

# Exploration of DPP-IV Inhibitory Peptide Design Rules Assisted by the Deep Learning Pipeline That Identifies the Restriction Enzyme Cutting Site

Changge Guan,<sup>#,\*</sup> Jiawei Luo,<sup>#</sup> Shucheng Li,<sup>#</sup> Zheng Lin Tan,<sup>#</sup> Yi Wang, Haihong Chen, Naoyuki Yamamoto, Chong Zhang, Yuan Lu, Junjie Chen, and Xin-Hui Xing<sup>\*</sup>



Cite This: *ACS Omega* 2023, 8, 39662–39672



Read Online

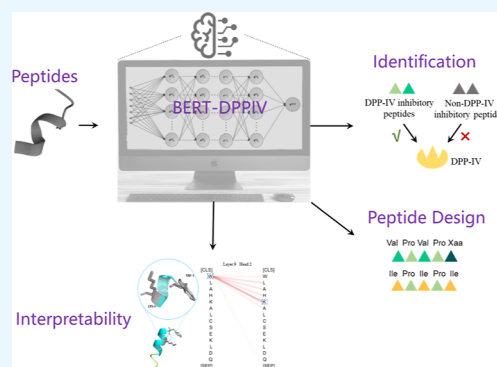
ACCESS |

Metrics & More

Article Recommendations

Supporting Information

**ABSTRACT:** The mining of antidiabetic dipeptidyl peptidase IV (DPP-IV) inhibitory peptides (DPP-IV-IPs) is currently a costly and laborious process. Due to the absence of rational peptide design rules, it relies on cumbersome screening of unknown enzyme hydrolysates. Here, we present an enhanced deep learning model called bidirectional encoder representation (BERT)–DPPIV, specifically designed to classify DPP-IV-IPs and explore their design rules to discover potent candidates. The end-to-end model utilizes a fine-tuned BERT architecture to extract structural/functional information from input peptides and accurately identify DPP-IV-IPs from input peptides. Experimental results in the benchmark data set showed BERT–DPPIV yielded state-of-the-art accuracy and MCC of 0.894 and 0.790, surpassing the 0.797 and 0.594 obtained by the sequence-feature model. Furthermore, we leveraged the attention mechanism to uncover that our model could recognize the restriction enzyme cutting site and specific residues that contribute to the inhibition of DPP-IV. Moreover, guided by BERT–DPPIV, proposed design rules for DPP-IV inhibitory tripeptides and pentapeptides were validated, and they can be used to screen potent DPP-IV-IPs.



## INTRODUCTION

Due to the side effects and the requirement of injection of commercial dipeptidyl peptidase IV (DPP-IV) inhibitors, which are used to treat approximately 537 million patients with type II diabetes, it is crucial to develop new DPP-IV inhibitory drugs and functional foods.<sup>1–3</sup> Peptides with DPP-IV inhibitory activities are a promising class of oral hypoglycemics without adverse effects that are derived from products of enzymatically hydrolyzed edible animal, plant, and macroalgal proteins.<sup>1</sup>

Current approaches for the discovery of DPP-IV inhibitory peptides (DPP-IV-IPs) are known to be labor- and cost-intensive. For instance, the conventional enzymatic hydrolysate screening method requires extensive separation and purification, mass spectrometric identification, and revalidation of peptide synthesis.<sup>2,3</sup> In addition, the lack of design rules and effective characterization methods makes it challenging to design DPP-IV-IPs and synthesize peptide libraries for DPP-IV-IPs screening.<sup>4</sup> These problems have limited the discovery of efficient DPP-IV-IPs. Therefore, it is an important task to develop a high-throughput method capable of rapidly identifying effective DPP-IV-IPs that can be used in functional foods or medicine innovation.<sup>1,2,5</sup>

Recently, powered by artificial intelligence, large language models based on natural language processing (NLP) have been successfully used to solve biological problems. Therefore,

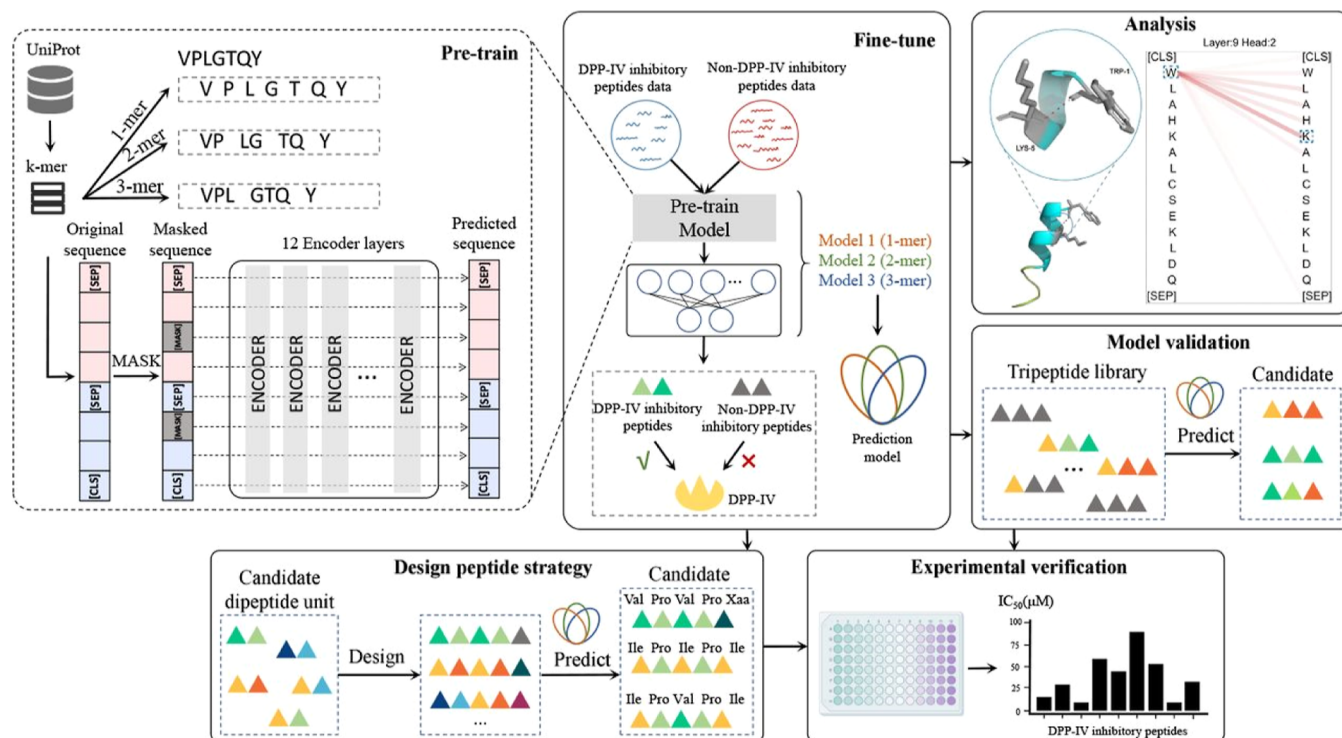
computation tools complement experimental studies in the DPP-IV-IPs discovery process. Over the past decade, several computational approaches have been used for the discovery of DPP-IV-IPs, including the quantitative structure–activity relationship method<sup>6–9</sup> and machine learning (ML).<sup>7,10–13</sup> The ML methods mainly include support vector machines and random forest.<sup>14–16</sup> However, such methods reported to date have some major limitations, leading to poor performance: (1) the performance of these models is heavily dependent on the quality of the features extracted by feature engineering, (2) they are poorly able to represent amino acid sequences, and (3) they fail to capture information hidden in the amino acid sequence itself, which renders these models inefficient. Furthermore, classifiers for DPP-IV-IPs developed to date have not been verified experimentally due to the low accuracy of the model, which will result in the high cost of experimental verification.

**Received:** July 30, 2023

**Accepted:** September 27, 2023

**Published:** October 13, 2023





**Figure 1.** Schematic representation of the study workflow. We first collected sequences to pretrain a BERT model (upper left). We then used the DPP-IV-IP data set for fine-tuning the pretrained BERT model to obtain three kinds of PLMs. The three models were combined to construct the BERT–DPPIV (upper middle). To determine what the model had learned, we analyzed and visualized the model’s attention (upper right, analysis module). We then used BERT–DPPIV to screen tripeptides that may inhibit DPP-IV to illustrate the performance of our model inhibition (upper right, model validation module). After the model screening, the functional peptides screened out by the model were verified by biological experiments (lower right). Having demonstrated the practical applicability of our model, we proposed a design strategy for DPP-IV-IPs based on dipeptide repeat units and demonstrated the feasibility of this strategy using models and experiments (lower left).

