

RESEARCH ARTICLE

MoCETSE: A mixture-of-convolutional experts and transformer-based model for predicting Gram-negative bacterial secreted effectors

Hua Shi^{1*}, Yihang Lin¹, Dachen Liu¹, Quan Zou^{2,3}

1 School of Opto-electronic and Communication Engineering, Xiamen University of Technology, Xiamen, Fujian, China, **2** Yangtze Delta Region Institute (Quzhou), University of Electronic Science and Technology of China, Quzhou, Zhejiang, China, **3** Institute of Fundamental and Frontier Sciences, University of Electronic Science and Technology of China, Chengdu, Sichuan, China

* shihua@xmut.edu.cn



Abstract

Identifying effector proteins of secretion systems in Gram-negative bacteria is crucial for deciphering their pathogenic mechanisms and guiding the development of antimicrobial strategies. Extracting evolutionary and sequence features using pre-trained protein language models (PLMs) has emerged as an effective approach to improve the performance of effector protein prediction. However, the high-dimensional features generated by PLMs contain extensive general biological information, making it difficult to focus on core features when applied directly to effector protein tasks, which in turn limits prediction performance. In this study, we propose MoCETSE, a deep learning model for predicting effector proteins in Gram-negative bacteria. Specifically, MoCETSE first extracts contextual representations of sequences using the pre-trained protein language model ESM-1b. Subsequently, it refines key functional features via a target preprocessing network to construct more expressive sequence representations. Finally, integrated with a transformer module incorporating relative positional encoding, MoCETSE explicitly models the relative spatial relationships between residues, enabling highly accurate prediction of secreted effector proteins. MoCETSE exhibits excellent and robust performance in both five-fold cross-validation and independent testing. Benchmark results demonstrate that it maintains strong competitiveness compared to existing binary and multi-class predictors. Additionally, the model can effectively perform genome-wide effector protein prediction, showing outstanding specificity and reliability. MoCETSE provides an efficient and robust computational framework for the accurate identification of bacterial effector substrates and offers key biological insights.

OPEN ACCESS

Citation: Shi H, Lin Y, Liu D, Zou Q (2026) MoCETSE: A mixture-of-convolutional experts and transformer-based model for predicting Gram-negative bacterial secreted effectors. *PLoS Comput Biol* 22(3): e1013397. <https://doi.org/10.1371/journal.pcbi.1013397>

Editor: Shugang Zhang, Ocean University of China, CHINA

Received: August 6, 2025

Accepted: March 3, 2026

Published: March 11, 2026

Copyright: © 2026 Shi et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data availability statement: The data, code, and models from this study are openly accessible on Github (<https://github.com/YihangLin123/MoCETSE>). This repository enables researchers to access the datasets, utilize the code, invoke the models developed, and conduct predictions.

Author summary

Secreted effector proteins are a class of key virulence factors in Gram-negative bacteria. After being injected into host cells, they interfere with normal cellular

Funding: This work was supported by the National Natural Science Foundation of China (No. 62372392 to HS). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

functions, leading to the development of diseases. Accurate identification of these virulence proteins is crucial for understanding bacterial pathogenic mechanisms and developing therapeutic strategies. However, existing methods suffer from issues such as feature redundancy and insufficient capture of long-range dependency signals. Here, we developed a novel computational framework called MoCETSE that enables end-to-end intelligent prediction of effector proteins directly from raw protein sequence information. The model leverages a pre-trained protein language model to extract deep biological information from raw sequences; a target preprocessing network then refines the extracted information to focus on features most relevant to effector protein identification. During the learning of secretion signal features, we introduced relative positional encoding to effectively capture associations between distant positions in the sequence. In cross-category prediction, MoCETSE outperformed tools such as DeepSecE. Furthermore, we provide interpretable biological mechanisms supporting the model, revealing which key sequence motifs and functional regions play core roles in distinguishing different types of effector proteins.

Introduction

Secretion systems are complex molecular apparatuses unique to Gram-negative pathogens. They provide a unique pathogenic mechanism by mediating the transmembrane transport of effector proteins for direct injection into the host cytoplasm [1,2]. Currently identified bacterial secretion systems are mainly classified into type I (T1SS), type II (T2SS), type III (T3SS), type IV (T4SS), and type VI (T6SS) based on their secretion mechanisms and molecular characteristics [3]. Although there is heterogeneity in the effector proteins carried by different pathogens, these secretion systems play an indispensable role in the infection process of various clinically important Gram-negative bacteria by regulating the delivery of key virulence factors [4,5]. For example, pathogenic bacteria represented by *Shigella* and *Citrobacter rodentium* can inject type III secreted effectors (T3SEs) into host cells through T3SS [6], directly disrupting the normal physiological activities of host cells. The pathogenicity of *Legionella pneumophila* depends on the establishment and maintenance of *Legionella-containing vacuoles* regulated by Dot/Icm T4SS [7]; this system specifically interferes with host vesicle transport, actin reorganization, and signal transduction by transporting T4SEs into the host cytoplasm. The T6SS of *Pseudomonas aeruginosa*, as an efficient protein transport apparatus, can directly deliver type VI secreted effectors (T6SEs) across the membrane to adjacent target cells [8], enhancing its pathogenic ability during infection. Therefore, predicting, identifying, and classifying bacterial secreted effectors and understanding the virulence processes of secreted effectors are of great significance for an in-depth understanding of microbial pathogenic mechanisms and the development of new anti-infective therapeutic strategies.

To date, a variety of predictive tools for secretion system effector proteins have been developed in the field, covering type I (T1SEs) [9], type III (T3SEs) [10–12],

type IV (T4SEs) [13–17], and type VI (T6SEs) [18,19] secreted effectors. Among these tools, Bastion3 and Bastion6 achieve accurate prediction of T3SEs and T6SEs, respectively, by integrating evolutionary conservation and structural constraint information characterized by the position-specific scoring matrix (PSSM) [10,18]. By contrast, CNN-T4SE adopts a multi-source feature fusion strategy, which significantly improves the recognition accuracy of T4SEs through integrating PSSM, protein secondary structure and solvent accessibility (PSSSA), as well as one-hot encoding [14]. However, the limited size of known datasets, coupled with high sequence homology and a lack of structural conservation among effector proteins, constitutes a major bottleneck for current predictive models. This renders it challenging for models to capture robust representative features based on small-scale datasets, thereby restricting the further improvement of their performance [20]. In response to these challenges, PLMs [21–23] have opened up a novel research avenue for effector protein prediction tasks. Built on the attention mechanism [24], these models mine potential biological patterns from datasets encompassing billions of protein sequences via deep neural networks in large-scale unsupervised learning tasks. PLMs are capable of converting raw amino acid sequences into distributed representations with profound biological significance [25,26], thereby capturing the structural, functional, sequence characteristics, and evolutionary relationships of proteins.

