
AMPCliff: quantitative definition and benchmarking of activity cliffs in antimicrobial peptides

Kewei Li¹, Yuqian Wu², Yutong Guo³, Yinheng Li⁴, Yusi Fan¹, Ruochi Zhang^{5,*}, Lan Huang¹, Fengfeng Zhou^{1,6,*}.

1 College of computer Science and Technology, and Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education, Jilin University, Changchun, Jilin, China, 130012.

2 School of Software, Jilin University, Changchun 130012, Jilin, China.

3. School of Life Sciences, Jilin University, Changchun 130012, Jilin, China.

4 Department of Computer Science, Columbia University, 116th and Broadway, New York City, New York, 10027, United States.

5 School of Artificial Intelligence, and Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education, Jilin University, Changchun, Jilin, China, 130012.

6 School of Biology and Engineering, Guizhou Medical University, Guiyang 550025, Guizhou, China.

* Correspondence may be addressed to Fengfeng Zhou: FengfengZhou@gmail.com or ffzhou@jlu.edu.cn. Correspondence may also be addressed to Ruochi Zhang: zrc720@gmail.com.

Emails: Kewei Li, kwbb1997@gmail.com; Yuqian Wu, 924476388@qq.com; Yutong Guo, 3291865187.qq.com; Yinheng Li, yl4039@columbia.edu; Yusi Fan, fan_yusi@163.com; Ruochi Zhang, zrc720@gmail.com; Lan Huang, huanglan@jlu.edu.cn; Fengfeng Zhou, FengfengZhou@gmail.com.

Abstract

Activity cliff (AC) is a phenomenon that a pair of similar molecules differ by a small structural alternation but exhibit a large difference in their biochemical activities. The AC of small molecules has been extensively investigated but limited knowledge is accumulated about the AC phenomenon in peptides with canonical amino acids. This study introduces a quantitative definition and benchmarking framework AMPCliff for the AC phenomenon in antimicrobial peptides (AMPs) composed by canonical amino acids. A comprehensive analysis of the existing AMP dataset reveals a significant prevalence of AC within AMPs. AMPCliff quantifies the activities of AMPs by the metric minimum inhibitory concentration (MIC), and defines 0.9 as the minimum threshold for the normalized BLOSUM62 similarity score between a pair of aligned peptides with at least two-fold MIC changes. This study establishes a benchmark dataset of paired AMPs in *Staphylococcus aureus* from the publicly available AMP dataset GRAMPA, and conducts a rigorous procedure to evaluate various AMP AC prediction models, including nine machine learning, four deep learning algorithms, four masked language models, and four generative language models. Our analysis reveals that these models are capable of detecting AMP AC events and the pre-trained protein language ESM2 model demonstrates superior performance across the evaluations. The predictive performance of AMP activity cliffs remains to be further improved, considering that ESM2 with 33 layers only achieves the Spearman correlation coefficient=0.50 for the regression task of the MIC values on the benchmark dataset. Source code and additional resources are available at <https://www.healthinformatics-lab.org/supp/> or <https://github.com/Kewei2023/AMPCliff-generation>.

Keywords: AMPCliff; activity cliff (AC); minimum inhibitory concentration (MIC); benchmark; antimicrobial peptide (AMP).

Introduction

The development of peptide drugs and the application of Quantitative Structure-Activity Relationship (QSAR) modeling, enhanced by advances in AI technologies like deep learning, have become prominent in scientific research and industry¹⁻³. However, despite its potential, deep learning faces challenges in interpretability, data quality

dependency, and generalization, especially evident in handling activity cliffs in drug design^{2,4-7}. For instance, in the Multi-Drug Resistance Protein 1 (MDR1) dataset, nearly half of the molecular pairs are activity cliffs. Compared to non-activity cliff pairs, the Root Mean Square Error (RMSE) for activity cliff pairs has increased by at least 0.6⁶. Addressing these challenges to improve how activity cliffs are computationally represented remains a critical and pressing issue in computer-aided drug design.

Activity cliffs, which are pairs of similar compounds that only differ by a small structural modification but exhibit a large difference in their binding affinity for a given target^{4,5,8-16}, have been a persistent challenge in QSAR modeling for over 30 years⁵. This phenomenon affects tasks from virtual screening to hit-to-lead and lead optimization in drug development^{15,17}. Recent efforts in QSAR modeling for antimicrobial peptides (AMP) include using RNN-based models by Junjie Huang et al., who developed a "classification-rank-regressive" pipeline¹⁸, and by Amir Pandi et al., who adapted regression models to classification by setting thresholds¹⁹. Despite these advancements, the models have shown limited success, with the best reported Pearson correlation of 0.77. We found out that removing activity cliffs from training data significantly improves model performance to 0.9 (no matter Pearson or Spearman metrics), which was the similar value with random split in the literature¹⁶. It suggests that activity cliffs might be a critical factor in the observed challenges of building effective QSAR models for AMPs.

The study of activity cliffs in small molecules is segmented into three phases: (i) identification, (ii) data-driven causal analysis, and (iii) predictive modeling. Since 2006, the concept has evolved from a broad to more refined definitions^{9-11,15,20} such as MMP-Cliff⁹ and 3D activity cliffs¹¹, along with various metrics^{8,13,21} like SARI²¹ and SALI⁸. The MMP-Cliff⁹, defined in 2012 by Xiaoyue Hu and Jürgen Bajorath's team, remains the most prevalent form, focusing on structural differences within molecule pairs. This team also explored 3D activity cliffs¹¹, which highlighted the challenges of applying 2D measures to 3D structures, emphasizing the importance of molecular representation in activity cliff characterization. Despite advances, the exploration of 3D cliffs has stalled since 2015¹⁴, possibly due to limited data availability. Currently, the MMP-Cliff is the primary model used in activity cliff prediction, defining as follows: for a qualifying MMP, the difference in size of the exchanged fragments was limited to at most eight non-hydrogen atoms, and the maximal size of an exchanged fragment was set to 13 non-hydrogen atoms. Furthermore, the potency difference between

compounds in an MMP meeting the structural criteria had to be at least 2 orders of magnitude.

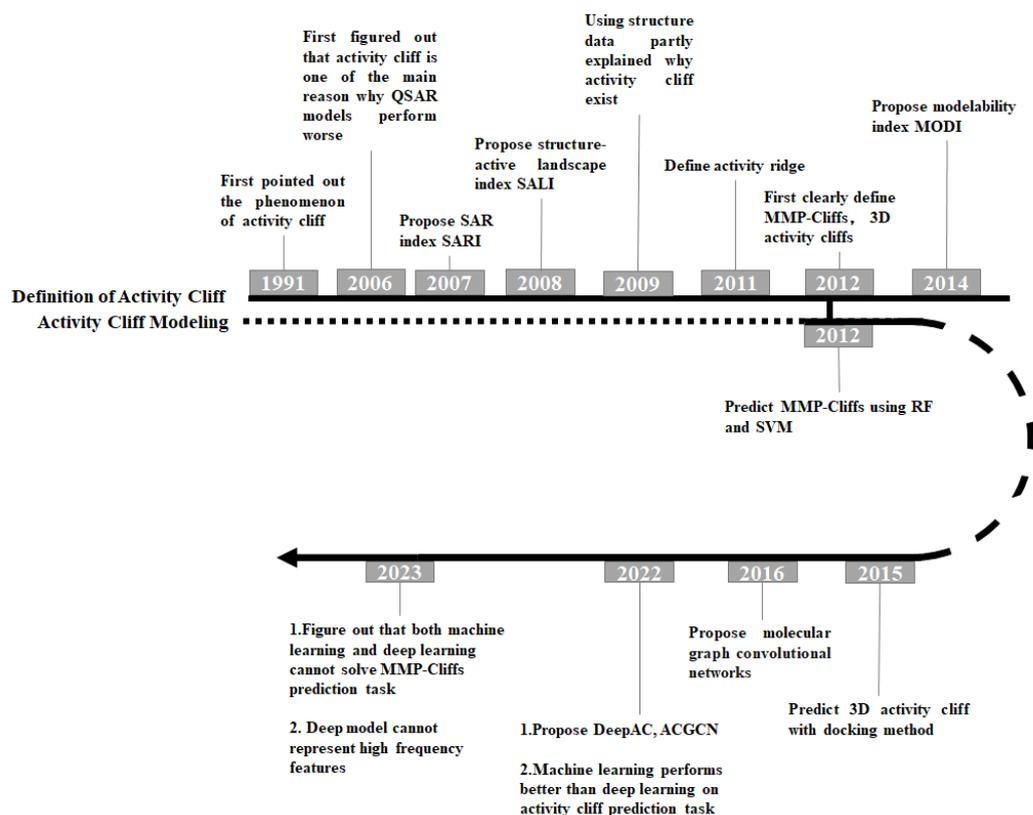


Fig. 1 A brief overview of the development of the activity cliff for molecular prediction task.

As is shown in **Fig. 1**. Before 2016, most research focused on analyzing the causes of activity cliff phenomena from the data itself. Since then, numerous prediction models for MMP-Cliffs, such as those utilizing random forests²² and support vector machines²³, have been developed^{5,15}. Recently, deep learning has been applied in this field, but has often failed to predict MMP-Cliffs accurately^{5,12} due to issues like feature fusion and a tendency to learn low-frequency signals, making it difficult to capture more critical high-frequency signals^{7,24-27}. On the computational view, activity cliffs are regarded as high-frequency signals. Enhanced methods combining machine learning with stereochemistry-aware features¹² like Extended Connectivity Fingerprints (ECFP) have proven more effective in characterizing MMP-Cliffs^{12,16}. Additionally, recent research has revisited QSAR models from an algorithmic perspective, identifying fundamental issues with deep learning in predicting small molecule properties beyond just the lack of big data⁷.

As far as we know, there was no work directly mentioned there were any activity cliffs in the antimicrobial peptides. However, we found out a similar phenomenon on the MIC prediction of AMP. So, we investigated the evolutionary of activity cliff conception in the molecular prediction and drug design, figured out a limitation in the definition of MMP-Cliff, as it does not account for the structural similarity of substructures undergoing transformations within a pair of molecules, nor does it consider the evolutionary information developed over time. Typically, substructures that exhibit structural similarities are likely to possess analogous functional characteristics. It was understandable since many molecules were synthesized artificially, which required a longer period of data accumulation. Luckily, many antimicrobial peptides in the published datasets were composed by canonical peptides, which made it possible to take it into account.

This paper comprehensively defined the Activity Cliff Phenomenon in antimicrobial peptides (AMPCliff), systematically evaluated the recent algorithms on AMPCliff prediction. Our contributions were listed as follows:

- **Discovery of AMPCliff:** The paper identifies the Activity Cliff Phenomenon in antimicrobial peptides (AMPCliff), characterized by structurally similar compounds with significant differences in their Minimum Inhibitory Concentrations (MIC) against the same bacteria.
- **Definition of AMPCliff:** A strict procedure for defining AMPCliff is established, utilizing the public AMP dataset GRAMPA, and illustrated using *S.aureus* as an example. The study also introduces a novel data partitioning method, AC-split.
- **Comprehensive Algorithm Evaluation:** The paper systematically assesses the performances of various machine learning, deep learning methods and LMs including masked and generative language models on the AMPCliff prediction task, revealing that the models can learn underlying patterns from low-frequency signals to predict high-frequency ones, with ESM2 showing superior performance in specific assessments.

Materials and Methods

Benchmark Dataset

It was important that building a QSAR model for antimicrobial peptides design, however, the difficulty of obtaining the MIC values on current public datasets like APD3, DRAMP, DBAASP, YADAMP etc. made it hard to build a trustworthy regressive model. Jacob Witten et al. successfully collected MIC values of these public datasets, and the data was script from Spring 2018²⁸. To the best of our knowledge, GRAMPA²⁸ was the first and only public dataset which can easily download MIC value of each peptide. Some works have used it for peptide design. Junjie Huang et al. directly use GRAMPA to build an RNN-based regressive model to filter high activity peptides¹⁸. Amir Pandi et al. also filtered some sequences from GRAMPA to build RNN-based and CNN-based regressive models¹⁹.

GRAMPA contains 6760 unique sequences, and 51345 total MIC measurements. Some peptide/bacteria pairs occurred multiple times due to overlap between databases and/or activity tested against multiple bacterial strains²⁸. In order to remove all potential peptides containing disulfide bonds, we remove any sequence containing Cysteine. Finally we got 35862 sequences, and we followed²⁸ to take the geometric mean when multiple measurements for a bacterium-AMP pair were present in the database. Our preprocessed data contained 3759 peptides with associated log MIC against *E.coli*, and 3373 peptides with associated log MIC against *S.aureus*. These two bacteria were the richest ones in the database. Since *S.aureus* was one of the pathogens that pose the greatest threat to human life due to the rapid increase in AMP resistance^{29,30}, this paper we used *S.aureus* as an example to discuss the AMPCliffs. Since the sequence length ranged from 7 to 25 occupied over 80% of the total data (see **Fig. 15**), we zoomed into peptides whose length range from 7 to 25 with associated log MIC against *S.aureus* for further analyzation.