Deep neural networks with advanced architecture can overcome these limitations and have the capability of automating feature learning to extract discriminating feature representations with minimal human effort.<sup>17</sup> In particular, deep learning (DL) methods, e.g., long–short-term memory, which is based on NLP and regards amino acid sequence as natural language, are a promising method for identifying functional peptides for extraction of discriminative feature representation, such as antimicrobial peptides.<sup>18–21</sup> However, DL-based models have been criticized for their interpretability and unexplainable characteristics, which are referred to as “black box”. Moreover, the limited volume of data is a great challenge to training DL models from scratch.

In this article, we present the first novel peptide language model (PLM)-based DL model, named BERT–DPPIV, to offer a promising solution to address the challenges of DPP-IV-IPs screening and design. BERT–DPPIV combines a pretraining process based on unlabeled protein data and a NLP BERT model with an attention mechanism to identify DPP-IV-IPs. By learning the discrimination task, BERT–DPPIV was trained to automatically learn the characteristics and features contained in the original peptide sequences and to distinguish and represent amino acid sequences in high-dimensional spaces, achieved state-of-the-art.

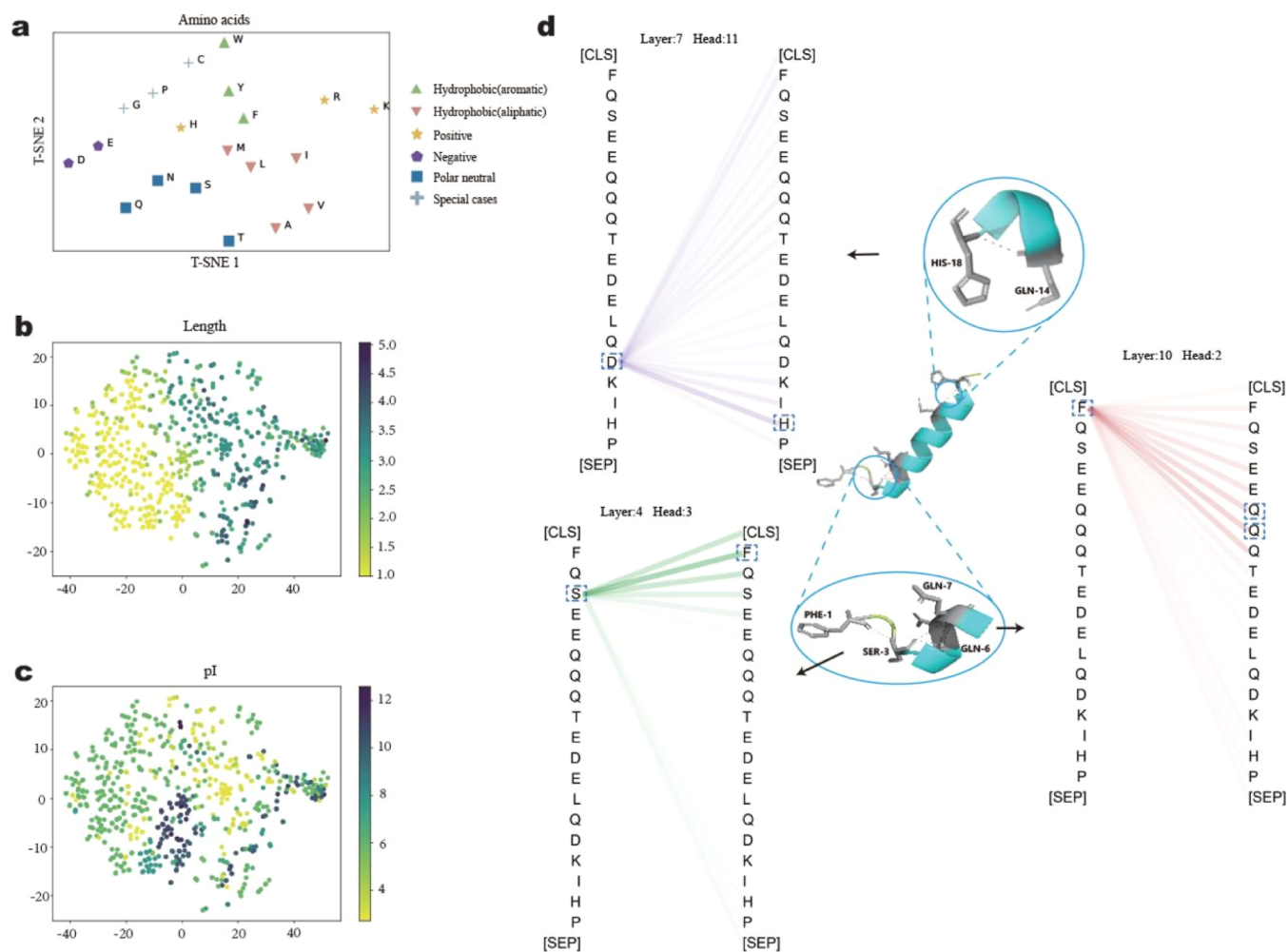
Benefiting from the attention mechanism, BERT–DPPIV has good interpretability, which is more advantageous than existing ML algorithms for DPP-IV-IPs identification. According to the attention analyses, we have found that our model can learn the tertiary structure properties of peptides except for

physicochemical properties. Particularly, we have found that BERT–DPPIV can identify cleavage sites of DPP-IV enzymes on the substrate polypeptides, which was first reported, and suggested that the NLP DL model may be used to analyze enzyme cleavage sites.

The accuracy of BERT–DPPIV was validated by a wet laboratory experiment based on measurement of the half maximal inhibitory concentration ( $IC_{50}$ ) of synthetic peptides predicted to be DPP-IV-IPs by BERT–DPPIV, and the results suggest that the prediction accuracy is in alignment with experimental data, thus extending the computational model to practical applications. By combining both the PLM model and biological assays, we have proposed a novel DPP-IV inhibitory pentapeptide design strategy based on the dipeptide repeat unit X-proline to provide a new idea for the design of DPP-IV-IPs. This design strategy has addressed the bottleneck in DPP-IV-IP discovery and demonstrated that the *in silico* evaluation method can aid peptide drug development.

## RESULTS

**Overview of BERT–DPPIV.** BERT–DPPIV is a PLM-based DL framework designed to mine novel DPP-IV-IPs based on their amino acid sequence (Figure 1). To realize this, the framework consists of four main components: pretraining, fine-tuning, analysis, and model validation and application. During the pretraining procedures, a total of 556,603 protein sequences from databases including UniProt, SWISS-PROT, TrEMBL, and PIR–PSD were used to pretrain 12 layers of BERT-based language models. We employed three different



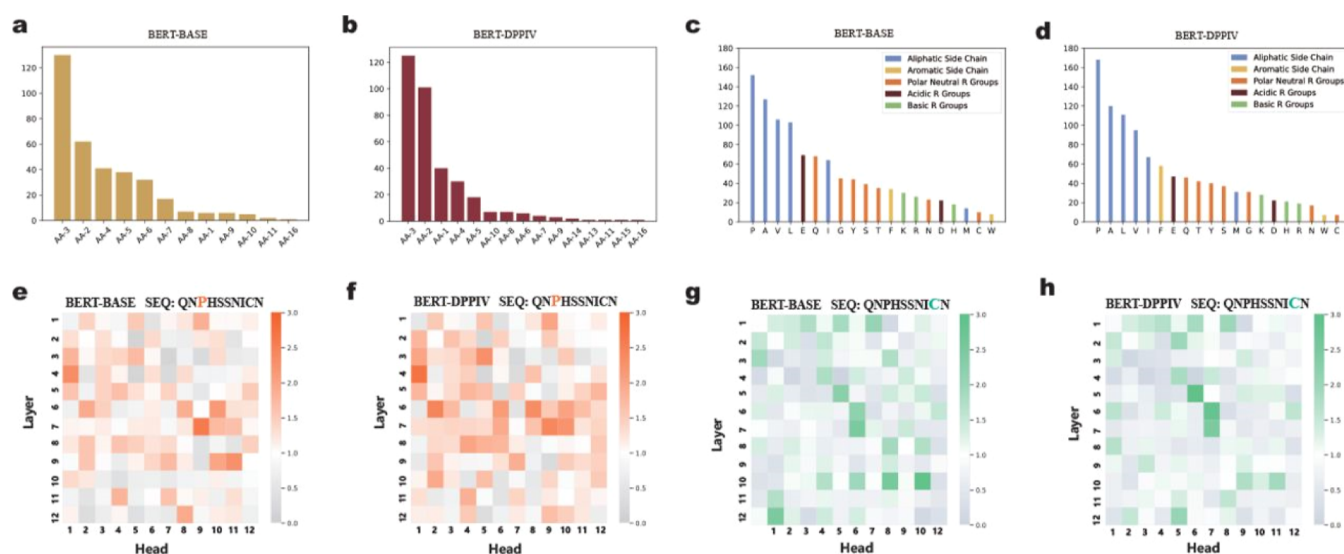
**Figure 2.** Sequence representation visualization and attention visualization of structural information learned by the model. A vector representation of the peptide sequence and of individual amino acids was extracted from the model and analyzed by *t*-SNE. The amino acid representation learned by the model (a) and the physicochemical properties contained in the model-learned representations of peptide sequences (b,c) are shown. Here, we show how the inner workings of the model's attention heads can be used to analyze a single peptide in more detail. Each attention head performs internally an all-against-all comparison to compute weighted sums for each token relative to all other tokens in the sequence. High scores indicate that the model learned to put more weight on certain residue pairs (upper left, lower left, and lower right). The structure of the peptide shows that Ser-3 can interact with Phe-1, Gln-6, and Gln-7 and that Gln-14 can interact with His-18. When visualizing the attention weights for these sites, we observed that some of the attention heads could capture this interaction information.