A growing body of research has demonstrated that applying pre-trained features derived from PLMs can significantly enhance model performance in secreted substrate recognition tasks. For instance, DeepSecE achieves accurate identification of effector proteins by integrating large-scale PLMs with a transformer module optimized for secretion signals [27]. Similarly, T4Seeker remarkably boosts the predictive capability for T4SEs through fusing PLM features, distance residue (DR) features, and long short-term memory (LSTM) networks [17]. However, while the high-dimensional features generated by PLMs encapsulate general biological semantic information, they inevitably contain a certain degree of information redundancy. Direct utilization of these features in effector classification tasks with diverse characteristics [28] often fails to precisely capture the task-specific key features, thereby limiting the model's predictive efficacy. Accordingly, the introduction of a dedicated preprocessing framework capable of task-oriented refinement of high-dimensional semantic embeddings and extraction of multi-scale local features is conducive to further improving model performance. Additionally, a hallmark feature of secreted proteins is the presence of a signal peptide, which is typically located at the N-terminus (the “classical” location) or occasionally at the C-terminus or internally (the “non-classical” locations) [3]. This signal peptide plays a crucial role in the secretion and translocation of proteins. Nevertheless, most existing methods fail to effectively capture information regarding this signal region, which not only constrains the performance of models in effector protein prediction but also compromises their interpretability.

To address these challenges, we constructed a target preprocessing network (TPN) designed to focus on the key sequence features embedded in the PLM outputs. In addition, we introduced a multi-head attention mechanism integrated with relative positional encoding to enhance the model's ability to capture the relative order and long-range functional dependencies among amino acid residues. Based on these designs, we developed a high-performance effector protein prediction model called MoCETSE, which is capable of accurately identifying five major types of secreted effector proteins (T1SE, T2SE, T3SE, T4SE, and T6SE). This model not only rivals current advanced binary classifiers (such as Bastion3 [10], CNN-T4SE [14], and Bastion6 [18]) in performance, but its prediction accuracy also surpasses the existing popular multi-class classifier DeepSecE [27]. Furthermore, MoCETSE can efficiently perform genome-wide predictive inference of secreted proteins. Our model provides a new perspective for extending the application of PLMs in bacterial effector protein prediction and is expected to advance our understanding of the pathogenic mechanisms associated with bacterial secretion systems.

Materials and methods

Data description

The non-redundant labeled training and independent test datasets used by the MoCETSE effector protein classifier were all derived from the multi-class effector protein prediction tool DeepSecE [27]. During the original construction of

DeepSecE dataset, the authors further applied CD-HIT (v4.8.1) [29] with a 60% sequence identity threshold to remove redundant and highly homologous sequences both within the training set and between the training and test sets. After obtaining the training dataset, we performed a homology redundancy check between our training dataset and the independent test dataset from Bastion3 [10], CNN-T4SE [14], and Bastion6 [18]. We used CD-HIT with a 90% sequence identity threshold for clustering analysis and removed the overlapping sequences between the training dataset and the benchmark test dataset to ensure the independence of the training dataset. Finally, the training dataset (Dataset S1) contains 1,577 non-effector proteins, 128 T1SEs, 68 T2SEs, 392 T3SEs, 507 T4SEs, and 232 T6SEs, totaling 2,904 samples; the test dataset (Dataset S2) contains 150 non-effector proteins, 20 T1SEs, 10 T2SEs, 30 T3SEs, 30 T4SEs, and 20 T6SEs, totaling 260 samples. The protein sequences in these training and independent test datasets span a wide range of lengths, from fewer than 100 amino acids to over 1,000 amino acids. As shown in S1 Fig, the lengths of these proteins are mainly distributed within the range of 100–600 amino acids.

During the benchmark testing phase, to comprehensively evaluate the performance of MoCETSE against existing mainstream binary and multi-class methods, we selected corresponding benchmark datasets for different types of secreted proteins. For type I secreted proteins, given the lack of publicly available independent test sets, we followed the methodology established in DeepSecE, selecting 20 T1SE samples and 150 non-secreted protein samples from the independent test dataset (Dataset S2) to construct a new benchmark dataset (Dataset S3). For type III, IV, and VI secreted proteins, we directly used the publicly available independent test datasets from Bastion3, CNN-T4SE, and Bastion6, which were designated as Dataset S4, Dataset S5, and Dataset S6, respectively. The samples in these test datasets were all curated by the original authors from primary literature, public databases, or experimental data (encompassing diverse bacterial species and studies). Additionally, the original authors applied CD-HIT with a sequence similarity threshold to remove duplicate or highly homologous samples from both the training and independent test datasets, preventing data leakage and ensuring dataset independence. The detailed specifications of the datasets used in this study are provided in S1 Table, and the distribution of different types of effector and non-effector proteins within the datasets is visually depicted in S2 Fig.

Overview framework of MoCETSE

Our proposed innovative framework, MoCETSE (Fig 1), starts with the input of effector protein sequence data. We leverage ESM-1b [21], a pre-trained protein language model, to encode amino acid sequences into informative feature vectors. Specifically, ESM-1b is constructed by stacking 33 layers of transformer modules, which alternate between self-attention mechanisms and feed-forward connections. Each transformer module contains multiple attention heads capable of capturing sequence information from different perspectives. We limit protein sequences to a maximum of 1,020 amino acids; these are converted into tokens and processed by ESM-1b, resulting in 1,280-dimensional token embeddings. To enhance the model's learning capacity for protein sequence features and training efficiency, a target preprocessing network (Fig 1B) is utilized. This network reduces the token embedding dimension from 1280 to 256 by leveraging eight convolutional experts that work collaboratively. Subsequently, the features are fed into a transformer module incorporating relative positional encoding (Fig 1C), configured with 4 attention heads and a feed-forward network employing GELU as the activation function [30]. After processing by the transformer module, the obtained feature representations pass through a fully connected layer for linear transformation, mapping them to a dimensional space corresponding to the number of classification categories. Finally, the softmax activation function is used to convert the output values into a probability distribution, ultimately achieving the classification of non-effector proteins and five types of effector proteins (T1SE–T4SE and T6SE):

$$P(y|x) = \text{softmax} \left(W^T \frac{1}{L} \sum_{i=1}^L \text{RPET}(\text{TPN}(\text{PLM}(x)))_i \right) \quad (1)$$