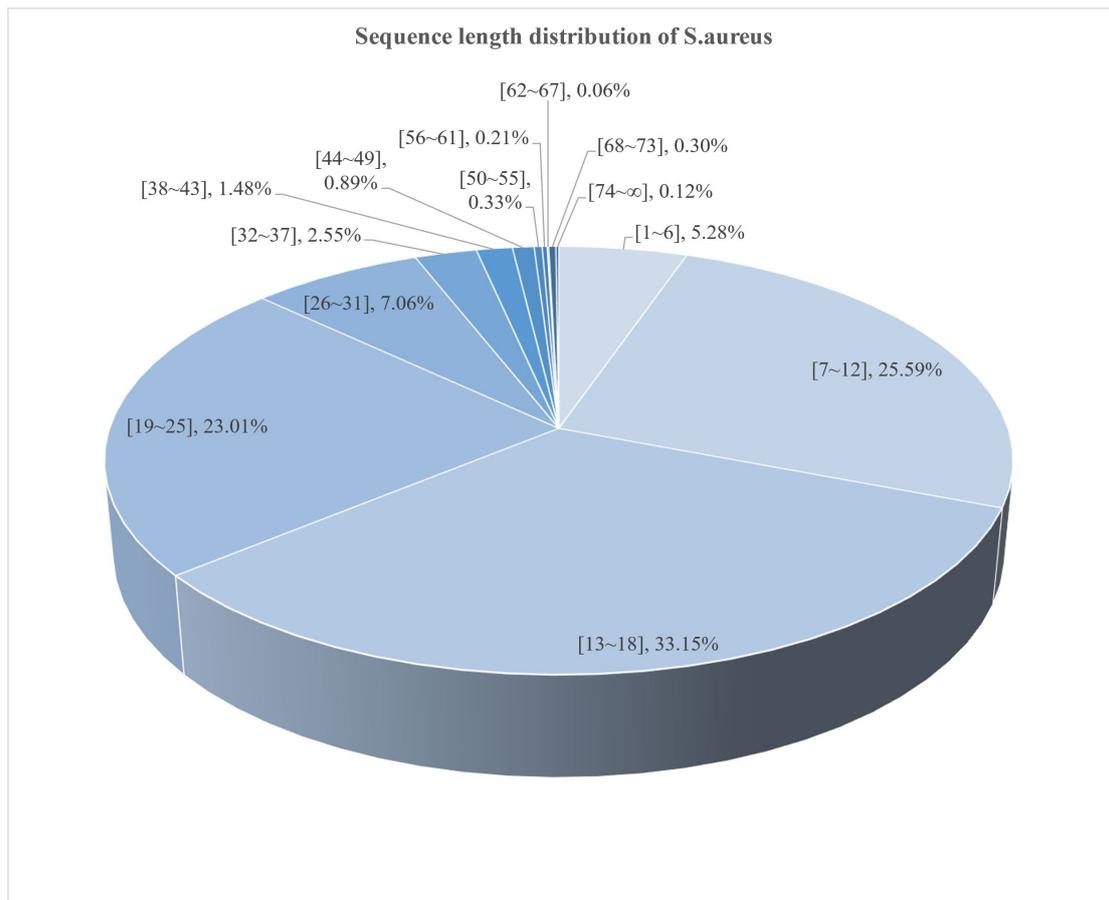


Fig. 2 the sequence length distribution of peptides with associated log MIC *S.aureus*. “[a~b], c%” denotes the percentage of peptides with length ranged from a to b was c%. Sequence length ranged from 13 to 18 got the highest 33.15%, 7-12 get the second 25.59%, 19-25 got the third 23.01%. The other regions were lower than 10%.

Feature Representations

Researchers categorize the input features of antimicrobial peptides sequences into two types: peptide level features (PELF) and amino acid level features (AALF)³¹. An overall overview of these two types of features is shown in **Fig. 16**.

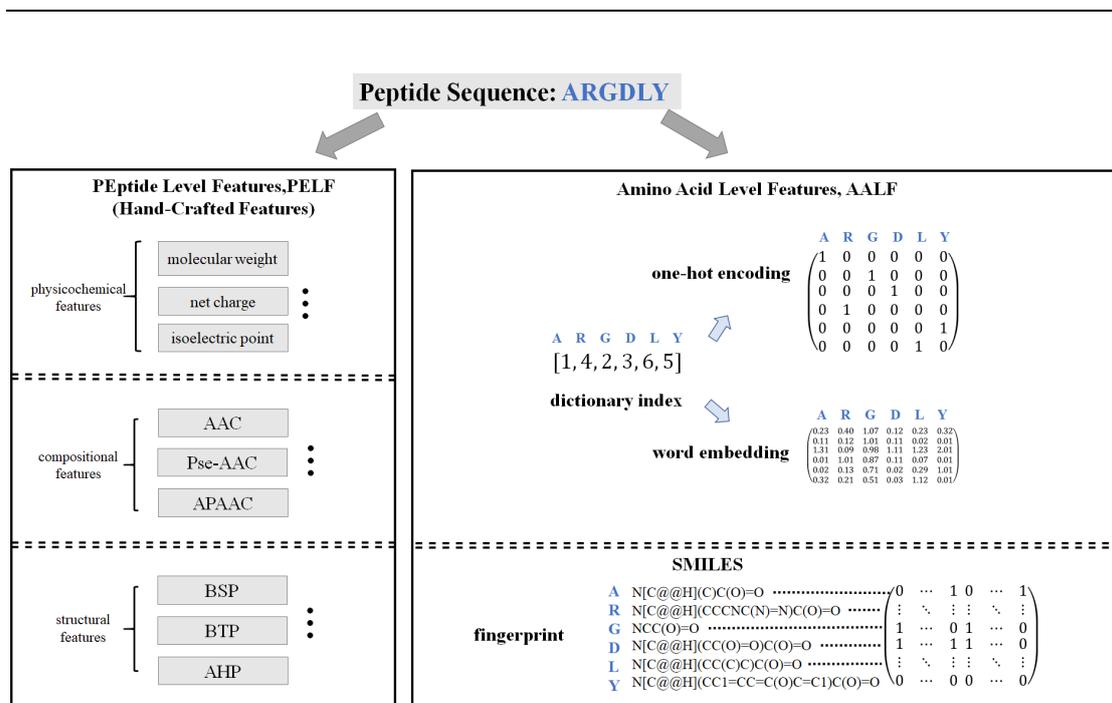


Fig. 3 An overview of two types of features PELF and AALF. In this paper, RELF is also called hand-crafted features.

Hand-crafted features

As Fig. 16 illustrates, PELF consists of (i) physicochemical properties, (ii) compositional components, and (iii) structural features³¹. Each feature category has its distinct calculation principles, focusing on the physicochemical properties of individual amino acids, mathematical statistics of amino acid sequences, and structure-related amino acid features derived from extensive structural database statistics. Next, the details of these three types of features will be introduced:

(i) physicochemical properties molecular weight include the molecular weight, net charge, isoelectric point, etc. of the peptide³². The above characteristics are calculated as follows:

1. Molecular weight is calculated by summing the atomic masses of all amino acids (considering water molecule loss due to peptide bond formation).
2. net charge at a specific pH, calculate the sum of the charges of all ionizable side chains and the N-terminal and C-terminal.

-
3. isoelectric point (pH where the peptide carries no net charge), This can be estimated using an algorithm that takes into account the pKa value of the ionizable group in the peptide.

(ii) Compositional features involve amino acid composition (AAC), pseudo AAC (Pse-AAC), and amphipathic Pse-AAC (APAAC), etc. of peptides³³. The above characteristics are calculated as follows:

1. AAC: The frequency or proportion of each amino acid type in the peptide sequence.
2. Pse-AAC: An extension of AAC that incorporates sequence order information by adding a set of correlation factors that describe the pattern of amino acids in the sequence.
3. APAAC: Similar to pseudo-AAC, but also includes the amphiphilic nature of the residues, which is critical for the interaction of AMPs with biological membranes.

(iii) Lastly, structural features like β -sheet (BSP), β -turn (BTP), and α -helix (AHP) propensities^{34,35} are estimated based on algorithms from structural databases, importance for peptide stability and function, and tendencies of amino acid residues to form specific secondary structures. The above characteristics are calculated as follows:

1. BSP: Use algorithms based on known structural databases to estimate the likelihood of amino acid residues forming β -sheet conformations.
2. BTP: Use a scale that measures the tendency of a sequence to form β -turns, which is important for the structural stability and function of the peptide.
3. AHP: Calculated by assessing the propensity of amino acid residues to become part of an α -helical structure based on empirical data from protein structure databases.

The features mentioned can be found detailed in literature^{19,36} and are part of traditional handcrafted features that, in some classification problems, can perform even better than deep learning-based sequence representation methods³⁷. Combining these with deep learning sequence representations can yield complementary effects³⁷. Based on these

experimental findings, traditional handcrafted sequence features continue to be used in predicting functional peptide.

Amino-acid features

As illustrated in **Fig. 16**, amino acid level features are categorized into four types: (i) dictionary index¹⁸, (ii) one-hot encoding¹⁹, (iii) word embedding³⁸, and (iv) molecular fingerprints³⁹. Dictionary index and one-hot encoding evolve similarly to text representation in natural language processing, while molecular fingerprints, derived from small molecule representation methods, follow a distinct path. Molecular fingerprints initially evolved from the molecule representation method. Molecules have smaller molecular weight than polypeptides¹. we briefly introduce the basic characteristics of these four types of features:

(i) dictionary index is to number each different amino acid in order starting from 1.

(ii) One-hot encoding is an advanced version of dictionary numbering, changing each amino acid feature from a numerical number into a vector consisting of only 0,1. For example, if the dictionary size is 33, then the feature vector dimension of each amino acid is 33. Only at the corresponding numbered position, the feature value in this dimension is 1, and other positions are 0.

However, the underlying assumption of one-hot encoding is that two different amino acids are independent of each other, and the correlation of physical and chemical properties between amino acids is not considered. Like the dictionary number, it is equivalent to the model simply remembering the frequency of occurrence of amino acids. Moreover, one-hot encoding has an important flaw in that it cannot characterize amino acids that are out-of-distribution. For example, non-canonical amino acids. If non-canonical amino acids that have not been seen in the training set appear in the test set, then one-hot encoding will not be able to represent them.

Considering that the mutual independence assumption of one-hot encoding is too strong, researchers began to try to use deep learning methods to let the model learn the universal representation of amino acids by itself, that is, (iii) the word embedding

representation method. That is, each amino acid is mapped into a high-dimensional vector. This feature construction idea originated from natural language processing. With the introduction of the BERT and GPT series of natural language models, research began to regard protein sequences as the "language of life" and use natural language models for modeling^{40,41}.

The most important flaw of the word embedding method is the same as that of one-hot encoding, which is that it cannot represent amino acids that are out-of-distribution. However, for protein and peptide sequences composed of natural amino acids, since natural amino acids are only composed of 20 kinds of amino acids, even if the word embedding method has such defects, it will not affect the modeling of natural amino acids and proteins. The antimicrobial peptides studied in this article are all composed of natural amino acids, so the research in this article does not need to consider this defect.

(iv) The most commonly used calculation method in molecular fingerprints is Extended Connectivity Fingerprint (ECFP). In cheminformatics, a fingerprinting method for representing molecular structures that generates a unique representation of a molecule by taking into account the connectivity of atoms in the molecule and information about neighboring atoms. This type of fingerprint is commonly used in compound similarity searches and drug discovery studies³⁹. This method only needs to know the SMILES expression (Simplified Molecular Input Line Entry System) of molecular characterization to calculate, and is very practical in the fields of small molecule characterization and unnatural amino acid characterization. Although it is not a characterization method specifically designed for natural amino acid sequences, it can also be characterized on natural amino acid sequences. You only need to record the SMILES expression of each amino acid on a certain sequence to calculate the extended connection fingerprint separately, and then splice these fingerprints together along the direction of the amino acid sequence³⁹, as shown in **Fig. 16**.

Because it is calculated based on SMILES expressions and molecular fingerprints, it overcomes the shortcomings of word embeddings in the representation of natural amino acids. For the natural amino acid sequences studied in this article, the relevant pre-training models are mostly based on word embedding methods. This may be because the chemical formula of natural amino acid sequences is much simpler than that of other drug molecules. An amino acid consists of an amino group (-NH₂), a carboxyl group

(-COOH) and a side chain. Different amino acids only differ in the chemical formulas on the side chains.

Results

Case Studies

Activity cliffs in antimicrobial peptides

We discovered a similar activity cliff phenomenon in the benchmark dataset GRAMPA, **Fig. 2** showed some examples of AMPCliffs in the benchmarking dataset GRAMPA. The corresponding sequence alignment scores were also illustrated next to the AMPCliff pairs. It is worthy to know that solely using the score of sequence alignment cannot directly measure the similarity of two sequences, on the one hand the substitution matrix BLOSUM62 isn't designed for measuring similarity, on the other hand the algorithm of the score doesn't take the sequence length into account.

AFHHI FRGIVHVGKTIHRLVTG	$5.5 \times 10^{-6} \text{ mol/L}$	$\updownarrow 10 \text{ fold}$ Smith-Waterman score: 110
FFHHA FRGIVHVGKTIHRLVTG	$50.0 \times 10^{-6} \text{ mol/L}$	
AGRGKQGGKVRAKAKTRSS	$8.2 \times 10^{-6} \text{ mol/L}$	$\updownarrow 7.9 \text{ fold}$ Smith-Waterman score: 89
DGRGKQGGKVRAKAKTRSS	$64.4 \times 10^{-6} \text{ mol/L}$	
AGRGKQGGKVRAKAKTRSS	$8.2 \times 10^{-6} \text{ mol/L}$	$\updownarrow 8.3 \text{ fold}$ Smith-Waterman score: 89
-GRGKQGGKVRAKAKTRSS	$68.4 \times 10^{-6} \text{ mol/L}$	
AGRGKQGGKVRAKAKTRSS	$8.2 \times 10^{-6} \text{ mol/L}$	$\updownarrow 5.2 \text{ fold}$ Smith-Waterman score: 89
KGRGKQGGKVRAKAKTRSS	$43.0 \times 10^{-6} \text{ mol/L}$	

Fig. 4 Examples of AMPCliffs in the benchmarking dataset GRAMPA.

Mutation on conservative region can intricate activity cliff

Additionally, we found out that some substitution of conservative amino acids can also cause activity cliff in antimicrobial peptides. We illustrated 10 sequence pairs as examples in **Table 1**. The complete AMPCliffs can be found in Supplementary Table S6. For example, FLPIIGKLLSGLL and FLPIVGKLLSGLL only had turned I to V on

the 5th position, while the difference of their MIC values increased from $4.45 \times 10^{-6} \text{ mol/L}$ to $59 \times 10^{-6} \text{ mol/L}$, with approximately 13-fold dilution difference.