word segmentation approaches, i.e.,  $k$ -mer = 1,  $k$ -mer = 2, and  $k$ -mer = 3, to generate three pretrained BERT models. Two major tasks, i.e., the masked language model (MLM) and next sentence prediction (NSP), are involved in the pretraining procedures for capturing word-level and sentence-level representations and learning the common features of protein sequences. The fine-tuning process aimed to construct the PLM model specifically for the task of identifying DPP-IV-IPs.

The framework's analysis component involved evaluating the performance and capabilities of BERT-DPP-IV. Various metrics were conducted to assess the model's ability to predict and identify DPP-IV-IPs accurately. Three models with different  $k$ -mers in the test set showed the highest accuracy (Acc) and Matthews correlation coefficient (MCC) of 0.891 and 0.784 were achieved for  $k$ -mer = 1, followed by 0.887 and 0.775 for  $k$ -mer = 2, and 0.842 and 0.684 for  $k$ -mer = 3, and BERT-DPP-IV combining three models achieved state-of-the-art, which resulted in an improved Acc and MCC of 0.894 and 0.790 (Supporting Information Table S1).

Our models ( $k$ -mer = 1,  $k$ -mer = 2, and  $k$ -mer = 3) outperform the first sequence-based iDPP-IV-SCM model, as evidenced by significant improvements in Acc, MCC, and AUC, respectively (Supporting Information Table S1). These results highlight the ability of our models to extract more valuable information from peptide sequences. To further compare with sequence-based physicochemical properties, the SVM model and the sequence- and structure-based StackDPP-IV model. Our model ( $k$ -mer = 1) achieves the same accuracy and MCC as the best model, StackDPP-IV, and a comparable AUC of 0.957 less than 0.961 of the StackDPP-IV. BERT-DPP-IV outperforms it by improvements of 0.003 and 0.006 in Acc and MCC, respectively, and achieves a comparable AUC of 0.960. These results suggest that our models are able to extract physicochemical properties or structure information, and BERT-DPP-IV can make good use of information from different models. Moreover, visualization of the model demonstrated these findings further.

Finally, the model was validated and applied to real-world scenarios, providing a tool for mining and discovering novel



**Figure 3.** Visualization of the position and importance of the amino acid species. To study the importance of position and amino acid species, attention to positions and amino acids was statistically analyzed. BERT-BASE and BERT-DPPIV represent the models before and after fine-tuning, respectively. The statistical importance of the position of the peptide sequence according to the attention of the model is shown (a,b). The statistical importance of amino acid species according to the attention of the model is shown (c,d). The attention of the model to proline and cysteine in the peptide sequence was visualized (e–h).

DPP-IV-IPs based on their amino acid sequences. Overall, we demonstrate the capability of BERT-DPPIV to predict important parameters for DPP-IV-IPs to guide DPP-IV-IPs screening, i.e., (1) inhibitory activity of peptides, (2) physicochemical properties of peptide sequences, (3) structural properties from peptide sequences, and (4) capture of cleavage site information. Furthermore, we have also shown the application of BERT-DPPIV in guiding the design of DPP-IV-IPs in addition to peptide screening. Pentapeptides were designed, and their DPP-IV inhibitory activities were predicted and verified.

**Sequence Representation and Its Properties Learned by BERT-DPPIV.** To investigate the information captured by our models, we selected  $k$ -mer = 1, which exhibits the highest Acc among our three  $k$ -mers models, for further study. First, we analyzed the representation of the peptide sequences learned by the pretrained BERT, as prior studies have shown that BERT models can effectively learn the representation of the sequence.<sup>18–21</sup> The vector representations of the peptide sequences during the training process of the model were extracted and projected to 2D by  $t$ -distributed stochastic neighbor embedding ( $t$ -SNE) for visualization (Supporting Information Figure S1), which showed that the model was able to identify DPP-IV-IPs by learning the vector representation of peptides.

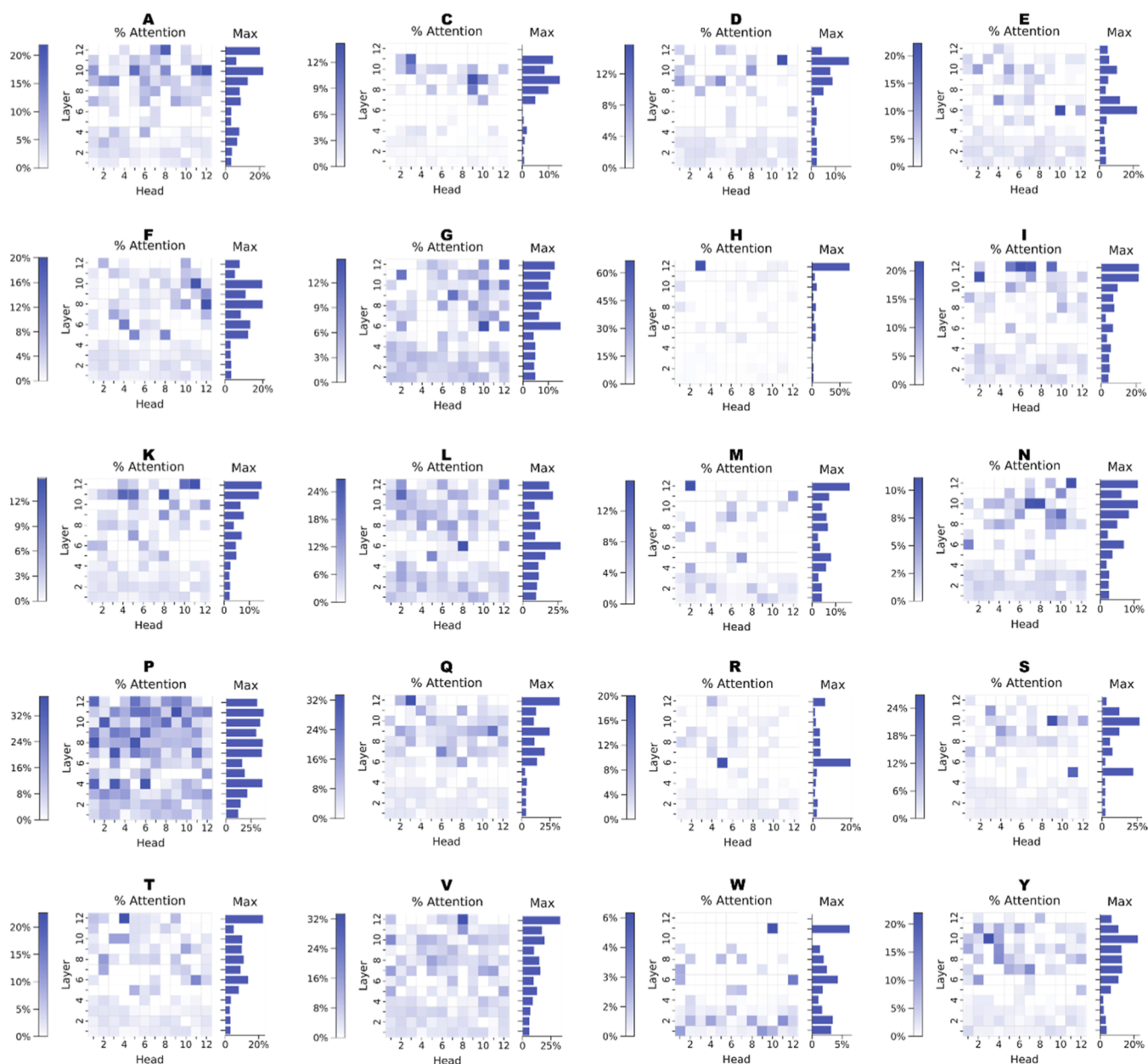
In addition, we have illustrated the ability of our model to learn the physicochemical properties of peptides by visualizing the distribution of amino acid and peptide representations. The vector representation learned by the model was well able to distinguish amino acids with different properties (Figure 2a). We then used the modLAMP package<sup>22</sup> to extract 10 kinds of physicochemical property information from peptide sequences. Based on the peptide sequence physicochemical properties, the vector representation of it was then analyzed. The vector representations learned by the model were able to distinguish among sequences with different physical and chemical properties (Figure 2b,c, Supporting Information

Figure S2). To this point, we have shown that the PLM can learn the physicochemical properties of the peptide sequences.