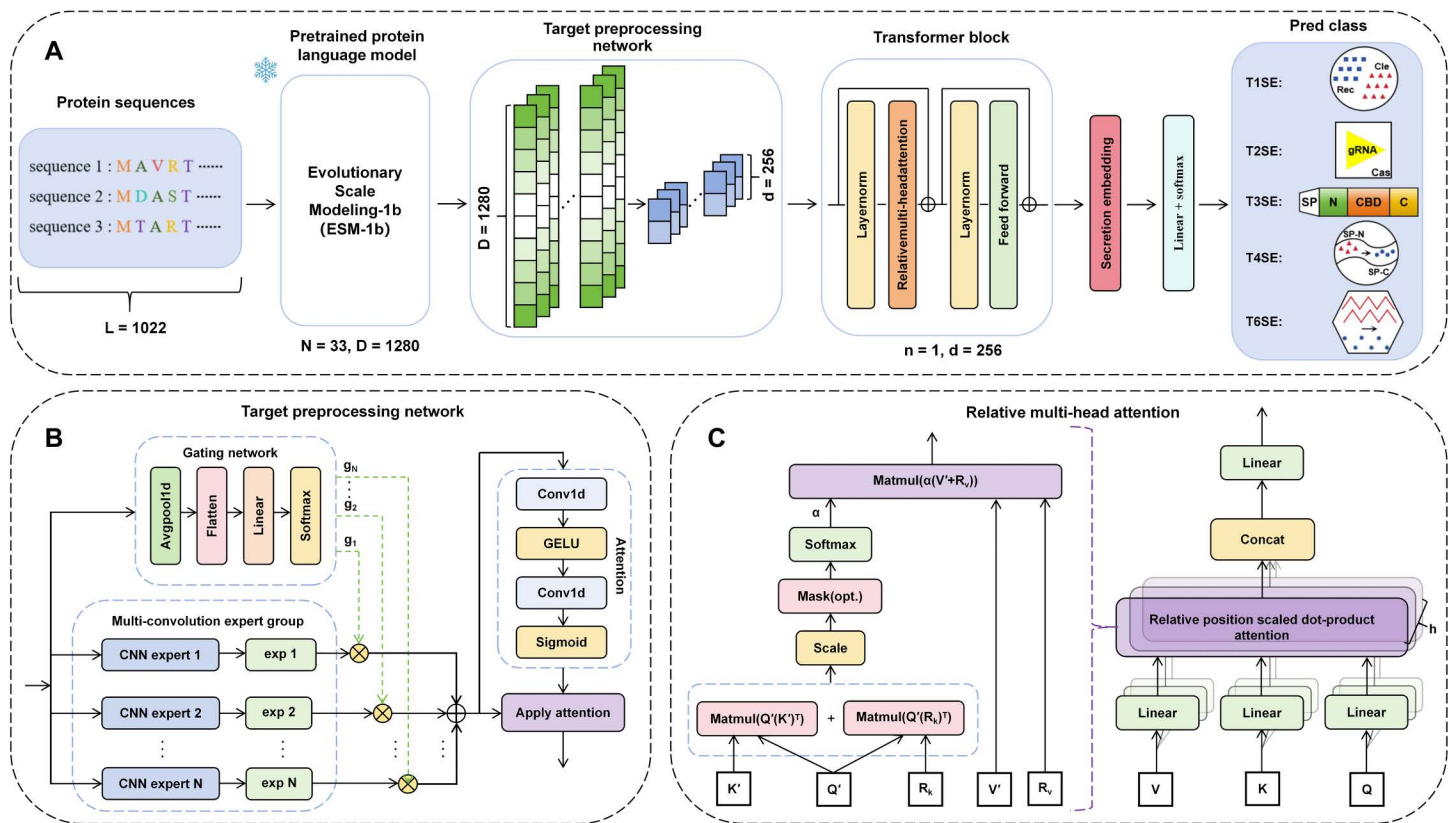


Fig 1. Overall workflow of MoCETSE. (A) MoCETSE uses the pre-trained protein language model ESM-1b, a target preprocessing network, and a transformer module to achieve effective learning of secreted protein features. The snowflake icon on the ESM-1b module indicates that its weights remain frozen during subsequent training. (B) The target preprocessing network is used to refine more critical discriminative features from secreted effector proteins. (C) The relative positional multi-head attention mechanism is employed to enhance the model's ability to identify key functional motifs in secreted effector proteins.

<https://doi.org/10.1371/journal.pcbi.1013397.g001>

Where x denotes the input protein sequence, and $P(y|x)$ represents the conditional probability of the predicted output y given the input sequence x , W^T denotes the weight matrix, which maps the model's output to the final class space, L indicates the sequence length, $(\cdot)_i$ represents the embedding vector of the i -th amino acid residue, PLM represents the pre-trained protein language model, TPN represents the target preprocessing network, and RPET represents the relative position-encoded transformer.

Protein language model

A multitude of studies have confirmed that the feature representations generated by protein language models can effectively capture essential sequence characteristics [31,32]. This advantage has led to their extensive application in predicting protein structures and functional properties. ESM-1b employs an unsupervised learning strategy to learn and extract various biological embedding features of proteins from the UniRef50 database [33], which contains approximately 250 million protein sequences. This database systematically integrates protein sequences from the UniProt Knowledgebase via a clustering algorithm, grouping sequences with at least 50% sequence similarity into the same cluster. The initial weights of ESM-1b were obtained from the pre-trained model parameter repository (<https://github.com/facebookresearch/esm>) and remained frozen without any fine-tuning during subsequent training.

Target preprocessing network

The target preprocessing network (TPN) adopts a “divide-and-conquer” strategy and is built upon a Mixture-of-Experts (MoE) architecture. Its core function is to preprocess high-dimensional embedding vectors generated by protein language models, facilitating dimensionality reduction while preserving and emphasizing key sequence features. By incorporating a multi-convolutional expert architecture, TPN synchronously captures critical features during the dimensionality reduction process, thereby further enhancing the specificity and representational power of feature extraction. The MoE framework, originally proposed by Jordan and Jacobs in 1991 [34], facilitates efficient modeling of complex data distributions by decomposing intricate tasks into sub-problems handled by specialized expert modules, with their outputs integrated via a gating mechanism. The efficacy of this paradigm has been validated by recent studies such as DeepSeekMoE [35], which achieves reduced parameter counts and competitive performance through refined expert partitioning, and DynamicMoE [36], which balances computational cost and precision using dynamic gating. As illustrated in Fig 1B, TPN in this study comprises three core components: the gating network module [37], the multi-convolutional expert group [38], and the attention module [24].

The gating network module is responsible for dynamically assigning weights to each convolutional expert module. By leveraging the specialized functions of different expert modules, it facilitates task decomposition and mitigates the parameter interference typically associated with global weight sharing in conventional single-network architectures. Specifically, the gating network generates gating weights g_1, g_2, \dots, g_N via a softmax activation function, constructing a probability distribution to adaptively regulate the contribution of each expert. The final output is derived as the weighted sum of the outputs from all experts, where the magnitude of the weights determines the relative influence of each expert. The output of the gating network module, $g(x)$, is defined as follows:

$$g(x) = \text{softmax}(x \cdot W_{gate}) \quad (2)$$

Where x represents the feature vector of raw input data after preprocessing via average pooling and flattening, and W_{gate} represents the learnable weight matrix.