Table 1 10 examples of AMPCliff pairs on conservative region.

seq1	seq2	MIC 1(M)	MIC 2(M)	B6 2Sc	B62 bit Sc	mutat ion	mutatio n value
FLPIIGKLLSGLL	FLPIVGKLLSGLL	4.45E-06	5.9E-05	61	2.084669	(5, I, V)	3
GLLKKIKWLL	GLLKRIKWLL	2.84E-05	3.23E-06	49	2.275529	(5, K, R)	2
GWLDVAKKIGKA AFNVAKNFI	GWLDVAKKIGKA AFNVAKNFL	2.23E-05	1.7E-06	107	2.083717	(21, I, L)	2
ILPWKWPWWKW RR	LLPWKWPWWKW RR	2.58E-06	1.61E-05	95	3.031742	(1, I, L)	2
INLKAIAALAKKL L	INLKAIAAMAKKL L	2.15E-05	1E-06	59	1.884033	(9, L, M)	2
KKKWLWLW	KRKWLWLW	6.74E-05	8.23E-06	53	3.02547	(2, K, R)	2
KQKWLWLW	KQRWLWLW	3.79E-05	3.7E-06	53	3.02547	(3, K, R)	2
RRLFRRILRWL	RRLFRRILRYL	5.5E-07	4.5E-06	49	2.068663	(10, W, Y)	2
RRWWRWWR	RRWYRWWR	1.8E-06	4.2E-05	55	3.115999	(4, W, Y)	2
RWWRWWR	RWWRYWR	1.4E-05	2.07E-06	50	3.302487	(5, W, Y)	2

B62Sc: BLOSUM62 score of the 2 aligned sequences

B62bitSc: BLOSUM62 bit score per column of the 2 aligned sequences, the parameter settings followed BLASTp⁴²

Mutation: (position, amino acid in seq1, amino acid in seq2)

Problem Formulation

AMPCliff definition

During our experimental investigations, we have identified a limitation in the definition of MMP-Cliff, as it does not account for the structural similarity of substructures undergoing transformations within a pair of molecules, nor does it consider the evolutionary information of amino acids developed over time. Typically, substructures that exhibit structural similarities are likely to possess analogous functional characteristics. Hence, we proposed a new similarity definition of activity cliff in antimicrobial peptides, called AMPCliff.

The entire fundamental principles of AMPCliff consistent with the common sense, that is, a pair of structurally similar compounds that are active against the same bacteria but significantly different in MIC values. In order to specifically define the AMPCliff, we defined the similarity of two sequences as followed:

For a given pair of sequences *sequence1*, *sequence2*, we denoted S_1, S_2 for short. Then we using sequence alignment algorithm to get the aligned sequence pair, denoted as S_{align}^1, S_{align}^2 with the same aligned sequence length L_a . Furthermore, we used a similarity matrix $\mathbf{SIMrx} \in \mathbf{R}^{N \times N}$ to measure the similarity of each 2 amino acids. For any two arbitrary amino acids p, q , we have their similarity value $\mathbf{SIMrx}(p, q) \in [0,1]$, and then we got the similarity of the sequences pair $\frac{1}{L_a} \sum_{i=1}^{L_a} \mathbf{SIMrx}(S_{align}^1(i), S_{align}^2(i))$. **Fig. 3** showed the procedure of AMPCliff definition. Noticed that this paper using negative log MIC (*mol/L*) for definition, model training and result analysis.

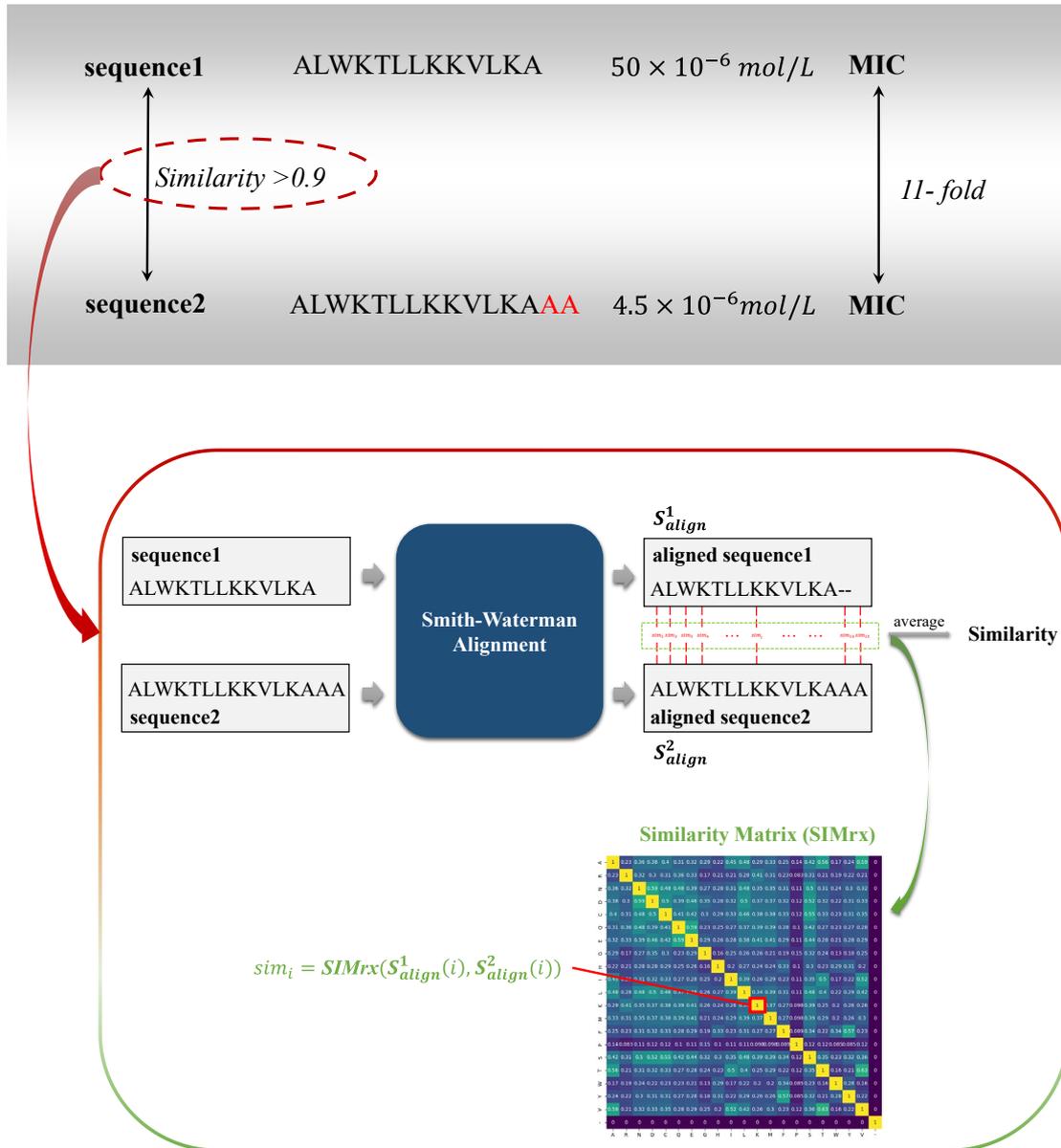


Fig. 5 the framework of identifying AMPCliff. Here we took ALWKTLKKVLKA as *sequence1*, ALWKTLKKVLKAAA as *sequence2* to illustrate the procedure. The main difference between AMPCliff and MMP-Cliff was the similarity measurement. AMPCliff firstly using sequence alignment to align 2 sequences. this paper we use Smith-Waterman alignment. Then using a similarity matrix to calculate the similarity of the amino acids from the two sequences at each position, using geometric mean value to get the final similarity. This figure showed the MIC difference of *sequence1* and *sequence2* was around 11-fold, larger than 2-fold to treat *sequence1* and *sequence2* as a pair of AMPCliffs. 2-fold was the minimum MIC dilution difference set in this paper.

Sequence alignment

Global alignment and local alignment are the 2 types of dynamic programming algorithm for sequence alignment. A global alignment is one that compares the two sequences over their entire lengths, and is appropriate for comparing sequences that are expected to share similarity over the whole length⁴³. The alignment maximises regions of similarity and minimises gaps using the scoring matrices and gap parameters provided to the program. While a local alignment searches for regions of local similarity and need not include the entire length of the sequences⁴³. Peptide sequences are much shorter than protein sequences, usually smaller than 50 amino acids. In this paper we used 7-25 length sequences for analysis, therefore it needed reconsidering which type of alignment was appropriate.

Felix Teufel et al. proposed GraphPart⁴⁴, which utilized global alignment, Needleman-Wunsch alignment, by EMBOSS module *needle* with default parameters. Whereas Javier L. Baylon et al.⁴⁵ established PepSeA, and they systematically evaluated the performance of global alignment and local alignment on therapeutic peptides composed by both canonical only and non-canonical amino acids, and they figured out that local alignment was more accurate. We found out that GraphPart truncated protein datasets SignalP 5.0⁴⁶ to obtain signal peptides with 70 N-terminal amino acids, and NetGPI⁴⁷ to obtain GPI anchors with 100 C-terminal amino acids. These peptides' lengths are fixed and longer than what we use in this paper (range from 7 to 25). Whereas the sequence length distribution in PepSeA⁴⁵ was much closer to our research, which was ranged from 9-18 in total. Therefore, we followed the work PepSeA⁴⁵ using local alignment algorithm. Since we only investigate AMPCliff on peptides composed by natural amino acids, we used Smith-Waterman alignment with BLOSUM62 substitution matrix and set penalty of gap open and gap extension -11 and -1, which are default parameters in BLASTp⁴².

Performance metrics

For the model evaluation, thoughtful consideration of statistical analysis, evaluation metrics, and task settings is critical as they impact the observed prediction performance⁶. Here we chose Spearman and recall as the evaluation metrics, since we cared more about the relative order of the predicted value than the specific prediction value of the

model. Besides, we calculated recall involved sorting both top 50 of the label values and the predicted values in descending order (since we used negative \log (MIC (mol/L)) for analysis, the larger the value, the better antimicrobial property peptide had.), followed by determining the count of their intersection.

Dataset Partition

Dataset Partition or data split problem is essential before developing a model for a specific property prediction task⁶. ACNet¹⁶ utilized random split and target split to evaluate the model performance, but random split merged the activity cliffs into training and test set, reducing the influence of the activity cliffs on the model performance and overestimated the capacity of the model on activity cliff prediction tasks. While target split introduced out-of-distribution (OOD) issue into the model and amplified the influence of activity cliffs on the model performance. So, we didn't use either of them for data partition. Instead, we kept the sequence feature distributions of the training set, validation set, and test set on the same domain. Furthermore, the key difference between test set and training, valid dataset was whether it contain AMPCliff pairs. We utilized the sequences without AMPCliff pairs for model training and validation, then test the model on the test set composed by AMPCliffs. We aimed to investigate whether the model can learn any pattern from the low-frequency data, which can transform to the high-frequency data (namely AMPCliffs), and the proposed data partition strategy called AC-split. For the criteria on splitting training and valid dataset, since we just used the valid dataset to find the best parameters during training, we used random split which consistent with the literature³⁸. The entire experiments in this paper used AC-split by default. **Fig. 4** gave a flowchart of AC-split for better illustration.

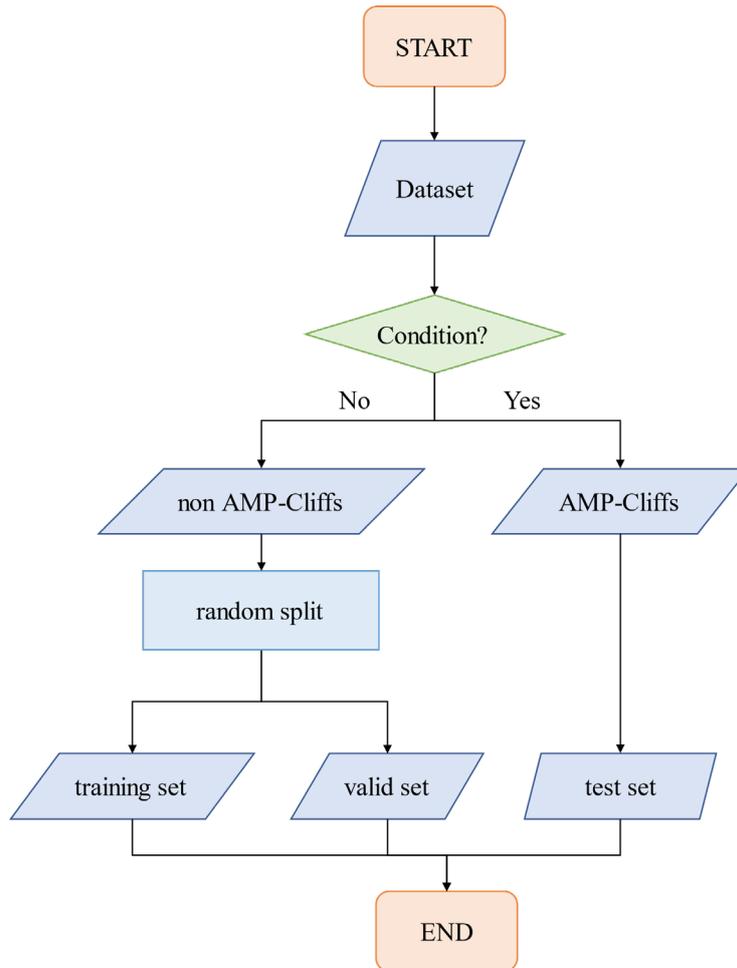


Fig. 6 the flowchart of data partition method AC-split. The decision node “Condition” was “BLOSUM62 average” or “Tanimoto average” in this paper. The details of the decision procedures please see the *definitions of AMPCliff* section.

In order to evaluate our proposed data split strategy whether eliminated the influence of OOD problem, we took ESM2 with 33 layers as an example to show the features of training, valid and test set by TSNE. The results of splitting test set by 2-fold dilution difference under each condition were shown in **Fig. 5**. Either under “Tanimoto average” condition or “BLOSUM62 average” condition, the samples from the three datasets were mixed together. It guaranteed the availability of our proposed dataset partition strategy.

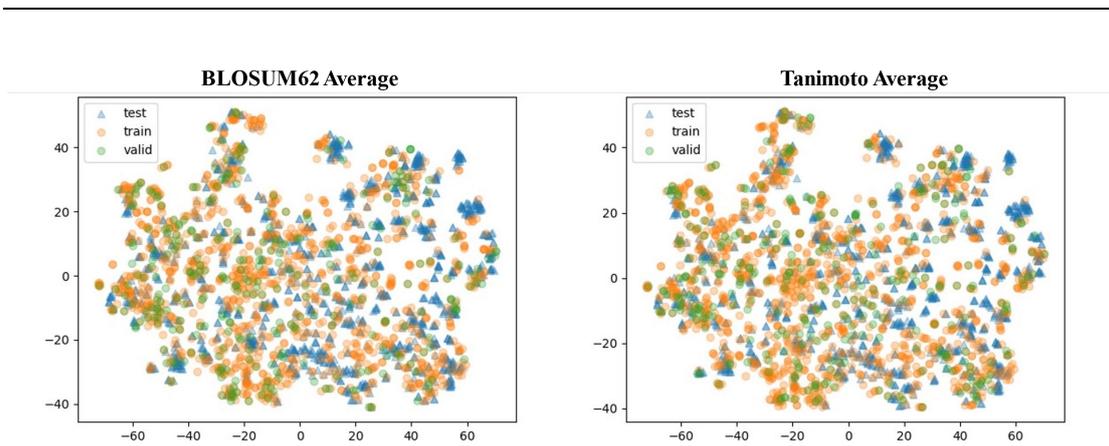


Fig. 7 TSNE visualization of EMS2 with 33 layers features for training, valid, and test datasets.