**Capturing the Structural Information and Attention Patterns Generated by the Models.** The attention mechanism provides a way to reveal the “black box” of DL models and features learned by it.<sup>23,24</sup> In this section, we investigated the capability of our model to learn structural information about peptide sequences with the help of an attention mechanism, which is a technology that mimics cognitive attention. First, the structure of the peptides with a length of  $\geq 5$  amino acids from positive samples was predicted by APPTTEST software,<sup>25</sup> and 12 kinds of peptides with  $\alpha$ -helix secondary structure were found (Supporting Information Figure S3). Three peptides were randomly selected from among these peptides for further structural confirmation using AlphaFold2.<sup>26</sup> The presence of  $\alpha$ -helical structures was confirmed.

Then, the attention of the model for each head at each layer of the model was visualized (12 layers, 12 heads per layer) (Supporting Information Figure S4). As shown in Figures 1, 2d, and Supporting Information S5, we observed that the attention was focused on the interacting sites of peptides, suggesting that this model was able to capture the structural information on peptide sequences. Furthermore, analysis of 144 attention mechanisms (Supporting Information Figure S4) showed that several learning mechanisms existed in the model, which include previous-word attention patterns (Supporting Information Figure S6a), next-word attention patterns (Supporting Information Figure S6b), delimiter-focused attention patterns (Supporting Information Figure S6c), specific-word attention patterns (Supporting Information Figure S6d,e), and related-word attention patterns (Supporting Information Figure S6f).

**Attentions of BERT-DPPIV Divert to DPP-IV Cleavage Sites and Proline.** In most cases, the attention learned by BERT-DPPIV focused on the second and third positions of the N terminus of the polypeptide sequences (Supporting Information Figure S6d–g), which are the cleavage sites of



**Figure 4.** Percentage of each attention head that is focused on the 20 natural amino acids.

DPP-IV.<sup>4</sup> This result suggested that our model can capture information about DPP-IV cleavage sites by learning the sequence of the polypeptide substrate.

To further illustrate the capability of cleavage site prediction, we statistically analyzed the attention sites of the sequences in the training set (Figure 3a,b). These results suggest that the attention of our model focused on the cleavage sites of DPP-IVs. Furthermore, a comparison of the changes in attention positions of the model before and after fine-tuning has supported the idea that our model could capture cleavage site information on DPP-IVs.

We further investigated the capability of our model to analyze the importance of each amino acid in the sequences. Statistical analyses have suggested that proline is the amino acid of major concern in our model, whereas cysteine is the amino acid of the least concern. An increase in attention toward proline and a decrease in attention toward cysteine were observed (Figure 3c,d).

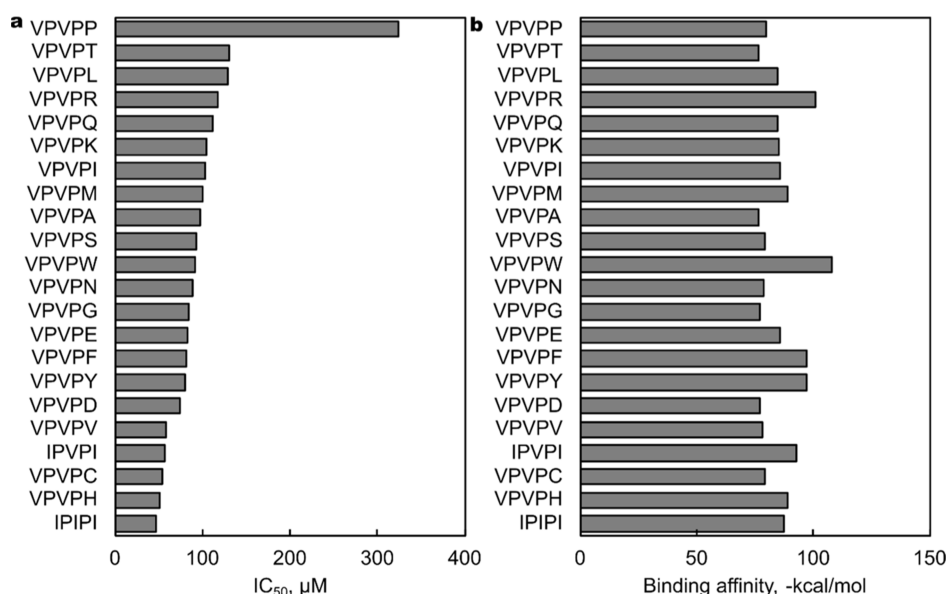
To investigate the importance of proline and cysteine in these amino acid sequences, we have visualized the attention of our model for each amino acid in a randomly selected sequence from  $\alpha$ -lactalbumin, “QNPHSSNICN”. We observed differences in all amino acids after training, which suggested that the model had directed its attention toward important amino acids, e.g., proline (which is an important amino acid for DPP-IV inhibition), whereas less important amino acids, e.g., cysteine, received less attention (Figure 3e–h). To illustrate the accuracy of this result, we have statistically analyzed the proportion of our model’s attention that was focused on each of the 20 standard amino acids (Figure 4) and compared it to the proportion of the pretrain model’s attention (Supporting Information Table S2). These results demonstrated that our model efficiently captured the amino acid types that play important roles in the amino acid sequences of a peptide.

**Application of BERT–DPP-IV Inhibitory Tripeptide Screening.** To illustrate the usefulness of

Table 1. Experimental Characterization of Tripeptides That Inhibit DPP-IV<sup>a</sup>

peptide	IC <sub>50</sub> , μM	peptide	IC <sub>50</sub> , μM	peptide	IC <sub>50</sub> , μM	peptide	IC <sub>50</sub> , μM
WRM	896.67	WPE	353.02	WAC	282.79	WPQ	196.18
WRP	741.50	WAG	352.16	WRH	277.20	WVE	192.50
WPC	739.03	WVD	343.00	WVV	275.33	WAH	189.55
WAR	664.68	WVA	341.25	WAA	274.53	WVF	188.33
WRW	603.14	WVG	340.63	WRG	269.85	WPW	182.96
WPS	588.31	WVT	335.63	WRS	265.09	WRE	160.15
WRT	555.60	WVK	331.38	WAL	260.10	WAF	153.32
WRV	541.17	WRI	330.4	WRF	257.07	WRL	150.75
WRA	534.67	WAM	329.85	WVH	256.86	WPG	134.98
WPV	511.47	WVM	317.00	WAV	249.79	WPN	128.59
WPD	480.61	WPH	314.57	WVS	247.00	WAY	117.40
WAE	461.79	WRN	313.20	WVW	238.63	WAW	103.66
WPL	461.32	WPY	303.79	WVC	231.78	WPI	
WAP	440.97	WPP	303.49	WPM	218.75	WPK	
WPT	435.08	WAN	301.78	WVY	214.88	WAK	
WAD	432.65	WRY	296.50	WRC	209.64	WAS	
WVP	402.00	WAT	296.33	WPR	208.94	WRK	
WVL	398.88	WAQ	289.19	WPA	206.01	WRR	
WVN	372.88	WPF	285.91	WRD	203.64	WVI	
WVQ	366.00	WRQ	285.50	WAI	198.88	WVR	

<sup>a</sup>IC<sub>50</sub>: half maximal inhibitory concentration. -: undetectable.