The multi-convolutional expert group module integrates the MoE architecture with convolutional neural networks to construct an efficient and heterogeneous feature extraction system. This module consists of multiple parallel convolutional sub-modules, each serving as an independent expert branch with its own dedicated parameter set. This configuration enables the extraction of features across various spatial scales and dimensions from the input data. Guided by the gating network, which assigns weights to each expert branch, the module executes a weighted fusion of their outputs. This design leverages collaborative multi-expert integration to adaptively combine diverse features, thereby enhancing the TPN's capacity to model complex data distributions.

For a given input x , let $g_i(x)$ represent the output of the gating network, and let $E_i(x)$ represent the output of the i -th convolutional expert network. The output F_{expert} of the multi-convolutional expert group module can be written in the following form:

$$F_{\text{expert}} = \sum_{i=1}^n g_i(x) \cdot E_i(x) \quad (3)$$

To enhance the model's selective focus on high-value features, the target preprocessing network incorporates a lightweight attention module, whose core function is to achieve feature optimization and selection through dynamic weight adjustment. The module takes the fused features F_{expert} as input, which are synergistically generated by the multi-convolutional expert group and the gating network. Through convolutional operations, the GELU activation function [30], and sigmoid normalization, the module generates an element-wise attention mask α . Subsequently, an

attention-weighting mechanism is employed to perform element-wise multiplication between α and F_{expert} . This process adaptively highlights task-relevant key information while suppressing irrelevant features and noise, thereby improving the network's robustness in processing complex data. The final output of TPN, F_{output} , can be defined as follows:

$$F_{\text{output}} = \alpha \odot F_{\text{expert}} \quad (4)$$

Relative multi-head attention

The transformer model, proposed by Vaswani et al. (2017) [24], represents a fundamental shift from the conventional recurrent neural networks (RNNs) [39] and convolutional neural networks (CNNs) [40] typically employed in sequence learning tasks. This architecture entirely dispenses with recurrent structures and convolutional operations, instead relying on a parallelized self-attention mechanism to directly capture global dependencies within a sequence. By doing so, it effectively overcomes the inherent limitations of traditional models in long-range modeling and computational parallelization.

However, this architectural design possesses an inherent limitation: due to the permutation-invariant nature of the self-attention mechanism, the model is unable to perceive the positional order of elements without the inclusion of auxiliary positional information. To address this deficiency, attention-based models typically incorporate relative positional encoding or learnable positional embeddings. As illustrated in Fig 1C, this study integrates relative positional encoding [41,42] transformer's multi-head attention mechanism. This modification enables the MoCETSE model to more accurately capture the relative spatial dependencies between key motifs within effector protein sequences, while simultaneously enhancing the model's generalization capability across sequences of varying lengths. The mathematical formulation of the relative position-aware scaled dot-product attention is defined as follows:

$$\text{Attention}_{\text{rel}}(Q, K, V, R_k, R_v) = \text{softmax}\left(\frac{QK^T + QR_k^T}{\sqrt{d_{\text{key}}}}\right) \cdot (V + R_v) \quad (5)$$

Where Q , K , and V are the matrices comprising the query, key, and value vectors, respectively; d_{key} denotes the dimension of the K and V matrices. Additionally, R_k and R_v refer to the relative positional encoding matrices representing the positional bias and positional correction, respectively.

The relative multi-head attention mechanism incorporates relative positional information into the calculation of attention weights by decoupling content-based and position-based components. This allows the attention scores to simultaneously reflect both the content correlation between the query and key vectors and their relative displacement. The output of the relative multi-head attention mechanism is defined as follows:

$$\text{head}_i = \text{Attention}_{\text{rel}}\left(QW_i^Q, KW_i^K, VW_i^V, R_k W_i^{R_k}, R_v W_i^{R_v}\right) \quad (6)$$

$$\text{RelMultiHead}(Q, K, V, R_k, R_v) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O \quad (7)$$

Where W_i^Q , W_i^K , W_i^V , $W_i^{R_k}$, $W_i^{R_v}$ and W^O are learnable weight matrices for the i -th head.

Model training

In this study, the MoCETSE model was trained using a five-fold cross-validation strategy. The model was developed using the PyTorch (v1.10.0) framework and executed on an NVIDIA A100 GPU (40 GB) with Python 3.9.7 and CUDA 11.3. Model weights were initialized using the Xavier initialization method, and the Adam optimizer was employed for gradient

updates. The initial learning rate was set to 5×10^{-5} with a weight decay of 4×10^{-5} . The batch size and maximum number of epochs were set to 32 and 30, respectively. To prevent overfitting, an early stopping mechanism with a patience of 5 epochs was implemented by monitoring the F1 score. Within the transformer blocks, the dropout rates for the attention mechanism, multi-head attention layers, and feed-forward networks were configured as 0.05, 0.4, and 0.4, respectively.

Performance assessment

This study adopts a multi-metric approach to evaluate model performance in both cross-validation and independent testing. Specifically, we employed macro-averaged prediction accuracy, F1 score, area under the receiver operating characteristic curve (AUC), and area under the precision-recall curve (AUPRC) for comparison. The macro-averaged prediction accuracy is the arithmetic mean of the accuracies across all classes. This metric assigns equal weight to each class, effectively mitigating bias caused by class imbalance and fairly reflecting the model's classification ability on minority classes. For the multiclass classification task of predicting five types of secreted effectors, the receiver operating characteristic (ROC) curve provides an intuitive visualization of the model's ability to distinguish between positive and negative samples across different threshold settings for each category. Additionally, the multi-class confusion matrix provides a detailed view of classification results for each protein category. A higher AUC value (approaching 1) indicates superior discrimination, while a higher AUPRC value reflects a stronger capability to identify positive instances in precision-sensitive contexts. When comparing performance with existing binary and multi-class methods on benchmark datasets, we utilized several standard metrics, including accuracy (ACC), recall (REC), precision (PR), F1 score (F1), and Matthews correlation coefficient (MCC), to ensure a comprehensive assessment.

$$ACC = \frac{TP+TN}{TP+FP+TN+FN} \quad (8)$$

$$REC = \frac{TP}{TP+FN} \quad (9)$$

$$PR = \frac{TP}{TP+FP} \quad (10)$$

$$F1 = \frac{2 \times PR \times REC}{PR+REC} \quad (11)$$

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP+FP) \times (TP+FN) \times (TN+FP) \times (TN+FN)}} \quad (12)$$

Where TP, TN, FP, and FN represent the numbers of true positives, true negatives, false positives, and false negatives, respectively. All evaluation metrics were calculated using the scikit-learn library in Python.

