioinformatics* **32**(12), 121–127 (2016).
39. Zhou, X., Hu, B., Lin, J., Xiang, Y. & Wang, X. ICRCHIT: A deep learning based comment sequence labeling system for answer selection challenge. In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, 210–214 (Association for Computational Linguistics, 2015).
40. Chen, T., Ruifeng, Xu., He, Y. & Wang, X. Improving sentiment analysis via sentence type classification using bilstm-crf and cnn. *Expert Syst. Appl.* **72**, 221–230 (2017).
41. Oh, J., Wang, J. & Wiens, J. Learning to exploit invariances in clinical time-series data using sequence transformer networks. *CoRR* **85**, 332–347 (2018).
42. Tyagi, A. *et al.* CancerPPD: A database of anticancer peptides and proteins. *Nucleic Acids Res.* **43**, D837 (2015).
43. Dwarampudi, M. & Reddy, N. V. *Effects of Padding on LSTMs and CNNs*. arXiv preprint. [arXiv:1903.07288](https://arxiv.org/abs/1903.07288) (2019).
44. Basak, P. *et al.* An evolutionary analysis identifies a conserved pentapeptide stretch containing the two essential lysine residues for rice l-myo-inositol 1-phosphate synthase catalytic activity. *PLoS ONE* **12**(9), e0185351 (2017).
45. Das, J. K., Das, P., Ray, K. K., Choudhury, P. P. & Jana, S. S. Mathematical characterization of protein sequences using patterns as chemical group combinations of amino acids. *PLoS ONE* **11**(12), e0167651 (2016).
46. Henikoff, S. & Henikoff, J. G. Amino acid substitution matrices from protein blocks. *Proc. Natl. Acad. Sci. USA* **89**(22), 10915–10919 (1992).
47. Koo, P. K. & Eddy, S. R. Representation learning of genomic sequence motifs with convolutional neural networks. *PLoS Comput. Biol.* **15**(12), e1007560 (2019).
48. Kohavi, R. A study of cross-validation and bootstrap for accuracy estimation and model selection. *Ijcai* **14**(2), 1137–1145 (1995).
49. Yarotsky, D. Error bounds for approximations with deep relu networks. *Neural Netw.* **94**, 103–114 (2017).
50. Kingma, D. P. & Ba, J. *Adam: A Method for Stochastic Optimization*. arXiv preprint. [arXiv:1412.6980](https://arxiv.org/abs/1412.6980)(2014)
51. Narayan, S. The generalized sigmoid activation function: Competitive supervised learning. *Inf. Sci.* **99**(1–2), 69–82 (1997).
52. Kukačka, J., Golkov, V., & Cremers, D. *Regularization for Deep Learning: A Taxonomy*. arXiv preprint. [arXiv:1710.10686](https://arxiv.org/abs/1710.10686) (2017)
53. Janocha, K., & Czarnecki, W. M. *On Loss Functions for Deep Neural Networks in Classification*. arXiv preprint. [arXiv:1702.05659](https://arxiv.org/abs/1702.05659) (2017)
54. Dipta, S. R. *et al.* SEMal: Accurate protein malonylation site predictor using structural and evolutionary information. *Comput. Biol. Med.* **125**, 104022 (2020).

55. Muhammad, R. *et al.* PyFeat: A Python-based effective feature generation tool for DNA, RNA and protein sequences. *Bioinformatics* **35**, 3831–3833 (2019).
56. Jani, M. R. *et al.* iRecSpot-EF: Effective sequence based features for recombination hotspot prediction. *Comput. Biol. Med.* **103**, 17–23 (2018).
57. He, W., Wang, Y., Cui, L., Su, R. & Wei, L. Learning embedding features based on multisense-scaled attention architecture to improve the predictive performance of anticancer peptides. *Bioinformatics* <https://doi.org/10.1093/bioinformatics/btab560> (2021).
58. Goldberg, Y. & Levy, O. *word2vec Explained: Deriving Mikolov et al.'s Negative-Sampling word-Embedding Method*. arXiv preprint. [arXiv:1402.3722](https://arxiv.org/abs/1402.3722) (2014).
59. Athiwaratkun, B., Wilson, A. G. & Anandkumar, A. *Probabilistic Fasttext for Multi-sense Word Embeddings*. arXiv preprint. [arXiv:1806.02901](https://arxiv.org/abs/1806.02901) (2018).
60. Strodthoff, N. *et al.* UDSPMProt: Universal deep sequence models for protein classification. *Bioinformatics* **36**(8), 2401–2409 (2020).
61. Brandes, N. *et al.* ProteinBERT: A universal deep-learning model of protein sequence and function. *bioRxiv* <https://doi.org/10.1101/2021.05.24.445464> (2021).

Author contributions

S.A. conceived and initiated this study. S.A., R.M.Z.H. performed the experiments. S.A., S.Ad., A.S., S.S. and A.D. wrote the manuscript. Z.H. helped with literature review and designed the online server. A.D., S.S. mentored and analytically reviewed the paper. All the authors reviewed the article.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-021-02703-3>.

Correspondence and requests for materials should be addressed to S.S. or A.D.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021