**Figure 5.** Characteristics of the activity of the DPP-IV inhibitory pentapeptides. The pentapeptides were characterized by (a) IC<sub>50</sub> and (b) binding affinity.

BERT–DPP-IV in the mining of DPP-IV-IPs, which can be used as antidiabetic drugs, we applied our model for the screening of DPP-IV inhibitory tripeptides. Prior work conducted by one of the authors (C.G.) classified DPP-IV inhibitory dipeptides into five classes and, from one class of peptides represented by VPX and IPX (where X represents one of the 20 amino acids), isolated nine tripeptides with efficient human DPP-IV (hDPP-IV) inhibitory activity.<sup>27</sup> However, it was unknown whether other kinds of tripeptides could efficiently inhibit hDPP-IV. In this section, one class of tripeptides that contains WPX, WAX, WRX, and WVX from five classes previously classified was selected for screening.

We predicted the DPP-IV inhibitory activity of the 80 kinds of tripeptides with our models and found that all of these

tripeptides are expected to possess DPP-IV inhibitory activity (Supporting Information Table S3). Then, these tripeptides were chemically synthesized, and the inhibitory activity was verified experimentally based on the IC<sub>50</sub>, which is inversely correlated to inhibitory activity. IC<sub>50</sub> values were detected for 72 out of the 80 kinds of tripeptides, and the DPP-IV inhibitory activity of each of these tripeptides was confirmed (Table 1). The Acc achieved in this study was 90%, which is consistent with our prediction (Supporting Information Table S1). The top three tripeptides for DPP-IV inhibitory activity were WAW, WAY, and WPN, with IC<sub>50</sub> values of 103.66, 117.40, and 128.59 μM, respectively. This result demonstrated that our model could predict DPP-IV-IPs with 90% accuracy.

**Proposed Strategy to Design DPP-IV-IPs Based on the Repeating Dipeptide Unit X-Proline.** Recently, the development of protein synthesis technology has enabled the screening of DPP-IV inhibitory activity among chemically synthesized peptides. These studies were conducted based on the property of preferential cleaving of X-proline or X-alanine from the N terminus of peptides by DPP-IV. However, the feasibility of designing the repeating dipeptide unit X-proline remained unclear. To address this problem, we used our model to investigate the potential of this design strategy.

As VPI has the highest hDPP-IV inhibitory activity,<sup>27</sup> the dipeptide unit VP was selected for further study. A library of 20 pentapeptides containing two VP repeats (VPVPX) was designed, and the inhibitory activity of these 20 kinds of pentapeptides was predicted by our model. In addition, IPIPI and IPVPI were also predicted by our model, as IPI has the highest porcine DPP-IV inhibitory activity.<sup>27</sup> The model predicted that 14 out of 22 pentapeptides exhibit inhibitory activity against DPP-IV, which has provided us with new insights into the properties of pentapeptides capable of inhibiting DPP-IV (Supporting Information Table S4).

To verify the properties of the peptides predicted, these peptides were chemically synthesized, and their IC<sub>50</sub> values were measured (Figure 5a,b, Supporting Information Table S5). Eight pentapeptides were identified as false negatives, whereas the IC<sub>50</sub> for the other peptides was correctly predicted. The top three peptides with the highest IC<sub>50</sub> were IPIPI, VPVPH, and VPVPC, which were 47.47, 51.11, and 54.77  $\mu$ M, respectively. The inhibitory activity of IPIPI was 2.18-fold higher than that of WAW, which is the tripeptide with the highest inhibitory activity detected in this study. This result indicates that pentapeptides with the repeating dipeptide unit VP exhibit higher DPP-IV inhibitory activity as compared with tripeptides and thus suggests that designing peptides with repeating units is a more efficient strategy.

To understand the DPP-IV inhibitory activity of these pentapeptides, their binding energy, which is inversely proportional to inhibitory activity, was analyzed with MDockPeP.<sup>28</sup> We calculated the binding affinities of pentapeptides to hDPP-IV and found that VPVPW, VPVPR, and VPVPY had the lowest affinity. However, a difference was observed in their ranking with respect to binding affinity relative to that for IC<sub>50</sub> (Supporting Information Table S5). This difference might have resulted from (1) the docking method, which reflects only the first step of the interaction between the polypeptide and DPP-IV and cannot reflect the subsequent two consecutive VP cleavage steps, or (2) the docking affinity being calculated based on DPP-IV in the static state, which thus did not consider the dynamic of DPP-IV.

These results suggested that a design strategy based on repeating X-proline units is effective and feasible, although the efficiency of peptides with repeating units of X-alanine should be investigated.

## DISCUSSION

In this study, we have successfully built a PLM-based model, BERT-DPPIV, for mining novel DPP-IV-IPs and applied it to guide biological experiments for screening and design. Our model, BERT-DPPIV, demonstrated an Acc of 0.894, which is higher than those previously reported, e.g., StackDPPIV.<sup>16</sup> Furthermore, this model can also capture more biologically relevant information from peptide sequences, notably the cleavage site information on DPP-IV and structure information

on peptide sequences (Figure 2a–d, Supporting Information Figure S5). Our model can automatically capture information from various aspects of sequence, structural feature information from the peptide sequence, and specific cleavage site features from the peptide data set. This renders our model more advantageous than StackDPPIV, which has been the best model described to date, as the latter must extract sequence and structure feature information by feature engineering.

Although our model achieved good performance, there is room for further improvement. As BERT models that use the word segmentation method with overlap, e.g., DNABERT,<sup>29</sup> typically show higher accuracy than those that do not, our model, which is based on the nonoverlapping word segmentation method, can be further improved by using an overlapping segmentation method.

The attention-based BERT model adopted in this study overcomes the poor interpretability of DL methods. Visualizing the model's attention during learning revealed what the model learned and how it had changed during training. Through this strategy, we discovered that our model can capture cleavage site information for DPP-IV in addition to the physicochemical properties and structural information on the peptide sequence (Figure 2a–d, Supporting Information Figure S5), which had not previously been reported as a classifier. This is the first use of the PLM BERT to discover enzyme cleavage site information from substrate information, and this model should be extremely helpful in understanding the biological implications of DPP-IV. Apparently, this feature can be extended to other enzymes after pretraining with respective databases, and it should be applicable to understanding the functions of different enzymes and to mining the cleavage site information on enzymes.

In contrast with the previous classifiers for inhibiting DPP-IV peptides<sup>14–16</sup> that were tested only on a data set and were not verified by actual screening experiments, we used the proposed model to screen tripeptides with DPP-IV inhibitory activity and showed experimentally that our model had an accuracy of 90%. DPP-IV inhibitory activity was identified in 72 out of 80 tripeptides predicted (Table 1). This result suggests that our model can be used to aid in the screening of functional peptides that inhibit DPP-IV.

The number of reported DPP-IV inhibitory polypeptides is still limited due to the absence of strategies for designing DPP-IV inhibitory polypeptide libraries. Therefore, we have also proposed a new peptide design strategy to design polypeptides that inhibit DPP-IV based on repeating dipeptide units to accelerate the mining of DPP-IV inhibitory polypeptides. This design strategy has been verified and shown to be feasible through the predictions of our model and subsequent biological assays. Twenty two pentapeptides containing VP or IP were discovered to inhibit DPP-IV efficiently. This design strategy can help researchers efficiently explore peptides that inhibit DPP-IV. The accuracy of the model for pentapeptide prediction was 63.6%, despite the lack of sequences with repeating dipeptide units in the training data set. It is likely that the model can be improved by increasing the availability of data in biological assays. Verification of the efficiency of the model in predicting long peptide sequences is in progress.

In conclusion, we constructed the first PLM-based DPP-IV classifier and verified its performance experimentally. Furthermore, based on the assistance of the classifier, 72 peptides were revealed to inhibit DPP-IV. A novel design strategy for peptides with DPP-IV inhibitory activity was proposed and

verified, based on which 22 pentapeptides were discovered to have DPP-IV inhibitory effects. This model could aid in the screening and design of DPP-IV-IPs, which will greatly accelerate the process and reduce the cost of antidiabetic drug development. However, the Acc of the model also needs to be improved by constructing a big data set; for example, the screened DPP-IV-IPs in this work should be integrated into the data set, and the IC<sub>50</sub> prediction model should be developed so that more potent DPP-IV-IPs can be found.