Benchmark with existing popular methods

To ensure a fair comparison with existing mainstream prediction methods, we selected 11 widely used effector secretion system prediction tools as comparative models and conducted unified evaluations through their official online web servers. The access links to these tools are provided in [S2 Table](#).

In the comparative experiments, we selected benchmark test datasets corresponding to the prediction task types of each method and submitted them to the online web servers of the baseline methods to obtain the prediction results for performance comparison. All methods were evaluated under exactly the same testing conditions. Specifically, for the

binary classification tasks, T1SEstacker [9] used Dataset S3, Bastion3 [10], T3SEpp [11] and EP3 [12] used Dataset S4, Bastion4 [13], CNN-T4SE [14], T4SEfinder [15], and T4SEpp [16] used Dataset S5, while Bastion6 [18] used Dataset S6. For the multi-class classification tasks, DeepSecE [27] and BastionX [43] were each evaluated separately on the above benchmark test sets.

UMAP-based embedding feature analysis

To visualize and compare the secreted effector embeddings learned by MoCETSE with the high-dimensional sequence features extracted by the pre-trained ESM-1b model, we employed the Uniform manifold approximation and projection (UMAP) algorithm via the umap-learn Python package for dimensionality reduction [44]. The local neighborhood size and the minimum distance between embedded points were set to 15 and 0.1, respectively. The UMAP experimental scripts are available at <https://github.com/YihangLin123/MoCETSE>.

Genome-wide effector prediction pipeline

Following the methodological framework of DeepSecE [27], we integrated Macsfinder [45], a dedicated tool for bacterial secretion system detection, with our MoCETSE model. Specifically, this pipeline first harnesses Macsfinder to search for structural component proteins of secretion systems based on hidden markov models (HMMs). Subsequently, it precisely identifies secretion system gene clusters by applying predefined spatial colocalization rules, including constraints on the inter-gene distance between adjacent components and the minimum component threshold. Based on these identified secretion system architectures, MoCETSE then scans all Coding Sequences (CDSs) within the bacterial genome, thereby enabling the accurate identification of potential secreted effector proteins.

Three representative Gram-negative bacterial strains, *Pseudomonas syringae* pv. tomato str. DC3000 (T3SE), *Legionella pneumophila* subsp. *pneumophila* str. Philadelphia 1 (T4SE), and *Pseudomonas aeruginosa* PAO1 (T6SE), were employed to evaluate the genome-wide predictive capability of MoCETSE for secretion system effectors. Recall was calculated as the proportion of experimentally validated effectors detected by MoCETSE:

$$\text{Recall} = \frac{N_{\text{overlap}}}{N_{\text{ver}}} \times 100\% \quad (13)$$

Where N_{ver} denotes the number of experimentally validated secreted proteins, and N_{overlap} represents the number of overlapping proteins between the experimentally identified and predicted secreted protein sets.

Sequence significance map

Leveraging the relative position multi-head attention mechanism embedded in the transformer module, the MoCETSE model enables quantitative assessment of the functional importance of amino acid residues within protein sequences. Specifically, by extracting attention weights from the multi-head attention matrix, the model transforms complex inter-residue interactions into quantifiable residue contribution scores. During this process, the model first aggregates the outputs from all attention heads and then calculates a global weighted average by incorporating relative positional weights, thereby characterizing the associative influence of specific residues on the entire sequence. The mathematical formulation is defined as follows:

$$\text{score} = \frac{1}{L} \sum_i \sum_h \sum_j (\text{attn}_{hij} \cdot R_{ij}) \quad (14)$$

Where R_{ij} denotes the relative positional weight between amino acids at positions i and j , and L represents the length of the protein sequence. The term attn_{hij} refers to the attention weight between positions i and j , as computed by the h -th attention head. A higher value of attn_{hij} indicates a stronger association between the two positions. To visualize the learned

positional importance, the Python toolkit Logomaker [46] is utilized to generate the sequence significance map. The map enables the intuitive characterization of key biological motifs within the secretory protein sequences.

Results

UMAP visualization of secretion protein embeddings

To evaluate the discriminative ability of the protein embedding features learned by the model, dimensionality reduction and visualization analysis were performed on the embedding vectors of secretory and non-secretory proteins in the training set (Dataset S1) [44]. UMAP maps the secretion embeddings generated by the MoCETSE model and the high-dimensional features generated by the pre-trained language model ESM-1b into a two-dimensional space, providing an intuitive visualization of the distribution relationships among the samples (Fig 2). The results showed that the embeddings of MoCETSE formed clear clusters in the low-dimensional space: the five different types of effector proteins and non-effector proteins each aggregated into distinct clusters, with clear boundaries between them (Fig 2B). Further observation revealed that samples within the same category were tightly clustered, indicating high consistency and homogeneity of sequence features. In contrast, the clusters of different categories displayed significant spatial separation, demonstrating that MoCETSE effectively distinguished the sequence features of varying protein types during the encoding process. This stands in sharp contrast to the ESM-1b projection, which showed considerable overlap between categories and poorly defined clustering boundaries (Fig 2A).

Furthermore, the secretion embeddings from DeepSecE (S3 Fig) exhibit noticeable overlaps between certain effector categories, particularly between T3SEs and T4SEs. Their samples show a significant intersection in the UMAP projection, leading to poorly defined cluster boundaries. In contrast, MoCETSE yields more distinct clustering and effectively