Similarity Measurement Discussion

BLOSUM62 vs Tanimoto similarity

In order to compare substitution matrix BLOSUM62(int values range from -4 to 11) and Tanimoto similarity of each amino acid with ECFP (float values range from 0 to 1) on the same scaling level, we used a normalization strategy similar to *MinMaxScaler* in the *sklearn* Python package but added another process for ensuring the symmetric of the new matrix, that is:

$$M'_{new} = \frac{M - \min(M)}{\max(M) - \min(M)} \quad (1)$$

$$M_{new} = \frac{M'_{new} + M'_{new}{}^T}{2} \quad (2)$$

Here M is the original matrix, M_{new} is the normalized matrix. **Fig. 6** shows the details of normalized BLOSUM62 metrics and Tanimoto similarity metrics. It was found that there were similarities in certain amino acid pairs. For instance, the normalized BLOSUM62 values for Phenylalanine(F) with Tryptophan(W) and Tyrosine(Y) are 0.42 and 0.65, and Phenylalanine(F) is most similar to Tyrosine(Y) when compared with all other amino acids; this was also reflected in the Tanimoto

similarity, where the similarity of Phenylalanine(F) with Tryptophan(W) and Tyrosine(Y) is 0.5 and 0.7, respectively, with Phenylalanine(F) being most similar to Tyrosine(Y) among all other amino acids.

However, there are many differences as well. For example, the normalized BLOSUM62 values between Cysteine(C) and Arginine(R), Asparagine(N), Aspartic(D) are all less than 0.1, approaching 0; whereas the Tanimoto similarities between Cysteine(C) and Arginine(R), Asparagine(N), Aspartic(D) are 0.3, 0.5, 0.52 respectively. When looking at the similarity between Cysteine(C) and other amino acids, these are relatively high values. The fundamental reason for this phenomenon is that the calculation of Tanimoto similarity does not take the structural similarity into account at the atomic level; it crudely assumes that the substitution distance between all atoms is the same.

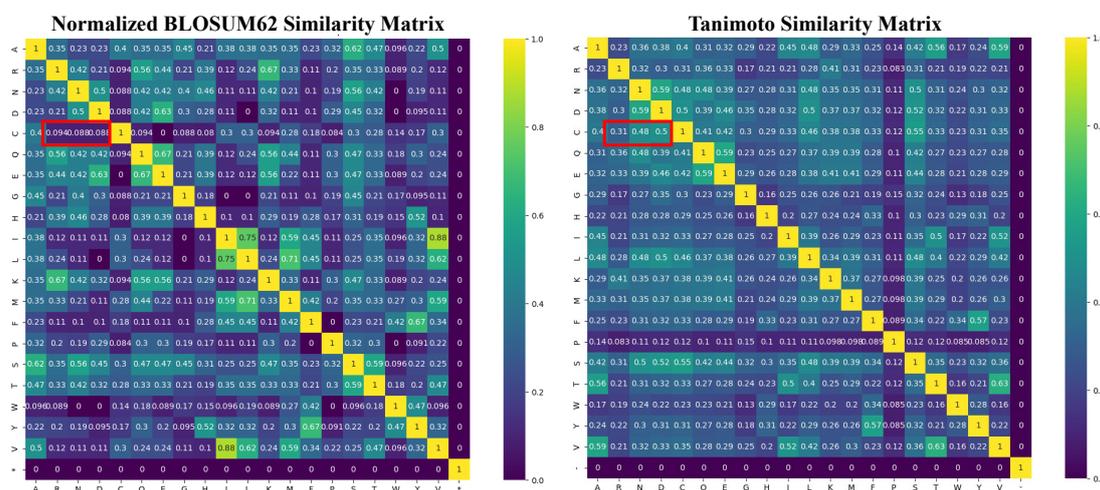


Fig. 8 The details of 2 normalized metrics: normalized BLOSUM62 and normalized Tanimoto similarity

As is shown in Fig. 7, Cysteine (C) is an amino acid containing sulfur (S), which typically tends to form disulfide bonds. On the other hand, Arginine (R), Asparagine (N), and Aspartic acid (D) have quite different side chains and chemical properties: Arginine (R) has an amino group (-NH₂) and is basic, Asparagine (N) contains an amide group (-CONH₂) and is amidic, while Aspartic acid (D) includes a carboxyl group (-COOH) and is acidic. These fundamental differences in their chemical nature result in distinct physicochemical properties. Whereas Tanimoto similarity failed to measure these functional differences among the amino acids. The definition of MMP-

Cliff has similar problem as Tanimoto similarity, they all regarded the atom-level difference the same.

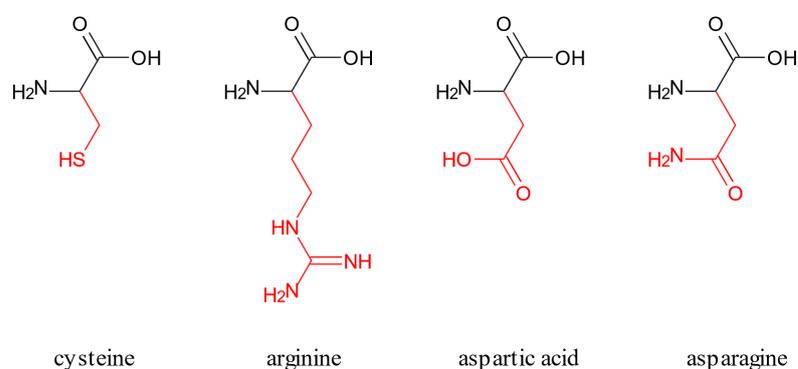


Fig. 9 The structures of cysteine(C), arginine(R), aspartic acid(D) and asparagine(N)

MMseqs2 is inappropriate for measuring the similarity of peptides

MMseqs2⁴⁸ is a popular clustering method on the worldwide large protein datasets UniProt⁴⁹. It clustered UniProt into three type of datasets UniRef100, UniRef90 and UniRef50 by setting the similarity threshold of sequence identity in 100, 90 and 50. Plenty of protein design work has used sequence identity to illustrate the novelty of their discovered proteins^{50,51}. So, we naturally thought about whether we could use sequence identity to measure the similarity between peptide sequences. After searching on the official UniProt website⁵², we surprisingly found that sequences shorter than 11 residues are excluded in UniRef90 and UniRef50. However, in the rational peptide drug design, peptides in 10-13 length are quite in common. Many computational algorithms were built upon the data partition by UniRef or MMseqs2, therefore, it is urgent to figure out the reason why MMseqs2 has poor performance in short peptides and design a new clustering method on short peptides.

Some researchers have figured out that MMseqs2 performed worse on short peptides^{44,53} and proposed some new cluster algorithms like GraphPart, GibbsCluster for short peptides dataset partition problem^{44,54}. But MMseqs2 is a comprehensive and customizable tool, none of them have pointed out a certain reason why MMseqs2 is incapable on the short peptide clustering problem. GraphPart figured out that MMseqs2 can be used for homology reduction by selecting only the representative sequences, but should not be used for homology partitioning, since the maximum similarity criterion

is only applied to the representative sequences, not to the other members of the clusters⁴⁴. Maria Hauser et al. have figured out that SWIPE is the most sensitive tool on the short peptides in 50 lengths among SWIPE, BLAST, MMseqs-sens, MMseqs-fast, DIAMOND-sens, RAPSearch2, UBLAST and DIAMOND⁵³. Whereas GibbsCluster is a clustering method especially to deconvolute multiple specificities in major histocompatibility complex (MHC) class I peptidome datasets⁵⁴. All of these work hard to tell the fundamental reason why MMseqs2 performs worse on short peptides. Although GraphPart have pointed out a potential reason, MMseqs2 had a *cluster* module for clustering indeed.

In order to clarify the vague cognition on the performance of MMseqs2 on short peptides, we dig deeply into the mathematical calculation of MMseqs2. It turns out that there are two main drawbacks of MMseqs2 on short peptides. The default module of MMseqs2 calculates *kmer* before the sequence alignment. It is a trade-off strategy which scarifies precision for speed. This is reasonable for protein clustering, since proteins are mostly long sequences, whereas short peptides don't need this trade-off, which is one of the key differences between short peptides and protein. Another calculation difference is sequence identity. MMseqs2 calculate sequence identity in 2 different ways⁵⁵: (i) If `--alignment-mode 3`, compute the number of identical aligned residues divided by the number of aligned columns including columns containing a gap in either sequence. (ii) Otherwise, the score per column by default, i.e., the local alignment bit score divided by the maximum length of the two aligned sequence segments. As is shown in the following:

$$\begin{aligned}
 & \text{sequence identity}(S_{align}^1, S_{align}^2) = \\
 & \left\{ \begin{array}{l} \frac{\sum_{i=1}^L \{1 | S_{align}^1[i] = S_{align}^2[i] \wedge S_{align}^1[i] \neq \text{gap} \wedge S_{align}^2[i] \neq \text{gap}\}}{L} \\ \frac{\text{bit}_{score}(S_{align}^1, S_{align}^2)}{\max(L_1, L_2)} \end{array} \right.
 \end{aligned}$$

3)

However, neither of them is directly suitable for measuring the similarity of short peptides. For the first calculation method, the result was sensitive to sequence length. For instance, if two aligned peptides both have length 8, they will get sequence identity=0.875 with only 1 residue difference. Then we set the similarity threshold to 0.9, these short peptides will be wrongly filtered. For the second measuring method, it

was more appropriate than the former one, since it also takes the degree of similarity between aligned amino acids and the number and length of gaps into account⁵⁵. However, the value of score per column doesn't range from 0 to 1, which make it hard to set a threshold. Therefore, we combined the benefits of both ways: we normalized the BOLSUM62 matrix by Eq.(1), Eq.(2) and calculate the average similarity of each column in the two aligned peptide sequences, which we called "BLOSUM62 average".

We also set an experiment to proof the former analyzation. We compared the number of AMPCliffs in different sequence length region measured by Levenshtein distance (Levenshtein), Levenshtein distance of aligned sequences (Levenshtein aligned), the first sequence identity method (sequence identity), the average similarity of each column in the two aligned peptide sequences (BLOSUM62 average) and the average Tanimoto similarity of each column in the two aligned peptide sequences (Tanimoto average). **Fig. 8** shows the number of matched AMPCliffs in different length distribution.

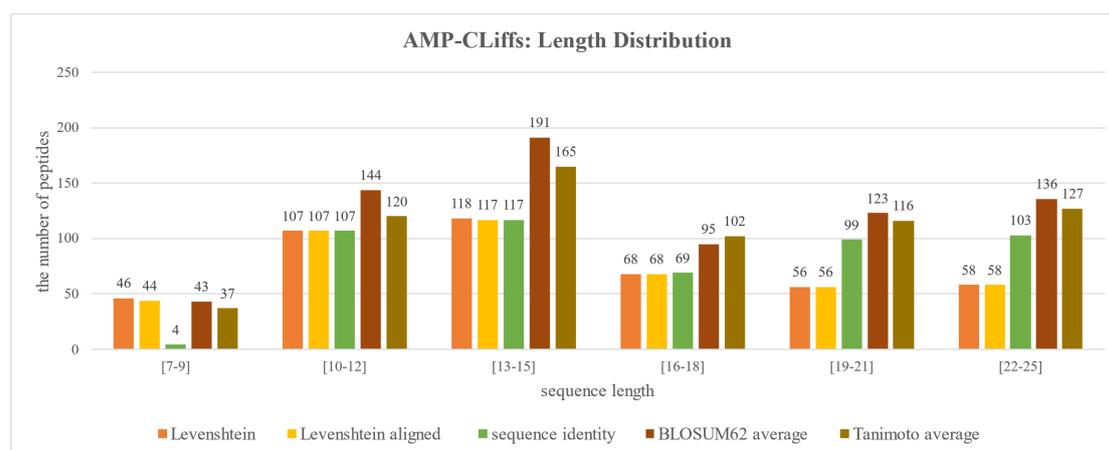


Fig. 10 the number of matched AMPCliffs in different length distribution. We compared the number of sequences measured by Levenshtein distance (Levenshtein), Levenshtein distance of aligned sequences (Levenshtein aligned), the first sequence identity method (sequence identity), the average similarity of each column in the two aligned peptide sequences (BLOSUM62 average) and the average Tanimoto similarity of each column in the two aligned peptide sequences (Tanimoto average). Here we set 1 for "Levenshtein" and "Levenshtein aligned", the other 3 methods with threshold 0.9.

The difference between “Levenshtein” and “Levenshtein aligned” lies in whether the sequences are aligned or not. There is minimal difference between the two in identifying AMPCliffs;

The distinction between “Levenshtein aligned” and “sequence identity” comes down to whether an average over the length is taken. It is precisely this operation that makes sequence identity more tend to identifying longer sequences. Consider two aligned sequences, S_{align}^1 and S_{align}^2 , differ by only one amino acid position, resulting in a numerator of 1. Given our similarity threshold of 0.9, which means that the aligned sequences must be at least 10 amino acids long to be detectable. This explains why the AMPCliffs identified by sequence identity tend to be longer and why there are hardly any sequences in the [7-9] range. At the same time, this also suggests that clustering short peptides using the MMSeqs2 algorithm, which is used for clustering in the large protein sequence database UniProt, is not suitable, a viewpoint that coincides with the experimental results in the literature^{52,53}.

In contrast, the distribution of AMPCliff numbers identified by “BLOSUM62 average” and “Tanimoto average” are more similar, with the former generally yielding a higher count, as illustrated in **Fig. 9** below. Each metric has its advantages and disadvantages: BLOSUM62 incorporates the prior knowledge of medicinal chemists and evolutionary information of amino acids but fails to extend to non-natural amino acid peptides; Tanimoto similarity indiscriminately treats different atomic substitution distances the same, losing evolutionary information but able to extend to non-natural amino acid peptides. Hence, this thesis presents two versions of the AMPCliff identification algorithm. For tasks involving natural amino acids, the BLOSUM62 version is recommended; for non-natural amino acid tasks, the Tanimoto version is recommended if a proper sequence alignment algorithm for non-canonical peptides established. As far as we knew, there was only one sequence alignment algorithm PepSeA⁴⁵ available. This paper mainly showed the results of AMPCliff prediction under “BLOSUM62 average” condition. The complete results were shown in Supplementary data.