## METHODS

**Database.** The pretrained protein data were downloaded from UniProt and can be accessed on Google Drive at [https://drive.google.com/file/d/1QeXWV5\\_OIKgms7u5ShNfvPEKPBLOC3UT/view?usp=sharing](https://drive.google.com/file/d/1QeXWV5_OIKgms7u5ShNfvPEKPBLOC3UT/view?usp=sharing). This data set contains 556,603 protein sequences.<sup>30</sup> The data set used for fine-tuning was previously described and was used to train and test our proposed model.<sup>14</sup> The benchmark data set includes 532 DPP-IV-IPs and the same number of peptides with no inhibitory activity against DPP-IV. The number of IPs and noninhibitory peptides in the independent data set was 133 peptides each.

**Reagents and Peptides.** hDPP-IV (>200 U/mL) was obtained from the ATGen company (South Korea). Gly-Pro-p-nitroanilide (Gly-Pro-pNA) was obtained from Cayman Chemical Co. (Michigan, US). Peptides containing VPVPX, WPX, WAX, WRX, and WVX (where X represents any one of 20 amino acids) with purity >95% were synthesized by Gen-Script (Suzhou, China).

**Training of the PLM.** BERT<sup>31</sup> is a pretrained language model developed by Google for natural language text applications and achieves state-of-the-art results on downstream NLP tasks through transfer learning. We pretrained 12-layer BERT-based language models on 556,603 protein sequences from the UniProt data set.<sup>30</sup>

To adapt our model to different sequence lengths, we used different word segmentation lengths,  $k$ -mer = 1,  $k$ -mer = 2, and  $k$ -mer = 3, to generate three pretrained BERT models. The pretraining process consisted of two pretraining tasks: MLM and NSP, so that our pretrained BERT models could capture word-level and sentence-level representations and learn the common features of protein sequences. MLM is able to model complex relationships between amino acids and capture evolutionary information on proteins.<sup>32</sup> During the MLM task, 15% of the tokens among the original protein sequence were randomly masked to obtain the masked sequence, and a special token [MASK] was used to replace the masked token 80% of the time; a random token was used 10% of the time, and the selected token was unchanged 10% of the time. The model then predicted the masked tokens based on the context of the unmasked sequence for language modeling. For training, we minimized the negative log likelihood of the true amino acid at each of the masked positions using eq 1

$$L_{\text{MLM}} = - \sum_{\hat{x} \in m(x)} \log q(\hat{x} | x_{m(x)}) \quad (1)$$

where  $L_{\text{MLM}}$  is the loss function of MLM,  $m(x)$  and  $x_{m(x)}$  denote the masked tokens from the protein sequences  $x$  and the remaining tokens, respectively. During the NSP task, the data were randomly divided into two equal parts A and B: choosing the protein sequence segments A and B for each pretraining example, 50% of the time B is the actual next

sentence that follows A (labeled as IsNext), and 50% of the time it is a random sentence from the corpus (labeled as NotNext). BERT was trained by identifying whether these protein sequence segment pairs were continuous. The loss function of the NSP task,  $L_{\text{NSP}}$ , was described as eq 2

$$L_{\text{NSP}} = - \log p(t|x, y) \quad (2)$$

where  $t = 1$  if  $x$  and  $y$  are continuous protein segments from the protein corpus.

Then MLM and NSP were trained together with the goal of minimizing the combined loss function of the two strategies, as defined in eq 3.

$$L = L_{\text{MLM}} + L_{\text{NSP}} \quad (3)$$

The pretraining hyperparameters included train steps of 10 million times, a learning rate of  $2 \times 10^{-5}$  and a batch size of 32 to train the BERT model.<sup>30</sup> After the pretraining process, we obtained three pretrained BERT models. To construct the PLM for a specific downstream task that identifies and predicts DPP-IV-IPs, we modified the pretrained BERT models by adding a classification layer on top of the BERT output for the [CLS] token. We fine-tuned the pretrained model on a benchmark data set containing DPP-IV-IPs without major architectural modifications. For the fine-tuning hyperparameters, we used a learning rate of  $2 \times 10^{-6}$ , a batch size of 32, a warm-up proportion of 0.1, and an average training time of 50 epochs.

**Performance Evaluation of the Model.** We used four general quantitative indicators to evaluate our model: Acc, sensitivity (Sn), specificity (Sp), and MCC, each of which is defined by eqs 4, 5, 6, and 7, respectively. TP (true positive) is the number of correctly predicted DPP-IV-IPs, FN (false negative) is the number of DPP-IV-IPs that were in fact predicted to be non-DPP-IV-IPs, true negative (TN) is the number of correctly predicted non-DPP-IV-IPs, and false positive (FP) is the number of non-DPP-IV-IPs that were in fact predicted to be DPP-IV-IPs.

$$\text{Acc} = \frac{\text{TP} + \text{TN}}{(\text{TP} + \text{FP} + \text{TN} + \text{FN})} \quad (4)$$

$$\text{Sn} = \frac{\text{TP}}{(\text{TP} + \text{FN})} \quad (5)$$

$$\text{Sp} = \frac{\text{TN}}{(\text{TN} + \text{FP})} \quad (6)$$

$$\text{MCC} = \frac{\text{TP} \times \text{TN} - \text{FP} \times \text{FN}}{\sqrt{(\text{TP} + \text{FP}) \times (\text{TN} + \text{FN}) \times (\text{TP} + \text{FN}) \times (\text{TN} + \text{FP})}} \quad (7)$$

Sn and Sp reflect the model's ability to recognize DPP-IV-IPs and non-DPP-IV-IPs, respectively, and Acc embodies the overall prediction effect of the model. The value range of the three is [0,1], and the larger the value, the more accurate the model's prediction. MCC is usually considered a balanced indicator and can be used, even if the sample is not balanced. Its value is between -1 and +1, reflecting the correlation between the true label of the sample in the testing set and the predicted result. The higher value indicates a greater correlation. When the value is close to +1, the classification performance of the model is excellent; when it is close to -1, the prediction result of the model is the opposite of the actual result; and when it is close to 0, the prediction result of the



model is similar to a random prediction. In addition, the area under the receiver operating characteristic (AUC) was used as another statistical metric. Considering these five evaluation indicators, the performance of the classification models can be better evaluated.

**Visualization of the Model.** Representation visualization was used to analyze the learning capacity of the model. To the original protein sequence,  $x = [x_1, \dots, x_n]$ , we added the special start and end tokens, CLS and SEP, respectively, and got the final input sequence  $x' = [\text{CLS}, x_1, \dots, x_n, \text{SEP}]$ . The resulting sequence representation,  $h_i$  was obtained from hidden states in the  $i$ th layer of BERT-DPPIV.

Analyzing the information contained in the representation learned by the model, we used the vector corresponding to CLS in each layer of BERT-DPPIV to represent the analytical sequence and analyzed the representation at the residue level and sequence level by  $t$ -SNE. At the residue level, amino acids were divided into six categories, consisting of aromatic (W, F, Y), aliphatic (M, L, I, A, V), positive (R, H, K), negative (D, E), polar neutral (Q, N, S, T), and special-case (G, P, C) residues. At the sequence level, 10 kinds of physicochemical property information were obtained from the modLAMP package.<sup>22</sup>

Attention visualization was used to analyze the information about which the model was specifically concerned and provide interpretable analysis. Each attention head in a model layer produces an attention matrix,  $\alpha$ , which indicates the degree of correlation between token pairs. For example,  $\alpha_{ij}$  indicates the attention from token  $i$  to token  $j$ , and the weight of token  $i$  for all tokens is 1, as shown in eq 8.

$$\sum_{j=1}^n \alpha_{i,j} = 1 \quad (8)$$

We analyzed attention through a multiscale visualization tool for the Transformer model.<sup>24</sup> The attention-head view in this tool expresses the self-attention matrix  $\alpha$  in the form of connected lines and displays the attention patterns generated by one or more attention heads in a given layer. The lines in the head view indicate how much of the hidden state information on the attending token (right) will flow to the attended token (left) (Figure 2d). The different colors of the lines indicate different attention heads, whereas the color depth of the line is related to the attention weight. The model view in the tool provides a global view of attention patterns across all layers and heads of the model. We then use a slightly simplified version of AlphaFold<sup>26</sup> to obtain structure information on the polypeptide sequences and explore the relationship between sites of interaction in polypeptide structures and model attention patterns.