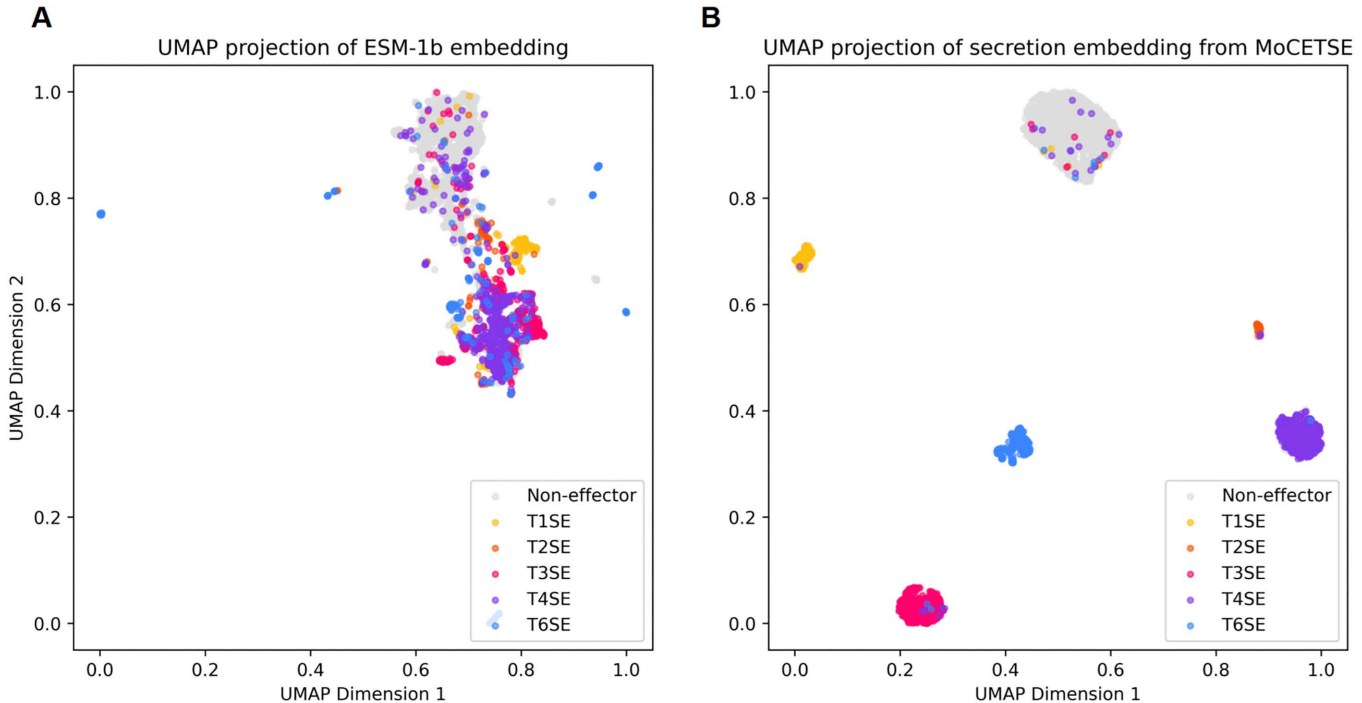


Fig 2. UMAP visualization of effector and non-effector protein clusters. (A) UMAP clustering projections of ESM-1b sequence embeddings. (B) UMAP clustering projections of secretion embeddings from MoCETSE. Different colors represent distinct effector types (T1SE–T4SE and T6SE) and non-effector proteins.

<https://doi.org/10.1371/journal.pcbi.1013397.g002>

minimizes such overlaps, demonstrating superior discriminative capacity across multiple effector classes. These results suggest that MoCETSE not only captures category-specific sequence features related to secretion systems but also substantially enhances the fine-grained classification of effector proteins.

Notably, as observed in [Figs 2](#) and [S3](#), the T2SE cluster (orange) partially overlaps with some T4SE samples in the visualization results of both MoCETSE and DeepSecE. This phenomenon can be attributed to the structural and evolutionary homology between the type II secretion system and the type IV pilus (T4P) system [\[3\]](#). Such intrinsic biological similarities may cause certain T2SEs to be situated in close proximity to T4SEs within the latent embedding space.

Performance evaluation of MoCETSE in effector protein prediction

We evaluated the performance of the MoCETSE model via five-fold cross-validation, where the validation set was partitioned from the training set (Dataset S1). To assess the model's generalization capability, we conducted additional testing on an independent test dataset (Dataset S2). In the model evaluation, we employed the one-vs-rest strategy and integrated the prediction outputs of the validation dataset during the cross-validation phase to separately plot the ROC curves for the five major types of secreted proteins and non-secreted proteins ([Fig 3A](#)). The AUC values for non-effector proteins and the five major secreted effector proteins ranged from 0.930 to 0.992, indicating the model's excellent classification ability. The multi-class confusion matrix (with percentages rounded to one decimal place) illustrates the model's prediction sensitivity and misclassification rates across all categories ([Fig 3B](#)). MoCETSE achieved a prediction sensitivity of approximately 95.0% for non-effector proteins, demonstrating its effectiveness in reducing the likelihood of false positives. For T1SEs, the sensitivity reached 90.6%, potentially due to their characteristically longer sequence lengths. Conversely, T2SEs exhibited a lower sensitivity of 76.5%, which can be attributed to their limited representation in the training set. Furthermore, we observed relatively higher misclassification rates between T3SE, T4SE, and T6SE; this may stem from their similar secretion mechanisms involving protein translocation across multiple bacterial membranes [\[47\]](#).

On the independent test dataset, MoCETSE exhibited robust generalization performance ([Fig 3C](#)), with AUC values for non-effector proteins and the five effector proteins ranging from 0.979 to 0.998. Except for non-secreted effector proteins, the prediction sensitivities for all other effector protein types exceeded 90% ([Fig 3D](#)). These results underscore the model's capability to reliably identify five types of secreted effector proteins while maintaining exceptional robustness and generalization.

Comparison of different foundation methods

In this study, we developed multiple models to conduct a comprehensive performance comparison. Initially, a PSSM-CNN prediction model based on PSSM features was developed using a three-layer convolutional neural network architecture. Subsequently, two pre-trained protein language models, TAPE [\[48\]](#) and ESM-1b [\[21\]](#), were employed as the base feature extraction modules, and a series of derivative models were further constructed based on ESM-1b. To investigate the impact of different training strategies on model performance, two distinct schemes were adopted: Linear Probing and Fine-tuning. Specifically, the Fine-tuning strategy involved unfreezing and optimizing only the final layer of the pre-trained model, whereas the Linear Probing strategy kept all parameters of the pre-trained model frozen and introduced only a linear classifier at the output layer. In addition, using the one-dimensional (1D) embedding features generated by ESM-1b, three classical machine learning algorithms, including the SVM, Random forest, and XGBoost, were employed to build classifiers. Moreover, the DeepSecE [\[27\]](#) model was constructed by integrating ESM-1b with a 1D convolutional layer and a secretion-specific transformer. We evaluated the performance of different models using three metrics: ACC, F1 score, and AUPRC. Experimental results demonstrated that MoCETSE, which incorporates a target preprocessing network and a secretion-specific transformer, significantly outperformed all other comparative models under both cross-validation (Dataset S1) and independent testing (Dataset S2) settings ([Table 1](#), [S4 Fig](#)).