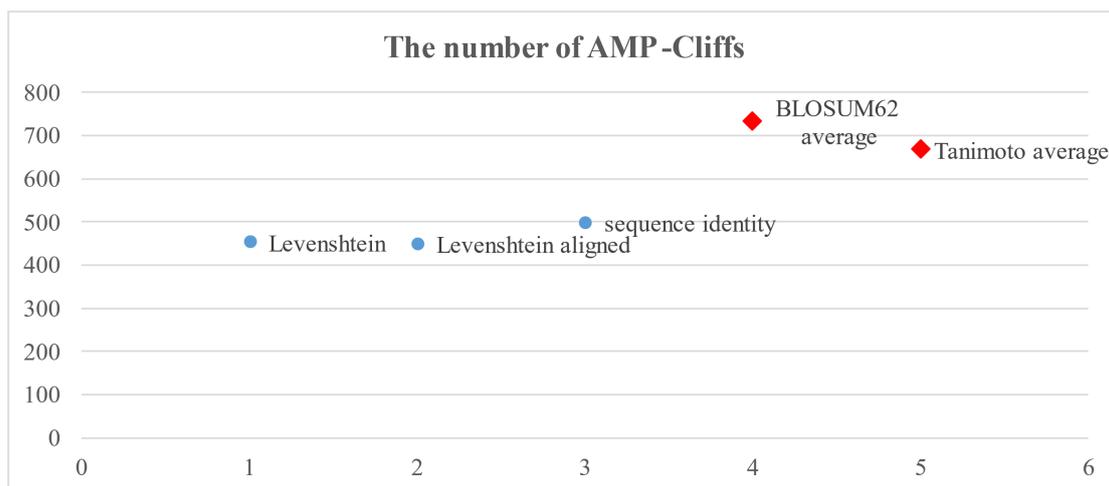


Fig. 11 The number of AMPCliffs detected by the 5 measurements

Quantify the Activity Difference

According to the definition of MMP-Cliff⁹, the potency difference between compounds in an MMP meeting the structural criteria had to be at least 2 orders of magnitude. The measurements of the activity in molecular and antimicrobial peptides are different, antimicrobial peptides usually used MIC value to measure the antimicrobial capacity. Johan W. Mouton et al.⁵⁶ pointed out the acceptance criteria for reproducibility of an MIC device following the ISO 20776-2 norm indicate an acceptable deviation of one dilution from the mode for 95% of cases, or a range of at least two 2-fold dilutions. Thus, we assumed that the minimum difference of the MIC in an AMPCliff pair should at least larger than 2-fold dilution. Hence, after the process of MIC by negative log MIC (*mol/L*), let the MIC (*mol/L*) value of any two arbitrary sequences $MIC_1, MIC_2, MIC_1 \geq MIC_2$ we got:

$$\frac{MIC_1}{MIC_2} \geq diff \quad (4)$$

$$-\log\left(\frac{MIC_1}{MIC_2}\right) = -\log MIC_1 - (-\log MIC_2) \leq -\log(diff) \quad (5)$$

$$|-\log MIC_1 - (-\log MIC_2)| \geq \log(diff) \quad (6)$$

Fig. 10 showed the number of AMPCliff peptides in “Tanimoto average” and “BLOSUM62 average” conditions of different dilution differences. From 2-fold to 5-fold, the data size reduced nearly a half. It illustrated there were quite numbers of AMPCliffs in the public AMP datasets. And **Fig. 11** gave show cases of the AMPCliffs as the dilution order increased.

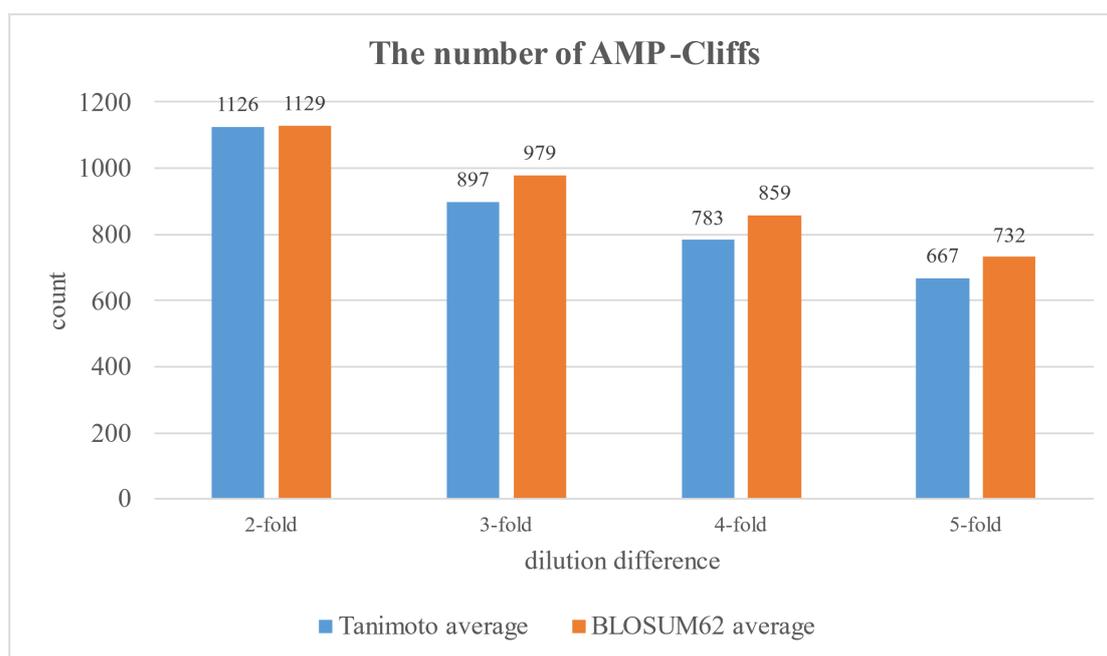


Fig. 12 the number of AMPCliff sequences on each condition of different dilution differences.

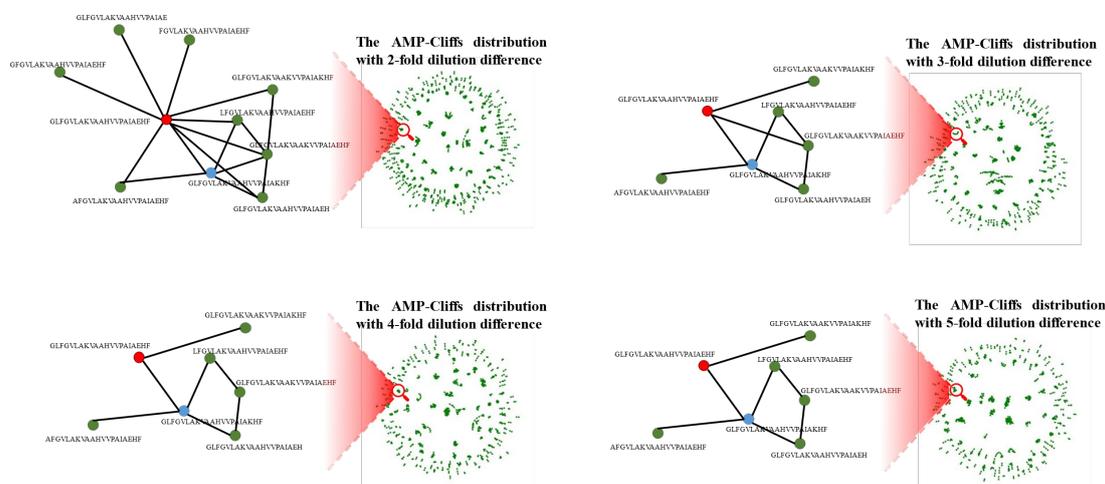


Fig. 13 Show cases of the variation of AMPCliffs as the dilution order increased.

Nodes represented sequences, if 2 nodes were marked as an AMPCliff pair, then they got an edge. The red node and the blue node represented the top 2 maximum degree nodes as dilution order increased.

Benchmarks

Experiment Settings

This paper both used hand-crafted features and amino-acid features. The hand-crafted features were used for machine learning methods, while the amino-acid features were used for deep learning models and LMs.

We followed the literature¹⁸ to construct hand-crafted features. The physicochemical properties contained molecular weight, charge characteristics, content of polar and nonpolar amino acids, hydrophobicity, and van der Waals volume, etc. The compositional features involve amino acid composition (AAC), dipeptide composition (DiC), pseudo AAC (Pse-AAC), and amphipathic Pse-AAC (APAAC), Moreau-Broto Autocorrelation, Quasi-Sequence Order and so on. And we additionally added a structural feature α -helix propensity.

For amino-acid features we used features including dictionary index, one-hot encoding, fingerprints and word embeddings. Except for word embeddings, the other three methods were all used for peptide design^{18,19,39}, so we kept them for comparison.

we listed the former mentioned features and their models in **Table 2**. If the corresponding methods were derived from literature, the column “citation” will list its abbreviation in this paper and give the original literature. The hyperparameters of each model were listed in Supplementary Table S1.

Table 2. the features and models used in this paper. If the corresponding methods were derived from literature, the column “citation” will list its abbreviation in this paper and give the original literature. Otherwise denote “this paper”.

category	Feature representation	models	citation
Machine Learning (ML)	Hand-crafted features	support vector machine (SVM), linear regression (LR), L1, L2, ElasticNet, random forest (RF), gradient boost regressor (GB), XGBoost, gaussian process regressor (GP)	This paper
Deep Learning (DL)	dictionary index	LSTM	AMPSpace ¹⁸
Deep Learning (DL)	One-hot encoding	LSTM	CellFree ¹⁹
Deep Learning (DL)	One-hot encoding	CNN	CellFree ¹⁹
Deep Learning (DL)	fingerprint	CNN	peptimizer ³⁹
Language Model (LM)	Word embedding	BERT, ESM2 with 6,12,33 layers, GPT2, ProGen2 small, base and medium version.	This paper

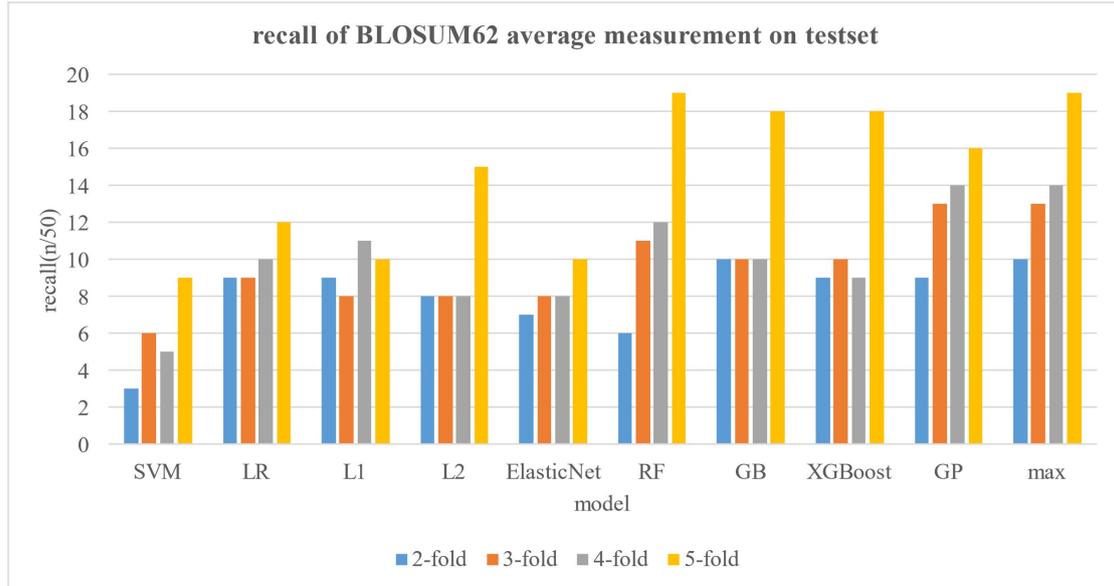
Performance of machine learning

Then we testify the capacity of machine learning methods on predicting AMPCliffs in 2-fold, 3-fold, 4-fold and 5-fold dilution differences. We used hand-crafted features with SVM, LR, L1, L2, ElasticNet, RF, GB, XGBoost and GP regressive models separately to train QSAR models in order to measure the capacity of these models on predicting AMPCliffs under each condition. **Fig. 12** showed the results of AMPCliffs under each condition. The performance metrics are recall and Spearman coefficient.

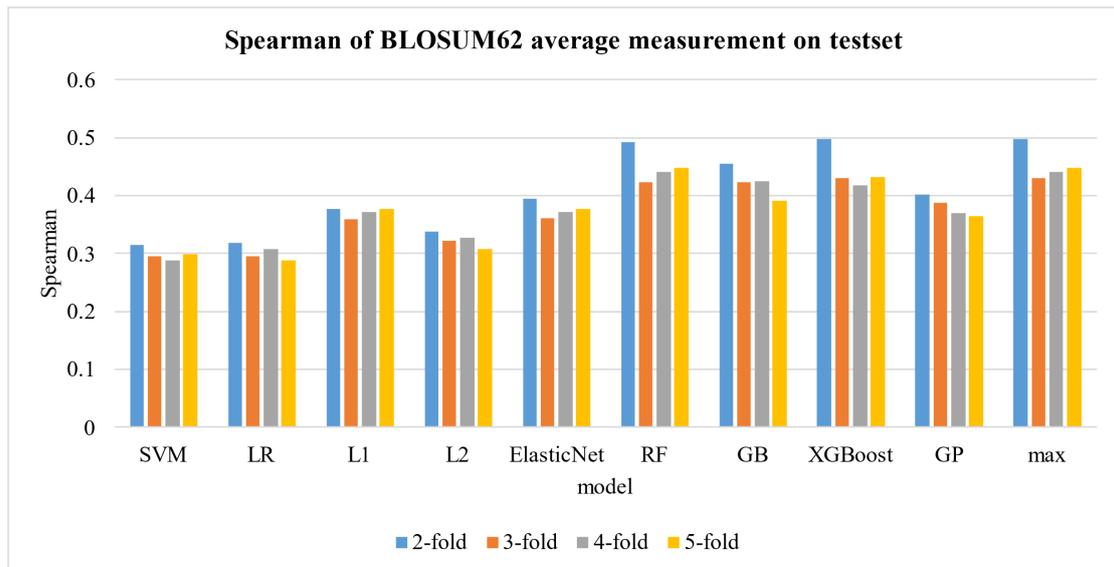
Fig. 12(a) and (c) showed the results of AMPCliffs' prediction split by "Tanimoto average" condition, (b) and (d) showed the results of AMPCliffs' prediction split by "BLOSUM62 average" condition. Both of them were poor at predicting AMPCliffs in 2-fold, needless to say 5-fold.