We then used the attention matrix to calculate the importance of each token,  $d_j$ , in the peptide sequence, as shown in eq 9.

$$d_j = \sum_{i=1}^n \alpha_{i,j} \quad (9)$$

We carried out a statistical analysis of 144 attention patterns, including all layers and all heads of each layer. We removed sequences of less than or equal to three amino acids. Then, we calculated the importance of all tokens, including the special tokens (CLS, SEP), and selected the three most important tokens from each sequence for each attention pattern. We

counted the distribution of the top three most important tokens across the different patterns for each sequence, with respect to their position and amino acid type. We also selected a specific residue in the sequence, calculated the importance of this residue in 144 attention patterns, and visualized the results using a heatmap. We then investigated the interaction between attention and particular amino acids.<sup>24</sup> We used eq 10 for amino acids and defined an indicator function,  $f(i, j)$ , which returns a value of 1 if this amino acid was present in the token  $j$  (e.g., if token  $j$  is a proline).

$$p_\alpha(f) = \frac{\sum_{x \in X} \sum_{i=1}^{|\text{CLS}|} \sum_{j=1}^{|\text{CLS}|} f(i, j) \cdot 1_{\alpha_{ij} > \theta}}{\sum_{x \in X} \sum_{i=1}^{|\text{CLS}|} \sum_{j=1}^{|\text{CLS}|} 1_{\alpha_{ij} > \theta}} \quad (10)$$

where  $p_\alpha(f)$  equals the proportion of attention directed to the amino acid. We computed the proportion of attention for the fine-tuned model and the pretrained model directed toward each of the 20 standard amino acids.

**DPP-IV Inhibitory Assay.** DPP-IV inhibitory activity was measured according to a slightly modified version of the method described previously.<sup>8</sup> hDPP-IV was used for the triplicate measurements. Peptide solution (20  $\mu\text{L}$ ); substrate solution, consisting of 2 mM Gly-Pro-pNA (20  $\mu\text{L}$ ); and Tris-HCl buffer (pH 8.0; 20  $\mu\text{L}$ ) were premixed in a 96-well microplate. The reaction was initiated by adding 40  $\mu\text{L}$  of hDPP-IV (final concentration, 0.025 U/ml) in 100  $\mu\text{L}$  of the above mixture. The samples were then incubated at 37  $^\circ\text{C}$  for 60 min. Absorbance at 405 nm was measured by using a microplate reader. The negative control mixture included Tris-HCl buffer (pH 8.0), Gly-Pro-pNA, hDPP-IV, and phosphate-buffered saline (pH 7.4). The hDPP-IV inhibitory ratio of the sample was calculated according to eq 11

$$\text{inhibition ratio} = \frac{A_{\text{control}} - A_{\text{sample}}}{A_{\text{control}}} \quad (11)$$

where  $A_{\text{control}}$  is the absorbance of the control sample with PBS, and  $A_{\text{sample}}$  is the absorbance of the peptide sample.

## ■ ASSOCIATED CONTENT

### Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acsomega.3c05571>.

Visualization of the data representation at the beginning and end of our model training, analysis of physicochemical properties of sequences based on model-based vector representation, structure prediction of peptides, attention-head view of YPSKPDNPGE, attentional visualization of structural information learned by the model, visualizing the model's attention patterns, classification results of BERT-DPPIV and other predictors, amino acids and the corresponding maximally attentive heads in the pre-trained model and fine-tuned model, predicted results of tripeptides and pentapeptides inhibiting DPP-IV, and characteristic results of pentapeptides inhibiting DPP-IV (PDF)

Model, Bert Embedding, Token embedding, Segment embedding, Positional embedding, BERT encoder, multi-head attention, feedforward, layer norm and dropout, pre-training process, visualization, statistical analysis of the attention matrix, and interaction between attention and particular amino acid (PDF)

Visualization of the training process (MP4)

## AUTHOR INFORMATION

### Corresponding Authors

**Change Guan** – Key Laboratory for Industrial Biocatalysis, Ministry of Education of China, Department of Chemical Engineering, Tsinghua University, Beijing 100084, China; [orcid.org/0000-0003-2531-6821](https://orcid.org/0000-0003-2531-6821); Email: [Change.Guan@Penmedicine.upenn.edu](mailto:Change.Guan@Penmedicine.upenn.edu)

**Xin-Hui Xing** – Key Laboratory for Industrial Biocatalysis, Ministry of Education of China, Department of Chemical Engineering and Center for Synthetic and Systems Biology, Tsinghua University, Beijing 100084, China; Institute of Biopharmaceutical and Health Engineering, Tsinghua Shenzhen International Graduate School, Shenzhen 518055, China; Institute of Biomedical Health Technology and Engineering, Shenzhen Bay Laboratory, Shenzhen 518118, China; Email: [xhxing@mail.tsinghua.edu.cn](mailto:xhxing@mail.tsinghua.edu.cn)

### Authors

**Jiawei Luo** – Department of Computer Science and Technology, Harbin Institute of Technology, Shenzhen 518055, China

**Shucheng Li** – Key Laboratory for Industrial Biocatalysis, Ministry of Education of China, Department of Chemical Engineering, Tsinghua University, Beijing 100084, China

**Zheng Lin Tan** – School of Life Science and Technology, Tokyo Institute of Technology, Yokohama, Kanagawa Prefecture 226-0026, Japan; [orcid.org/0000-0001-6447-3788](https://orcid.org/0000-0001-6447-3788)

**Yi Wang** – Key Laboratory for Industrial Biocatalysis, Ministry of Education of China, Department of Chemical Engineering, Tsinghua University, Beijing 100084, China

**Haihong Chen** – Institute of Biopharmaceutical and Health Engineering, Tsinghua Shenzhen International Graduate School, Shenzhen 518055, China; Institute of Biomedical Health Technology and Engineering, Shenzhen Bay Laboratory, Shenzhen 518118, China

**Naoyuki Yamamoto** – School of Life Science and Technology, Tokyo Institute of Technology, Yokohama, Kanagawa Prefecture 226-0026, Japan; [orcid.org/0000-0001-6546-1404](https://orcid.org/0000-0001-6546-1404)

**Chong Zhang** – Key Laboratory for Industrial Biocatalysis, Ministry of Education of China, Department of Chemical Engineering and Center for Synthetic and Systems Biology, Tsinghua University, Beijing 100084, China; [orcid.org/0000-0001-9609-8855](https://orcid.org/0000-0001-9609-8855)

**Yuan Lu** – Key Laboratory for Industrial Biocatalysis, Ministry of Education of China, Department of Chemical Engineering, Tsinghua University, Beijing 100084, China; [orcid.org/0000-0003-4500-6230](https://orcid.org/0000-0003-4500-6230)

**Junjie Chen** – Department of Computer Science and Technology, Harbin Institute of Technology, Shenzhen 518055, China

Complete contact information is available at: <https://pubs.acs.org/10.1021/acsomega.3c05571>

### Author Contributions

<sup>#</sup>C.G., J.L., S.L., and Z.L.T. contributed equally to this work. X.-H.X., N.Y., Y.L., and C.G. conceived the project and designed the study. C.G. designed the tripeptide and pentapeptide screening experiments. J.C. and C.G. designed the model construction and analysis; C.G. and S.L. performed wet experiments. C.G. and J.Luo performed model construction. C.G. and Z.L.T. performed data analysis and

interpretation. C.G., J.Luo, J.C., and Z.L.T. drafted and revised the paper. H.C., Y.W., and C.Z. provided consultation for the project. All authors contributed to the revision of the paper.

### Funding

This work was funded by the Shenzhen Science and Technology Innovation Commission (KCXFZ20201221173207022), Natural Science Foundation of China (62102118).

### Notes

The authors declare no competing financial interest.

## ACKNOWLEDGMENTS

The authors would like to thank Jiahao Li and Zourun Wu for their useful discussion.