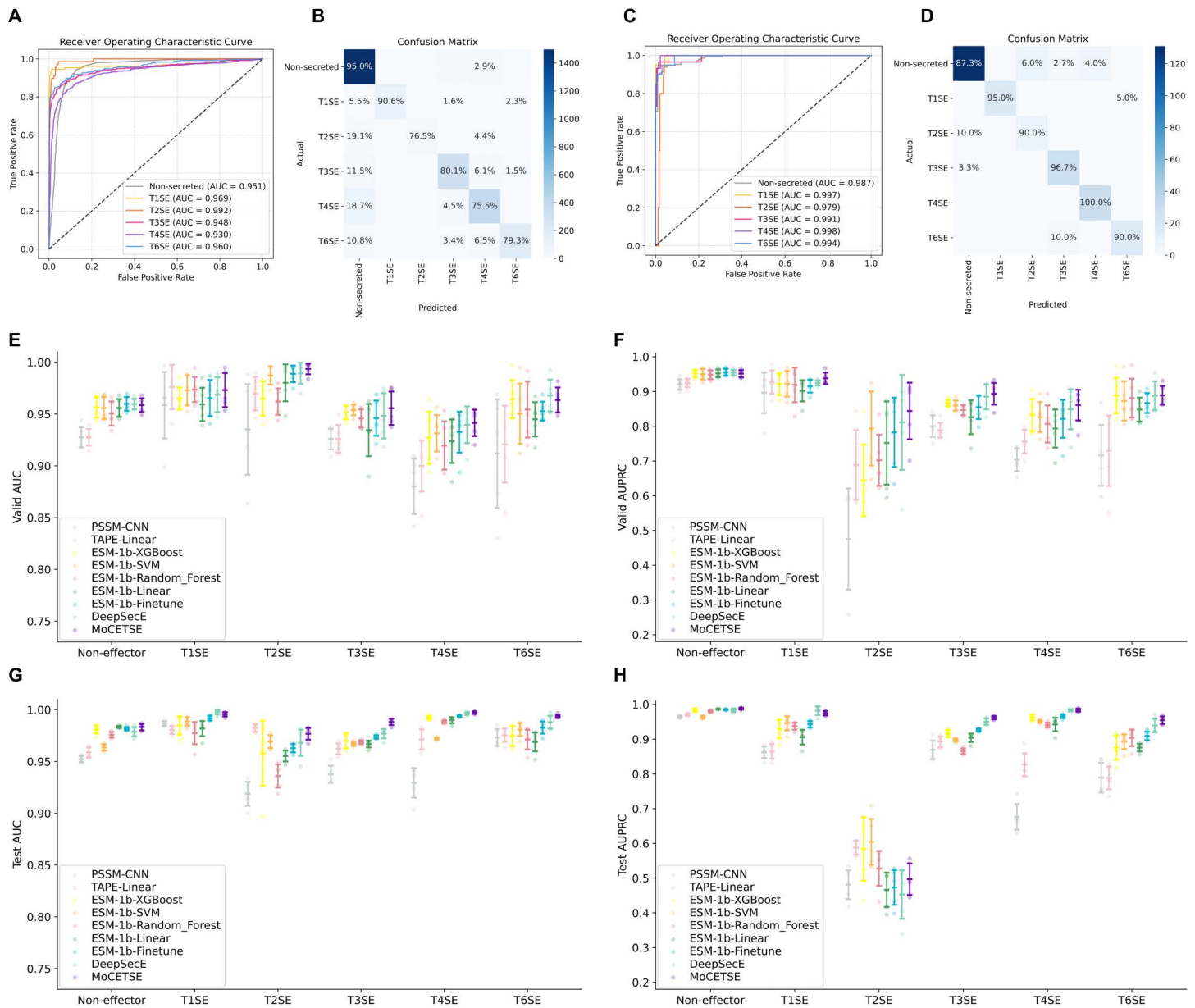


Fig 3. Performance evaluation of MoCETSE for secretion system effector protein prediction. (A to D) Model performance assessed via five-fold cross-validation (A,B) and independent testing (C,D). ROC curves illustrate the classification performance for each effector type. Sensitivity values for each class are indicated along the diagonals of the corresponding confusion matrices. (E–H) Comparison of AUC and AUPRC metrics across different model architectures or training strategies for the five protein categories under five-fold cross-validation (E,F) and independent testing (G,H). Error bars represent the 95% confidence intervals.

<https://doi.org/10.1371/journal.pcbi.1013397.g003>

In cross-validation testing, the PSSM-CNN model achieved an accuracy of 0.799 (95% CI: 0.772–0.826) and an F1 score of 0.712 (95% CI: 0.649–0.774). Incorporating pre-trained protein language models led to clear performance improvements. The TAPEBert-Linear model attained an accuracy of 0.816 (95% CI: 0.781–0.851) and an F1 score of 0.764 (95% CI: 0.686–0.842), while the larger-scale ESM-1b-Linear further improved accuracy to 0.847 (95% CI: 0.803–0.891) and F1 score to 0.771 (95% CI: 0.663–0.879). Notably, the ESM-1b Fine-tuning strategy did not result in a

Table 1. Model architecture performances across five-fold cross-validation (CV) and independent (Ind) tests.

Pre-trained model	Preprocessing module	Method	ACC		F1		AUPRC	
			CV	Ind	CV	Ind	CV	Ind
None	None	PSSM-CNN	0.799	0.822	0.712	0.724	0.752	0.774
TAPEBert	None	Linear probing	0.816	0.838	0.764	0.770	0.802	0.822
ESM-1b	None	SVM	0.845	0.888	0.751	0.776	0.871	0.876
ESM-1b	None	Random forest	0.835	0.872	0.735	0.762	0.851	0.859
ESM-1b	None	XGBoost	0.858	0.883	0.786	0.801	0.851	0.875
ESM-1b	None	Linear probing	0.847	0.874	0.771	0.784	0.846	0.845
ESM-1b	None	Fine-tuning	0.861	0.864	0.813	0.801	0.866	0.867
ESM-1b	1D convolutional layer	Secretion-specific transformer	0.873	0.883	0.843	0.836	0.885	0.882
ESM-1b	Target preprocessing network	Secretion-specific transformer	0.878	0.905	0.850	0.867	0.896	0.893

<https://doi.org/10.1371/journal.pcbi.1013397.t001>

significant performance improvement. Furthermore, the classifiers developed by integrating ESM-1b with three machine learning algorithms further underscore the exceptional feature extraction capabilities of ESM-1b. Specifically, the ESM-1b-XGBoost demonstrated superior performance, achieving an accuracy of 0.858 (95% CI: 0.830–0.886) and an F1 score of 0.786 (95% CI: 0.718–0.854). The DeepSecE model, leveraging a secretion-specific transformer, further enhanced the accuracy to 0.873 (95% CI: 0.835–0.911) and F1 score to 0.843 (95% CI: 0.785–0.901), substantiating the robust feature learning of specialized deep learning modules. Our proposed MoCETSE model outperformed all other models, with an accuracy of 0.878 (95% CI: 0.856–0.897) and an F1 score of 0.850 (95% CI: 0.811–0.889). When evaluated on the independent test set, MoCETSE improved the accuracy from 0.883 (95% CI: 0.836–0.931) to 0.905 (95% CI: 0.861–0.949) and the F1 score from 0.836 (95% CI: 0.782–0.889) to 0.867 (95% CI: 0.823–0.911) compared to DeepSecE, which utilizes a 1D convolutional layer. These results demonstrate that the TPN effectively refines high-dimensional PLM embeddings by filtering redundancy and capturing critical features for effector classification, while simultaneously highlighting the strong generalization capability of the MoCETSE framework.