Besides, generally tree-based models such as GB, RF, XGBoost performed better than the others. RF got the best spearman under the "BLOSUM62 average" condition with 4-fold and 5-fold dilution difference and under the "Tanimoto average" condition with 2-fold and 4-fold dilution difference. XGBoost got the best spearman under the "BLOSUM62 average" condition with 2-fold and 3-fold dilution difference. And GB got the best spearman matrix under the "Tanimoto average" condition with 3-fold and 5-fold dilution difference. Whereas RF and GB got the best recall under the "BLOSUM62 average" condition with 5-fold dilution and 2-fold dilution difference respectively. And XGBoost got the best recall under "Tanimoto average" condition with 2-fold, 4-fold and 5-fold dilution differences, while GP got the best recall under the "BLOSUM62 average" condition with 3-fold, 4-fold dilution difference and under "Tanimoto average" condition with 3-fold dilution difference.

When we gradually add some lower frequency signals like no bigger than 3-fold, 4-fold AMPCliffs into the training set, although the spearman decreased, the recall raised gradually. Although Spearman didn't follow the same trend as recall, it was worthy to know that except for the 2-fold dilution, the others had a similar trend. It indicated that models can learned some potential common patterns from the low-frequency signals to successfully predict the high-frequency signals. The details of results for each model can be found in supplementary Table S2.



(a)



(b)

Fig. 14 the recall and Spearman value of AMPCliffs prediction under Tanimoto and BLOSUM62 average condition. The x-axis showed the names of regressive models, each color of bar represented the activity difference by 2-fold(diff2), 3-fold(diff3), 4-fold(diff4) and 5-fold(diff5) respectively. Recall was calculated by counting the number of true predicted samples in the top 50 orders. (a) and (b) showed the results of AMPCliffs prediction split by “BLOSUM62 average” condition. Both of them were poor at predicting AMPCliffs in 2-fold (with Spearman less than 0.5),

needless to say 5-fold. However, when we gradually add some lower frequency signals like no bigger than 3-fold, 4-fold AMPCliffs into the training set, although the Spearman of test set decreased, the recall raised gradually.

Performance of deep learning

Here we select 3 mentioned tree-based models RF, GB, XGBoost models and GP models for the comparison with deep learning models. There were 4 deep learning regressive models which were established for peptide design^{18,19,39}. And we used the default parameters set in their original papers. The results under “BLOSUM62 average” condition was shown in **Table 3**. “Tanimoto average” condition had a similar trend, the details of the results under “Tanimoto average” condition can be found in supplementary Table S3 and Table S4. **Table 3(a)** and **(b)** told us that these deep learning methods cannot best machine learning methods on AMPCliff prediction tasks, which was consistent to literature⁷.

Table 3 (a)Spearman of deep learning and machine learning methods under BLOSUM62 average condition on the test set. The best Spearman values were highlighted in **bold**. The last row showed the best Spearman value for each dilution difference strategy.

BLOSUM62 average	Models	2-fold	3-fold	4-fold	5-fold
Deep learning	AMPSpace ¹⁸	0.405613	0.324516	0.345483	0.332168
	CellFree-cnn ¹⁹	0.368759	0.335063	0.320002	0.304561
	CellFree-rnn ¹⁹	0.337597	0.313291	0.285324	0.280016
	peptimizer ³⁹	0.292802	0.235428	0.249914	0.290861
Machine Learning	RF	0.492446	0.423145	0.441208	0.447867
	GB	0.454497	0.423129	0.424383	0.391368

	XGBoost	0.497805	0.43087	0.418753	0.431563
	GP	0.401721	0.387301	0.369817	0.364153
	max	0.497805	0.43087	0.441208	0.447867

(b) Recall of deep learning and machine learning methods under BLOSUM62 average condition on the test set. The best recall values were highlighted in bold. The last row showed the best Spearman value for each dilution difference strategy.

BLOSUM62 average	Models	2-fold	3-fold	4-fold	5-fold
Deep learning	AMPSpace ¹⁸	8	8	7	14
	CellFree-cnn ¹⁹	10	12	11	12
	CellFree-rnn ¹⁹	8	12	14	12
	peptimizer ³⁹	6	8	7	12
Machine Learning	RF	6	11	12	19
	GB	10	10	10	18
	XGBoost	9	10	9	18
	GP	9	13	14	16
	max	10	13	14	19

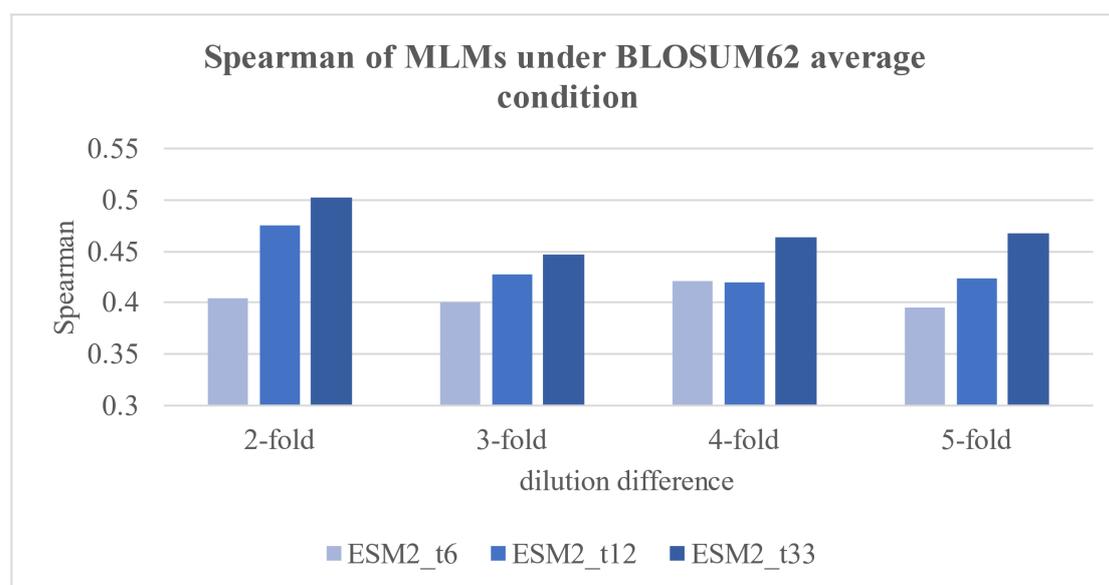
We also tested the recent performances of LMs on AMPCliffs prediction task. For MLMs, we selected original English BERT, ESM2 with 6, 12, and 33 layers (Since ProtBERT performed worse than ESM2, we eliminated it). For GLMs, we selected

GPT2, Progen2 small version, base version and large version. All of them were finetuned by the same hyperparameters, which can be found in Supplementary Table S5.

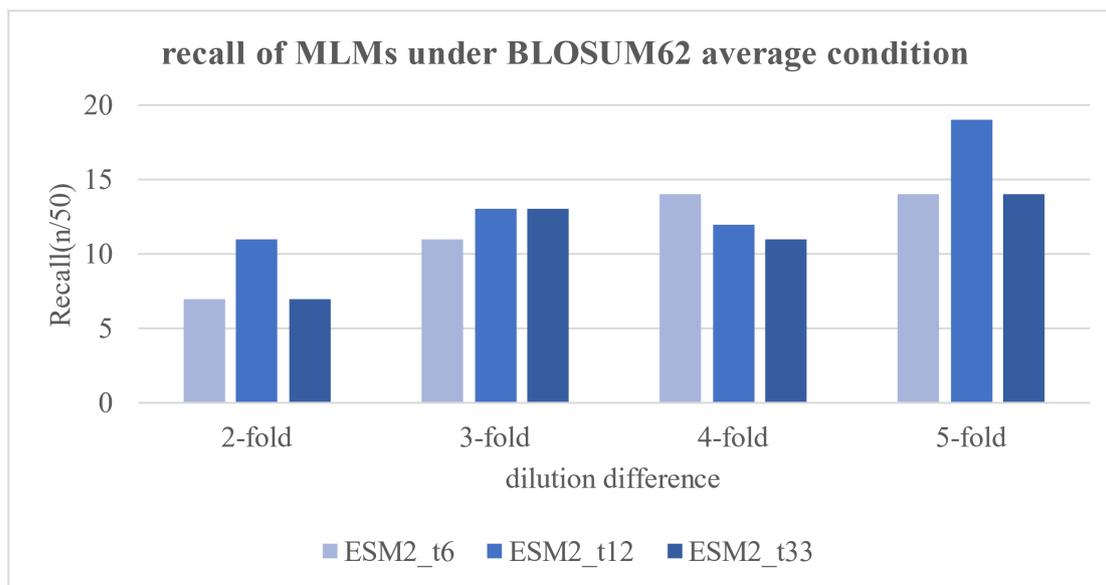
Again, here we took the results of AMPCliff under “BLOSUM62 average” condition prediction as an example, the results of “Tanimoto average” can be found in Supplementary Table S7 and Table S8.

Performances of MLMs

Generally, it had an obvious Spearman improvement as the model size of ESM2 increased under the “BLOSUM62 average” condition, seen the results in **Fig. 13(a)**. This trend existed almost in all kinds of dilution differences. Whereas it disappeared in the recall matrix (see **Fig. 13(b)**), which suggested that the top 50 samples in the test dataset didn’t perform better as the model size increased.



(a)



(b)

Fig. 15 The performances of ESM2 on different dilution difference strategies as the model size increased. The x-axis denoted the dilution difference. (a) The y-axis denoted the value of Spearman index, (b) the y-axis denoted the value of recall.

In addition, we compared the MLMs with the most effective machine learning methods mentioned in the former section, the results were shown in **Table 4**, and found out that the largest model ESM2 with 33 layers beat the machine learning methods and reached the best performances on Spearman among all dilution differences (see **Table 4(a)**). At the meantime, in terms of the recall, although MLMs especially ESM2 didn't follow the same trend as it was in Spearman, they got the better recall compared with the machine learning methods in general (see **Table 4(b)**). Another noteworthy result was the performance of BERT pretrained by English text. Although its Spearman got the worst among all of the MLMs, its recall had competitive performances with ESM2.

Table 4 (a)Spearman of MLMs and machine learning methods under BLOSUM62 average condition on the test set. The best Spearman values were highlighted in **bold**, and the second-best values were highlighted in underscore. The last row showed the best Spearman value for each dilution difference strategy.

	Models	2-fold	3-fold	4-fold	5-fold
Mask Language Model	BERT-base	0.412027	0.362385	0.357889	0.354876
	ESM2_t6	0.403675	0.400395	0.421556	0.394871
	ESM2_t12	0.47584	0.427978	0.419698	0.42403
	ESM2_t33	0.502899	0.446841	0.463059	0.466916
Machine Learning	RF	0.492446	0.423145	<u>0.441208</u>	<u>0.447867</u>
	GB	0.454497	0.423129	0.424383	0.391368
	XGBoost	<u>0.497805</u>	<u>0.43087</u>	0.418753	0.431563
	GP	0.401721	0.387301	0.369817	0.364153
	max	0.502899	0.446841	0.463059	0.466916

(b) Recall of MLMs and machine learning methods under BLOSUM62 average condition on the test set. The best recall values were highlighted in **bold**, and the second-best values were highlighted in underscore. The last row showed the best Spearman value for each dilution difference strategy.

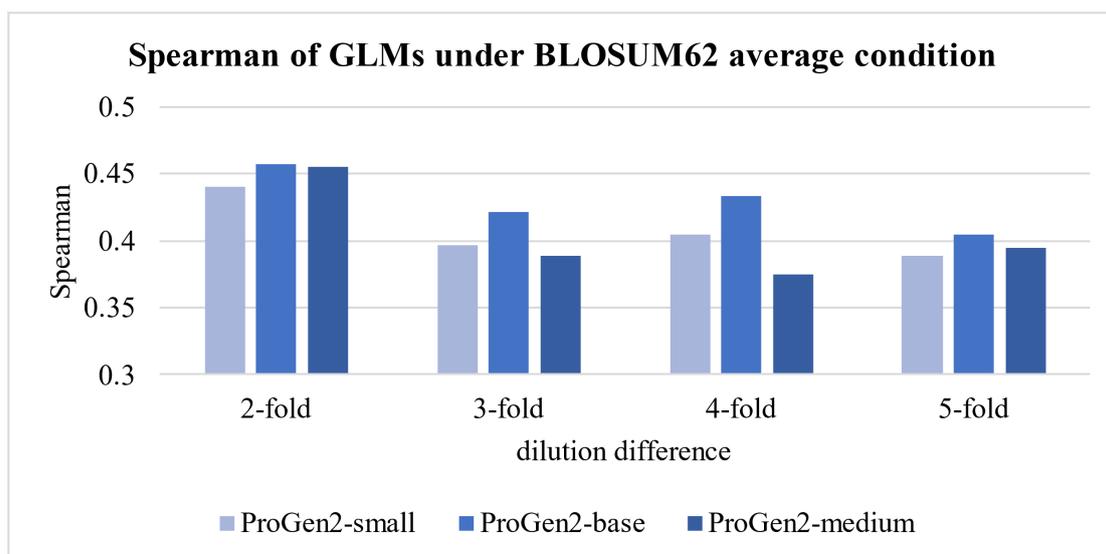
	Models	2-fold	3-fold	4-fold	5-fold
Mask Language Model	BERT-base	<u>10</u>	14	<u>12</u>	13
	ESM2_t6	7	11	14	14
	ESM2_t12	11	<u>13</u>	<u>12</u>	19
	ESM2_t33	7	<u>13</u>	11	14

Machine Learning	RF	6	11	<u>12</u>	19
	GB	<u>10</u>	10	10	<u>18</u>
	XGBoost	9	10	9	<u>18</u>
	GP	9	<u>13</u>	14	16
	max	11	14	14	19