## REFERENCES

- (1) Barnett, A. DPP-4 Inhibitors and Their Potential Role in the Management of Type 2 Diabetes: DPP 4 INHIBITORS AND MANAGEMENT OF TYPE 2 DIABETES. *Int. J. Clin. Pract.* **2006**, *60* (11), 1454–1470.
- (2) Nongonierma, A. B.; FitzGerald, R. J. Features of Dipeptidyl Peptidase IV (DPP-IV) Inhibitory Peptides from Dietary Proteins. *J. Food Biochem.* **2019**, *43* (1), No. e12451.
- (3) Wang, Y.; Li, S.; Guan, C.; He, D.; Liao, X.; Wang, Y.; Chen, H.; Zhang, C.; Xing, X.-H. Functional discovery and production technology for natural bioactive peptides. *Chin. J. Biotechnol.* **2021**, *37* (6), 2166–2180.
- (4) Mulvihill, E. E.; Drucker, D. J. Pharmacology, Physiology, and Mechanisms of Action of Dipeptidyl Peptidase-4 Inhibitors. *Endocr. Rev.* **2014**, *35* (6), 992–1019.
- (5) Duez, H.; Cariou, B.; Staels, B. DPP-4 Inhibitors in the Treatment of Type 2 Diabetes. *Biochem. Pharmacol.* **2012**, *83* (7), 823–832.
- (6) Nongonierma, A. B.; FitzGerald, R. J. An in Silico Model to Predict the Potential of Dietary Proteins as Sources of Dipeptidyl Peptidase IV (DPP-IV) Inhibitory Peptides. *Food Chem.* **2014**, *165*, 489–498.
- (7) Nongonierma, A. B.; Mooney, C.; Shields, D. C.; FitzGerald, R. J. In silico approaches to predict the potential of milk protein-derived peptides as dipeptidyl peptidase IV (DPP-IV) inhibitors. *Peptides* **2014**, *57*, 43–51.
- (8) Nongonierma, A. B.; FitzGerald, R. J. Structure Activity Relationship Modelling of Milk Protein-Derived Peptides with Dipeptidyl Peptidase IV (DPP-IV) Inhibitory Activity. *Peptides* **2016**, *79*, 1–7.
- (9) Nongonierma, A. B.; Paoella, S.; Mudgil, P.; Maqsood, S.; FitzGerald, R. J. Identification of Novel Dipeptidyl Peptidase IV (DPP-IV) Inhibitory Peptides in Camel Milk Protein Hydrolysates. *Food Chem.* **2018**, *244*, 340–348.
- (10) Nongonierma, A. B.; FitzGerald, R. J. Learnings from Quantitative Structure–Activity Relationship (QSAR) Studies with Respect to Food Protein-Derived Bioactive Peptides: A Review. *RSC Adv.* **2016**, *6* (79), 75400–75413.
- (11) Shoombuatong, W.; Prathipati, P.; Owasirikul, W.; Worachartcheewan, A.; Simeon, S.; Anuwongcharoen, N.; Wikberg, J. E. S.; Nantasenamat, C. Towards the Revival of Interpretable QSAR Models. In *Advances in QSAR Modeling*; Roy, K., Ed.; Challenges and Advances in Computational Chemistry and Physics; Springer International Publishing: Cham, 2017; Vol. 24, pp 3–55.
- (12) Hellberg, S.; Sjoestroem, M.; Skagerberg, B.; Wold, S. Peptide Quantitative Structure-Activity Relationships, a Multivariate Approach. *J. Med. Chem.* **1987**, *30* (7), 1126–1135.
- (13) Basith, S.; Manavalan, B.; Shin, T. H.; Lee, G. Machine Intelligence in Peptide Therapeutics: A Next-generation Tool for Rapid Disease Screening. *Med. Res. Rev.* **2020**, *40* (4), 1276–1314.
- (14) Charoenkwan, P.; Kanthawong, S.; Nantasenamat, C.; Hasan, M. M.; Shoombuatong, W. IDPP-IV-SCM: A Sequence-Based

Predictor for Identifying and Analyzing Dipeptidyl Peptidase IV (DPP-IV) Inhibitory Peptides Using a Scoring Card Method. *J. Proteome Res.* **2020**, *19* (10), 4125–4136.

(15) Zou, H.; Yin, Z. Identifying Dipeptidyl Peptidase-IV Inhibitory Peptides Based on Correlation Information of Physicochemical Properties. *Int. J. Pept. Res. Ther.* **2021**, *27* (4), 2651–2659.

(16) Charoenkwan, P.; Nantasenamat, C.; Hasan, M. M.; Moni, M. A.; Lio, P.; Manavalan, B.; Shoombuatong, W. StackDPP-IV: A Novel Computational Approach for Accurate Prediction of Dipeptidyl Peptidase IV (DPP-IV) Inhibitory Peptides. *Methods* **2022**, *204*, 189–198.

(17) Janiesch, C.; Zschech, P.; Heinrich, K. Machine learning and deep learning. *Machine Learning and Deep Learning. Electron. Mark.* **2021**, *31* (3), 685–695.

(18) Armenteros, J. J. A.; Tsirigos, K. D.; Sønderby, C. K.; Petersen, T. N.; Winther, O.; Brunak, S.; von Heijne, G.; Nielsen, H. SignalP 5.0 Improves Signal Peptide Predictions Using Deep Neural Networks. *Nat. Biotechnol.* **2019**, *37* (4), 420–423.

(19) Youmans, M.; Spainhour, C.; Qiu, P. Long Short-Term Memory Recurrent Neural Networks for Antibacterial Peptide Identification. In *2017 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*; IEEE: Kansas City, MO, 2017, pp 498–502.

(20) Yan, J.; Bhadra, P.; Li, A.; Sethiya, P.; Qin, L.; Tai, H. K.; Wong, K. H.; Siu, S. W. I. Deep-AmPEP30: Improve Short Antimicrobial Peptides Prediction with Deep Learning. *Mol. Ther. Nucleic Acids* **2020**, *20*, 882–894.

(21) Wu, C.; Gao, R.; Zhang, Y.; De Marinis, Y. PTPD: Predicting Therapeutic Peptides by Deep Learning and Word2vec. *BMC Bioinf.* **2019**, *20* (1), 456.

(22) Müller, A. T.; Gabernet, G.; Hiss, J. A.; Schneider, G. ModLAMP: Python for Antimicrobial Peptides. *Bioinformatics* **2017**, *33* (17), 2753–2755.

(23) Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; Polosukin, I. Attention Is All You Need. In *NIPS'17: Proceedings of the 31st International Conference on Neural Information Processing Systems*; Long Beach: USA, 2017, pp 6000–6010.

(24) Vig, J. A Multiscale Visualization of Attention in the Transformer Model. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*; Association for Computational Linguistics: Florence, Italy, 2019, pp 37–42.

(25) Timmons, P. B.; Hewage, C. M. APPTTEST Is a Novel Protocol for the Automatic Prediction of Peptide Tertiary Structures. *Briefings Bioinf.* **2021**, *22* (6), bbab308.

(26) Jumper, J.; Evans, R.; Pritzel, A.; Green, T.; Figurnov, M.; Ronneberger, O.; Tunyasuvunakool, K.; Bates, R.; Žídek, A.; Potapenko, A.; Bridgland, A.; Meyer, C.; Kohl, S. A. A.; Ballard, A. J.; Cowie, A.; Romera-Paredes, B.; Nikolov, S.; Jain, R.; Adler, J.; Back, T.; Petersen, S.; Reiman, D.; Clancy, E.; Zielinski, M.; Steinegger, M.; Pacholska, M.; Berghammer, T.; Bodenstern, S.; Silver, D.; Vinyals, O.; Senior, A. W.; Kavukcuoglu, K.; Kohli, P.; Hassabis, D. Highly Accurate Protein Structure Prediction with AlphaFold. *Nature* **2021**, *596* (7873), 583–589.

(27) Guan, C.; Iwatani, S.; Xing, X.; Yamamoto, N. Strategic Preparations of DPP-IV Inhibitory Peptides from Val-Pro-Xaa and Ile-Pro-Xaa Peptide Mixtures. *Int. J. Pept. Res. Ther.* **2021**, *27* (1), 735–743.

(28) Xu, X.; Yan, C.; Zou, X. MDockPeP: An Ab-initio Protein–Peptide Docking Server. *J. Comput. Chem.* **2018**, *39* (28), 2409–2413.

(29) Ji, Y.; Zhou, Z.; Liu, H.; Davuluri, R. V. DNABERT: Pre-Trained Bidirectional Encoder Representations from Transformers Model for DNA-Language in Genome. *Bioinformatics* **2021**, *37* (15), 2112–2120.

(30) Zhang, Y.; Lin, J.; Zhao, L.; Zeng, X.; Liu, X. A Novel Antibacterial Peptide Recognition Algorithm Based on BERT. *Briefings Bioinf.* **2021**, *22* (6), bbab200.

(31) Devlin, J.; Chang, M.-W.; Lee, K.; Toutanova, K. BERT: Pre-Training of Deep Bidirectional Transformers for Language Under-

standing. In *Proceedings of the 2019 Conference of the North American Association for Computational Linguistics: Minneapolis, MN, 2019*, pp 4171–4186.

(32) Elnaggar, A.; Heinzinger, M.; Dallago, C.; Rehawi, G.; Wang, Y.; Jones, L.; Gibbs, T.; Feher, T.; Angerer, C.; Steinegger, M.; Bhowmik, D.; Rost, B. ProtTrans: Toward Understanding the Language of Life Through Self-Supervised Learning. *IEEE Trans. Pattern Anal. Mach. Intell.* **2022**, *44*, 7112–7127.