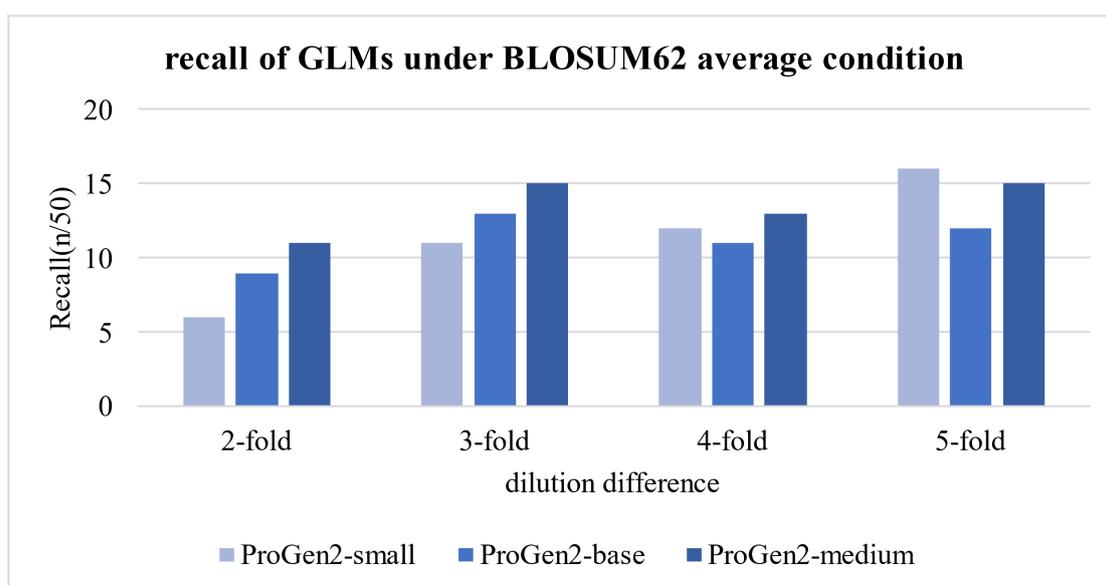
Performances of GLMs

Interestingly, GLMs had a counterpart trend against MLMs, that is, the general trend of recall improvement as the model size of ProGen2 increased (see **Fig. 14(a)**) disappeared in the Spearman matrix (see **Fig. 14(b)**). This suggested that GLMs were good at predicting the top 50 samples in the test dataset other than the general performance. However, when we compared GLMs with the machine learning methods, it turned out GLMs were poorer than machine learning methods in Spearman matrix (see **Table 5(a)**), and only the recall of the medium size of ProGen2 beat machine learning methods on the 2-fold and 3-fold dilution difference (see **Table 5(b)**), and the second-best recall on the 2-fold and 3-fold dilution difference were still the machine learning methods GB and GP respectively. Besides, the best recall of MLMs beat the one of GLMs, which can be found in in **Table 4(b)** and **Table 5(b)**.

If we compared the performance of finetuned BERT and GPT2, which were both pretrained in English text,



(a)



(b)

Fig. 16 The performances of ProGen2 on different dilution difference strategies as the model size increased. The x-axis denoted the dilution difference. (a) The y-axis denoted the value of Spearman index, (b) the y-axis denoted the value of recall.

Table 5 (a)Spearman of GLMs and machine learning methods under BLOSUM62 average condition on the test set. The best Spearman values were highlighted in bold,

and the second-best values were highlighted in underscore. The last row showed the best Spearman value for each dilution difference strategy.

	Models	2-fold	3-fold	4-fold	5-fold
Generative Language Model	GPT2-base	0.402676	0.326114	0.313672	0.317871
	ProGen2-small	0.440555	0.396926	0.405041	0.388975
	ProGen2-base	0.457303	0.421568	<u>0.433841</u>	0.405001
	ProGen2-medium	0.455452	0.388613	0.374782	0.394716
Machine Learning	RF	<u>0.492446</u>	<u>0.423145</u>	0.441208	0.447867
	GB	0.454497	0.423129	0.424383	0.391368
	XGBoost	0.497805	0.43087	0.418753	<u>0.431563</u>
	GP	0.401721	0.387301	0.369817	0.364153
	max	0.497805	0.43087	0.441208	0.447867

(b) Recall of GLMs and machine learning methods under BLOSUM62 average condition on the test set. The best recall values were highlighted in **bold**, and the second-best values were highlighted in underscore. The last row showed the best Spearman value for each dilution difference strategy.

	Models	2-fold	3-fold	4-fold	5-fold
	GPT2-base	4	9	11	12

Generative Language Model	ProGen2-small	6	11	12	16
	ProGen2-base	9	<u>13</u>	11	12
	ProGen2- medium	11	15	<u>13</u>	15
Machine Learning	RF	6	11	12	19
	GB	<u>10</u>	10	10	<u>18</u>
	XGBoost	9	10	9	<u>18</u>
	GP	9	<u>13</u>	14	16
	max	11	15	14	19

Conclusion

This paper we discovered the activity cliff phenomenon in antimicrobial peptides, called AMPCliff. The AMPCliff was defined followed the common sense of activity cliff in molecular, a pair of structurally similar compounds that are active against the same bacteria but significantly different in MIC values. Then we systematically introduced a strict procedure of defining the AMPCliff by using the public AMP dataset containing MIC value GRAMPA, taking one of the richest bacteria *S.aureus* as an example for illustration, and testing the recent algorithms to compare their performance on AMPCliff prediction task.

Following the key elements underlying molecular property prediction proposed by Jianyuan Deng et al.⁶ We proposed a novel data partition method AC-split and recommended recall as another evaluation matrix, then systematically assessing the performances of 9 machine learning methods and 4 deep learning methods proposed by several peptide design work, 4 masked language models and 4 generative language models as the MIC dilution difference increased from 2 to 5.

In general, all of the models suggested that models can learned some potential common patterns from the low-frequency signals to successfully predict the high-frequency signals. ESM2 were the recent best models on AMPCliff prediction task compared with machine learning and GLMs. Its general performance evaluated by Spearman exceeded as the model size became larger, but its recall on the top 50 samples evaluated by recall didn't follow the same trend.

References

- 1 Wang, L. *et al.* Therapeutic peptides: current applications and future directions. *Signal Transduct Target Ther* 7, 48 (2022). <https://doi.org/10.1038/s41392-022-00904-4>
- 2 Tropsha, A., Isayev, O., Varnek, A., Schneider, G. & Cherkasov, A. Integrating QSAR modelling and deep learning in drug discovery: the emergence of deep QSAR. *Nature Reviews Drug Discovery* (2023). <https://doi.org/10.1038/s41573-023-00832-0>

-
- 3 Sadybekov, A. V. & Katritch, V. Computational approaches streamlining drug discovery. *Nature* **616**, 673-685 (2023). <https://doi.org:10.1038/s41586-023-05905-z>
 - 4 van Tilborg, D., Alenicheva, A. & Grisoni, F. Exposing the Limitations of Molecular Machine Learning with Activity Cliffs. *J Chem Inf Model* **62**, 5938-5951 (2022). <https://doi.org:10.1021/acs.jcim.2c01073>
 - 5 Dablander, M., Hanser, T., Lambiotte, R. & Morris, G. M. Exploring QSAR models for activity-cliff prediction. *J Cheminform* **15**, 47 (2023). <https://doi.org:10.1186/s13321-023-00708-w>
 - 6 Deng, J. *et al.* A systematic study of key elements underlying molecular property prediction. *Nat Commun* **14**, 6395 (2023). <https://doi.org:10.1038/s41467-023-41948-6>
 - 7 Xia, J. *et al.* in *Thirty-seventh Conference on Neural Information Processing Systems*.
 - 8 Guha, R. & Van Drie, J. H. Structure–Activity Landscape Index: Identifying and Quantifying Activity Cliffs. *Journal of Chemical Information and Modeling* **48**, 646-658 (2008). <https://doi.org:10.1021/ci7004093>
 - 9 Hu, X., Hu, Y., Vogt, M., Stumpfe, D. & Bajorath, J. MMP-Cliffs: systematic identification of activity cliffs on the basis of matched molecular pairs. *J Chem Inf Model* **52**, 1138-1145 (2012). <https://doi.org:10.1021/ci3001138>
 - 10 Hu, Y. & Bajorath, J. Exploration of 3D Activity Cliffs on the Basis of Compound Binding Modes and Comparison of 2D and 3D Cliffs. *Journal of Chemical Information and Modeling* **52**, 670-677 (2012). <https://doi.org:10.1021/ci300033e>
 - 11 Hu, Y., Furtmann, N., Gütschow, M. & Bajorath, J. Systematic Identification and Classification of Three-Dimensional Activity Cliffs. *Journal of Chemical Information and Modeling* **52**, 1490-1498 (2012). <https://doi.org:10.1021/ci300158v>
 - 12 Muratov, E. N. *et al.* QSAR without borders. *Chem Soc Rev* **49**, 3525-3564 (2020). <https://doi.org:10.1039/d0cs00098a>
 - 13 Seebeck, B., Wagener, M. & Rarey, M. From Activity Cliffs to Target-Specific Scoring Models and Pharmacophore Hypotheses. *ChemMedChem* **6**, 1630-1639 (2011). <https://doi.org:https://doi.org/10.1002/cmdc.201100179>

-
- 14 Stumpfe, D., Hu, H. & Bajorath, J. Evolving Concept of Activity Cliffs. *ACS Omega* **4**, 14360-14368 (2019). <https://doi.org:10.1021/acsomega.9b02221>
 - 15 Stumpfe, D., Hu, H. & Bajorath, J. Advances in exploring activity cliffs. *Journal of Computer-Aided Molecular Design* **34**, 929-942 (2020). <https://doi.org:10.1007/s10822-020-00315-z>
 - 16 Zhang, Z., Zhao, B., Xie, A., Bian, Y. & Zhou, S. Activity Cliff Prediction: Dataset and Benchmark. *arXiv preprint arXiv:2302.07541* (2023).
 - 17 Sisay, M. T., Peltason, L. & Bajorath, J. Structural Interpretation of Activity Cliffs Revealed by Systematic Analysis of Structure–Activity Relationships in Analog Series. *Journal of Chemical Information and Modeling* **49**, 2179-2189 (2009). <https://doi.org:10.1021/ci900243a>
 - 18 Huang, J. *et al.* Identification of potent antimicrobial peptides via a machine-learning pipeline that mines the entire space of peptide sequences. *Nat Biomed Eng* **7**, 797-810 (2023). <https://doi.org:10.1038/s41551-022-00991-2>
 - 19 Pandi, A. *et al.* Cell-free biosynthesis combined with deep learning accelerates de novo-development of antimicrobial peptides. *Nat Commun* **14**, 7197 (2023). <https://doi.org:10.1038/s41467-023-42434-9>
 - 20 Medina-Franco, J. L. *et al.* Characterization of Activity Landscapes Using 2D and 3D Similarity Methods: Consensus Activity Cliffs. *Journal of Chemical Information and Modeling* **49**, 477-491 (2009). <https://doi.org:10.1021/ci800379q>
 - 21 Peltason, L. & Bajorath, J. SAR Index: Quantifying the Nature of Structure–Activity Relationships. *Journal of Medicinal Chemistry* **50**, 5571-5578 (2007). <https://doi.org:10.1021/jm0705713>
 - 22 Guha, R. Exploring uncharted territories: predicting activity cliffs in structure-activity landscapes. *J Chem Inf Model* **52**, 2181-2191 (2012). <https://doi.org:10.1021/ci300047k>
 - 23 Heikamp, K., Hu, X., Yan, A. & Bajorath, J. Prediction of activity cliffs using support vector machines. *J Chem Inf Model* **52**, 2354-2365 (2012). <https://doi.org:10.1021/ci300306a>
 - 24 Tancik, M. *et al.* Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems* **33**, 7537-7547 (2020).

-
- 25 Jacot, A., Gabriel, F. & Hongler, C. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems* **31** (2018).
 - 26 Rahaman, N. *et al.* in *International Conference on Machine Learning*. 5301-5310 (PMLR).
 - 27 Ronen, B., Jacobs, D., Kasten, Y. & Kritchman, S. The convergence rate of neural networks for learned functions of different frequencies. *Advances in Neural Information Processing Systems* **32** (2019).
 - 28 Witten, J. & Witten, Z. Deep learning regression model for antimicrobial peptide design. *BioRxiv*, 692681 (2019).
 - 29 Pendleton, J. N., Gorman, S. P. & Gilmore, B. F. Clinical relevance of the ESKAPE pathogens. *Expert review of anti-infective therapy* **11**, 297-308 (2013).
 - 30 王猛 *et al.* 以核糖体蛋白 L7/L12 为分子标志物精准检测肺炎链球菌的研究. *中国生物工程杂志* **40**, 34-41 (2020).
 - 31 Singh, V., Shrivastava, S., Kumar Singh, S., Kumar, A. & Saxena, S. StaBle-ABPpred: a stacked ensemble predictor based on biLSTM and attention mechanism for accelerated discovery of antibacterial peptides. *Brief Bioinform* **23** (2022). <https://doi.org:10.1093/bib/bbab439>
 - 32 Dziuba, B. & Dziuba, M. New milk protein-derived peptides with potential antimicrobial activity: An approach based on bioinformatic studies. *International journal of molecular sciences* **15**, 14531-14545 (2014).
 - 33 Dhall, A., Patiyal, S., Sharma, N., Usmani, S. S. & Raghava, G. P. Computer-aided prediction and design of IL-6 inducing peptides: IL-6 plays a crucial role in COVID-19. *Briefings in bioinformatics* **22**, 936-945 (2021).
 - 34 Levitt, M. Conformational preferences of amino acids in globular proteins. *Biochemistry* **17**, 4277-4285 (1978).
 - 35 陈瞿硕. GM-Pep: 从头设计功能性肽序列的高效策略. (2022).
 - 36 Das, P. *et al.* Accelerated antimicrobial discovery via deep generative models and molecular dynamics simulations. *Nat Biomed Eng* **5**, 613-623 (2021). <https://doi.org:10.1038/s41551-021-00689-x>

-
- 37 Garcia-Jacas, C. R., Garcia-Gonzalez, L. A., Martinez-Rios, F., Tapia-Contreras, I. P. & Brizuela, C. A. Handcrafted versus non-handcrafted (self-supervised) features for the classification of antimicrobial peptides: complementary or redundant? *Brief Bioinform* **23** (2022). <https://doi.org/10.1093/bib/bbac428>
- 38 Fernandez-Diaz, R. *et al.* AutoPeptideML: Automated Machine Learning for Building Trustworthy Peptide Bioactivity Predictors. *bioRxiv*, 2023.2011.2013.566825 (2023).
- 39 Schissel, C. K. *et al.* Deep learning to design nuclear-targeting abiotic miniproteins. *Nat Chem* **13**, 992-1000 (2021). <https://doi.org/10.1038/s41557-021-00766-3>
- 40 Elnaggar, A. *et al.* Prottrans: Toward understanding the language of life through self-supervised learning. *IEEE transactions on pattern analysis and machine intelligence* **44**, 7112-7127 (2021).
- 41 Ferruz, N., Schmidt, S. & Hocker, B. ProtGPT2 is a deep unsupervised language model for protein design. *Nat Commun* **13**, 4348 (2022). <https://doi.org/10.1038/s41467-022-32007-7>
- 42 Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic acids research* **25**, 3389-3402 (1997).
- 43 Mullan, L. EMBOSS|User docs|Introductory Bioinformatics Course. (2000).
- 44 Teufel, F. *et al.* GraphPart: homology partitioning for biological sequence analysis. *NAR Genom Bioinform* **5**, lqad088 (2023). <https://doi.org/10.1093/nargab/lqad088>
- 45 Baylon, J. L. *et al.* PepSeA: Peptide Sequence Alignment and Visualization Tools to Enable Lead Optimization. *J Chem Inf Model* **62**, 1259-1267 (2022). <https://doi.org/10.1021/acs.jcim.1c01360>
- 46 Almagro Armenteros, J. J. *et al.* SignalP 5.0 improves signal peptide predictions using deep neural networks. *Nature biotechnology* **37**, 420-423 (2019).
- 47 Gíslason, M. H., Nielsen, H., Armenteros, J. J. A. & Johansen, A. R. Prediction of GPI-anchored proteins with pointer neural networks. *Current Research in Biotechnology* **3**, 6-13 (2021).
- 48 Steinegger, M. & Söding, J. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nature biotechnology* **35**, 1026-1028 (2017).

-
- 49 Consortium, U. UniProt: a worldwide hub of protein knowledge. *Nucleic acids research* **47**, D506-D515 (2019).
- 50 Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science* **379**, 1123-1130 (2023).
- 51 Madani, A. *et al.* Large language models generate functional protein sequences across diverse families. *Nature Biotechnology* **41**, 1099-1106 (2023).
- 52 UniProt|Help UniRef, <<https://www.uniprot.org/help/uniref>> (2024).
- 53 Hauser, M., Steinegger, M. & Soding, J. MMseqs software suite for fast and deep clustering and searching of large protein sequence sets. *Bioinformatics* **32**, 1323-1330 (2016). <https://doi.org:10.1093/bioinformatics/btw006>
- 54 Andreatta, M., Alvarez, B. & Nielsen, M. GibbsCluster: unsupervised clustering and alignment of peptide sequences. *Nucleic Acids Res* **45**, W458-W463 (2017). <https://doi.org:10.1093/nar/gkx248>
- 55 Martin Steinegger, M. M., Eli Levy Karin, Lars von den Driesch, Clovis Galiez, Johannes Söding. *MMseqs2 User Guide*, <<https://github.com/soedinglab/mmseqs2/wiki>> (2023).
- 56 Mouton, J. W. *et al.* MIC-based dose adjustment: facts and fables. *J Antimicrob Chemother* **73**, 564-568 (2018). <https://doi.org:10.1093/jac/dkx427>

AMPCliff: quantitative definition and benchmarking of activity cliffs in antimicrobial peptides

Kewei Li¹, Yuqian Wu², Yutong Guo³, Yinheng Li⁴, Yusi Fan¹, Ruochi Zhang^{5,*}, Lan Huang¹, Fengfeng Zhou^{1,6,*}.

1 College of computer Science and Technology, and Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education, Jilin University, Changchun, Jilin, China, 130012.

2 School of Software, Jilin University, Changchun 130012, Jilin, China.

3. School of Life Sciences, Jilin University, Changchun 130012, Jilin, China.

4 Department of Computer Science, Columbia University, 116th and Broadway, New York City, New York, 10027, United States.

5 School of Artificial Intelligence, and Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education, Jilin University, Changchun, Jilin, China, 130012.

6 School of Biology and Engineering, Guizhou Medical University, Guiyang 550025, Guizhou, China.

* Correspondence may be addressed to Fengfeng Zhou: FengfengZhou@gmail.com or ffzhou@jlu.edu.cn. Correspondence may also be addressed to Ruochi Zhang: zrc720@gmail.com.

Emails: Kewei Li, kwbb1997@gmail.com; Yuqian Wu, 924476388@qq.com; Yutong Guo, 3291865187.qq.com; Yinheng Li, yl4039@columbia.edu; Yusi Fan, fan_yusi@163.com; Ruochi Zhang, zrc720@gmail.com; Lan Huang, huanglan@jlu.edu.cn; Fengfeng Zhou, FengfengZhou@gmail.com.

Supplementary Table S1

Hyperparameter settings for deep learning methods. Except for the initial learning rate, all of the model parameter setting followed the original papers

AMPSpace ^[33]	CellFree-cnn ^[30]	CellFree-rnn ^[30]	peptimizer ^[35]
Initial learning_rate=2e-3	Initial learning_rate=2e-3	Initial learning_rate=2e-3	Initial learning_rate=2e-3
embedding_dim=50,	kernel_size=5	hidden_dim=500	kernel_size=2,

hidden_dim=128, num_layer=2, bidirectional=False, dropout=0.7	conv_layers=2 channel=64 activate_func=relu linear_dim1=100 linear_dim2=100 linear_dim3=20	linear_dim=100	conv_layers=3 channels=256, dropout=0.1, linear_dim=256 radius=3 n_bits=2-048
--	---	----------------	--

Supplementary Table S2

The details of the results of each machine learning methods

Recall tanimoto average

models	diff2	diff3	diff4	diff5
SVM	4	7	6	7
LR	8	9	11	10
L1	9	10	10	10
L2	6	8	9	9
ElasticNet	7	6	7	9
RF	8	12	14	17
GB	11	9	10	13
XGBoost	13	10	15	17
GP	8	13	13	15
max	13	13	15	17

Spearman tanimoto average

	diff2	diff3	diff4	diff5
SVM	0.311585	0.271239	0.246765	0.231313
LR	0.283188	0.273977	0.309939	0.297302
L1	0.386778	0.38181	0.392317	0.390961
L2	0.312565	0.305443	0.334129	0.320098
ElasticNet	0.410903	0.384667	0.398571	0.394745
RF	0.492287	0.435887	0.451537	0.418049
GB	0.490807	0.452679	0.449903	0.424813
XGBoost	0.498402	0.443067	0.441797	0.413957
GP	0.394693	0.383093	0.352759	0.320407
max	0.498402	0.452679	0.451537	0.424813

Recall blosum62 average

	diff2	diff3	diff4	diff5
SVM	3	6	5	9
LR	9	9	10	12
L1	9	8	11	10

L2	8	8	8	15
ElasticNet	7	8	8	10
RF	6	11	12	19
GB	10	10	10	18
XGBoost	9	10	9	18
GP	9	13	14	16
max	10	13	14	19

Spearman blosum62 average

	diff2	diff3	diff4	diff5
SVM	0.315043	0.295966	0.287912	0.299459
LR	0.318579	0.295257	0.308216	0.288063
L1	0.376329	0.360152	0.372086	0.376945
L2	0.338543	0.321677	0.327059	0.308657
ElasticNet	0.394322	0.360766	0.372294	0.377859
RF	0.492446	0.423145	0.441208	0.447867
GB	0.454497	0.423129	0.424383	0.391368
XGBoost	0.497805	0.43087	0.418753	0.431563
GP	0.401721	0.387301	0.369817	0.364153
max	0.497805	0.43087	0.441208	0.447867

Supplementary Table S3

(a)Spearman of deep learning and machine learning methods under blosum62 average condition on the test set. The best Spearman values were highlighted in **bold**. The last row showed the best Spearman value for each dilution difference strategy.

blosum62 average	Models	diff2	diff3	diff4	diff5
Deep learning	AMPSSpace	0.405613	0.324516	0.345483	0.332168
	CellFree-cnn	0.368759	0.335063	0.320002	0.304561
	CellFree-rnn	0.337597	0.313291	0.285324	0.280016
	peptimizer	0.292802	0.235428	0.249914	0.290861
Machine learning	RF	0.492446	0.423145	0.441208	0.447867
	GB	0.454497	0.423129	0.424383	0.391368
	XGBoost	0.497805	0.43087	0.418753	0.431563
	GP	0.401721	0.387301	0.369817	0.364153
	max	0.497805	0.43087	0.441208	0.447867

(b)Recall of deep learning and machine learning methods under blosum62 average condition on the test set. The best recall values were highlighted in **bold**. The last row showed the best Spearman value for each dilution difference strategy.

blosum62 average	Models	diff2	diff3	diff4	diff5
------------------	--------	-------	-------	-------	-------

Deep learning	AMPSpace	8	8	7	14
	CellFree-cnn	10	12	11	12
	CellFree-rnn	8	12	14	12
	peptimizer	6	8	7	12
Machine learning	RF	6	11	12	19
	GB	10	10	10	18
	XGBoost	9	10	9	18
	GP	9	13	14	16
	max	10	13	14	19

Supplementary Table S4

(a) Spearman of deep learning and machine learning methods under tanimoto average condition on the test set. The best Spearman values were highlighted in **bold**. The last row showed the best Spearman value for each dilution difference strategy.

Tanimoto average	models	diff2	diff3	diff4	diff5
Deep learning	AMPSpace	0.513214	0.39439	0.376107	0.330334
	CellFree-cnn	0.501194	0.404835	0.357729	0.318542
	CellFree-rnn	0.464495	0.347788	0.349593	0.328807
	peptimizer	0.416705	0.262073	0.232673	0.284782
Machine learning	RF	0.492287	0.435887	0.451537	0.418049
	GB	0.490807	0.452679	0.449903	0.424813
	XGBoost	0.498402	0.443067	0.441797	0.413957
	GP	0.394693	0.383093	0.352759	0.320407
	max	0.513214	0.404835	0.376107	0.330334

(b) Recall of deep learning and machine learning methods under tanimoto average condition on the test set. The best recall values were highlighted in **bold**. The last row showed the best Spearman value for each dilution difference strategy.

Tanimoto average	models	diff2	diff3	diff4	diff5
Deep learning	AMPSpace	14	11	14	13
	CellFree-cnn	11	8	9	12
	CellFree-rnn	13	13	10	10
	peptimizer	9	9	10	9
Machine learning	RF	8	12	14	17
	GB	11	9	10	13
	XGBoost	13	10	15	17
	GP	8	13	13	15
	max	14	13	14	13

Supplementary Table S5

The hyperparameters of LMs

Parameter	setting
random_seed	0
num_epoch	50
batch_size	4
num_workers	1
learning_rate	1e-5
adam_epsilon	1e-8
weight_decay	0.01
loss	mse
lr_scheduler	cosine(T_0=10, T_mult=2, eta_min=1e-6), batch

Supplementary Table S7

(a) Spearman of MLMs and machine learning methods under Tanimoto average condition on the test set. The best Spearman values were highlighted in **bold**, and the second-best values were highlighted in underscore. The last row showed the best Spearman value for each dilution difference strategy.

	Models	2-fold	3-fold	4-fold	5-fold
Mask Language Model	BERT-base	0.577995	0.486886	0.458012	0.418275
	ESM2_t6	0.597729	0.517428	<u>0.476475</u>	0.434274
	ESM2_t12	<u>0.649133</u>	<u>0.562948</u>	0.52798	0.487287
	ESM2_t33	0.676772	0.581751	0.458346	<u>0.481701</u>
Machine Learning	RF	0.492287	0.435887	0.451537	0.418049
	GB	0.490807	0.452679	0.449903	0.424813
	XGBoost	0.498402	0.443067	0.441797	0.413957
	GP	0.394693	0.383093	0.352759	0.320407
	max	0.676772	0.581751	0.52798	0.487287

(b) Recall of MLMs and machine learning methods under Tanimoto average condition on the test set. The best recall values were highlighted in **bold**, and the second-best values were highlighted in underscore. The last row showed the best Spearman value for each dilution difference strategy.

	Models	2-fold	3-fold	4-fold	5-fold
Mask Language Model	BERT-base	14	16	17	15
	ESM2_t6	<u>18</u>	<u>17</u>	17	16
	ESM2_t12	20	18	17	20
	ESM2_t33	<u>18</u>	16	<u>16</u>	16
Machine Learning	RF	8	12	14	<u>17</u>
	GB	11	9	10	13

	XGBoost	13	10	15	<u>17</u>
	GP	8	13	13	15
	max	20	18	17	20

Supplementary Table S8

(a) Spearman of GLMs and machine learning methods under Tanimoto average condition on the test set. The best Spearman values were highlighted in **bold**, and the second-best values were highlighted in underscore. The last row showed the best Spearman value for each dilution difference strategy.

	Models	2-fold	3-fold	4-fold	5-fold
Generative Language Model	GPT2-base	0.509833	0.392148	0.361492	0.366459
	ProGen2-small	0.630527	<u>0.521846</u>	0.481506	0.415372
	ProGen2-base	<u>0.633118</u>	0.510379	0.41904	0.463478
	ProGen2-medium	0.637873	0.532595	0.379728	0.379319
Machine Learning	RF	0.492287	0.435887	<u>0.451537</u>	0.418049
	GB	0.490807	0.452679	0.449903	<u>0.424813</u>
	XGBoost	0.498402	0.443067	0.441797	0.413957
	GP	0.394693	0.383093	0.352759	0.320407
	max	0.637873	0.532595	0.481506	0.463478

(b) Recall of GLMs and machine learning methods under Tanimoto average condition on the test set. The best recall values were highlighted in **bold**, and the second-best values were highlighted in underscore. The last row showed the best Spearman value for each dilution difference strategy.

	Models	2-fold	3-fold	4-fold	5-fold
Generative Language Model	GPT2-base	10	14	13	13
	ProGen2-small	<u>21</u>	16	17	11
	ProGen2-base	19	16	<u>15</u>	14
	ProGen2-medium	24	21	13	9
Machine Learning	RF	8	12	14	17
	GB	11	9	10	13
	XGBoost	13	10	<u>15</u>	17
	GP	8	13	13	15
	max	24	21	17	17