To further evaluate the generalization capability of the model, 110 effector proteins in the independent test set were stratified into six groups based on their sequence identity relative to the training dataset. Local BLASTP was employed to determine the homology, and sequences were assigned to discrete identity intervals (<20%, 20–30%, 30–40%, 40–50%, 50–60%, and ≥60%) based on the highest identity score achieved under an E-value threshold of ≤0.01. As illustrated in [S5 Fig](#), MoCETSE achieved accuracy scores of 0.874, 0.892, 0.947, 0.970, 0.957, and 0.957 across the six sequence identity intervals, respectively. Notably, MoCETSE secured the top performance in five out of the six groups, with the sole exception being the 50%–60% interval, where its accuracy was slightly lower than that of DeepSecE (0.957 vs. 0.985). In both cross-validation and independent testing, this study compared the AUC and AUPRC metrics of different models for non-secretory proteins and five categories of secretory effector proteins ([Fig 3E–H](#)). The results indicated that using a pre-trained protein language model as a sequence feature extractor can improve the models' classification performance for various types of effector proteins. Overall, MoCETSE performed comparably to DeepSecE in secretory protein recognition tasks and exhibited a slight advantage over other derivative models based on the ESM-1b pre-trained language model.

Finally, we evaluated the effectiveness of the state-of-the-art pre-trained protein language model, ESM-2 [21], within our proposed framework. As illustrated in [S3 Table](#), replacing ESM-1b with ESM-2 did not yield the anticipated performance improvements. This observation is likely attributable to the relatively constrained size of our training dataset (comprising fewer than 3,000 sequences). Although ESM-2 features a more extensive parameter scale and enhanced representative power, its inherent complexity typically necessitates a larger and more diverse data volume to fully manifest its advantages. On a limited dataset, such a high-capacity model is more susceptible to overfitting. In contrast, ESM-1b provides more robust and generalizable feature representations that better align with our current data scale.

Performance comparison against existing popular models

The MoCETSE model aims to predict the secretion system effector proteins (T1SEs–T4SEs and T6SEs) in Gram-negative bacteria. To objectively assess its predictive performance, the present study compares MoCETSE with two mainstream multi-class models, DeepSecE [27] and BastionX [43], as well as nine representative binary classification models, using benchmark datasets (Dataset S3–S6). Specifically, T1SEstacker [9] was used for T1SE prediction; Bastion3 [10], T3SEpp [11], and EP3 [12] for T3SE; Bastion4 [13], CNN-T4SE [14], T4SEfinder [15], and T4SEpp [16] for T4SE; and Bastion6 [18] for T6SE. Evaluation metrics included ACC, REC, PR, F1 score, and MCC.

The experimental results demonstrate that MoCETSE outperforms existing multi-class models in overall performance. For T3SE and T4SE prediction, its performance is comparable to the state-of-the-art binary classification models, while for T1SE and T6SE prediction, MoCETSE exhibits superior performance (S4 Table). In the T1SE classification task (Dataset S3), MoCETSE outperforms T1SEstacker, which relies on amino acid composition without RTX C-terminal motifs (accuracy 98.2% vs 92.9%, F1 score 0.947 vs 0.727, MCC 0.942 vs 0.691) (S6 FigA). For the T3SE classification task (Dataset S4), MoCETSE achieved an accuracy of 91.7%, F1 score of 0.918, and MCC of 0.835. Its performance is comparable to T3SEpp (a method integrating various biological features) (Fig 4A). In the T4SE classification task (Dataset S5), MoCETSE outperforms DeepSecE, which also uses a pre-trained protein language model (accuracy 98.3% vs 97.8%, F1 score 0.949 vs 0.935, MCC 0.939 vs 0.911), and is second only to CNN-T4SE (Fig 4B). For the T6SE classification task (Dataset S6), MoCETSE demonstrated superior performance compared to Bastion6, which relies on PSSM predictors (ACC: 94.3%, F1: 0.946, MCC: 0.892), achieving an accuracy of 98.6%, F1 score of 0.957, and MCC of 0.920 (Fig 4C).

Although MoCETSE has not yet surpassed the popular T3SE prediction model Bastion3 and the T4SE prediction model CNN-T4SE in terms of benchmark performance, it demonstrates a significant advantage in minimizing cross-type misclassification on the independent test set (Dataset S2). Binary classifiers often mistakenly predict other types of secretion effector proteins as the target type. In contrast, the MoCETSE model effectively controls the false positive rate when confronted with multiple effector proteins, reducing the likelihood of misclassifying other types of secretion effector proteins (Figs 4D–F, S6C). For instance, although Bastion3 accurately identifies all T3SEs in the independent test dataset, it misclassifies 9 T1SEs and 12 T4SEs as T3SEs (Fig 4D). Similarly, CNN-T4SE misclassifies 10 T1SEs and 4 T3SEs as T4SEs (Fig 4E). Furthermore, we evaluated computational efficiency by comparing the runtime of MoCETSE with other classification methods on Dataset S2. The results indicate that MoCETSE completes predictions in less than one minute, demonstrating superior computational efficiency compared to existing models (Fig 4G).

Ablation experiments

To evaluate the contribution of each core component to the overall performance of the proposed model, we designed and conducted a series of ablation experiments. Using a framework consisting of a 1D convolutional layer (1D Conv) and multi-head attention (MHA) as the baseline architecture, we systematically integrated TPN and relative position encoding multi-head attention (RPE-MHA). The functional efficacy of each module was quantified by evaluating four model configurations across both five-fold cross-validation (CV) and an independent test set (Ind) (Table 2).

The experimental results demonstrate that MoCETSE, the final model integrating both TPN and RPE-MHA, significantly outperforms the baseline (1D Conv + MHA) across all evaluation metrics. In the five-fold cross-validation, MoCETSE achieved the superior performance with an ACC of 0.878 (95% CI: 0.856–0.897) and an F1 score of 0.850 (95% CI: 0.811–0.889). Notably, on the independent test set, the generalization performance of MoCETSE improved further, reaching an ACC of 0.905 (95% CI: 0.861–0.949) and an F1 score of 0.867 (95% CI: 0.823–0.911), representing a substantial gain over the baseline (ACC: 0.883, F1: 0.836).

Component decomposition analysis revealed that replacing the conventional 1D Conv with TPN, while the MHA remained unchanged, resulted in improvements across all performance metrics. Specifically, the ACC in CV increased

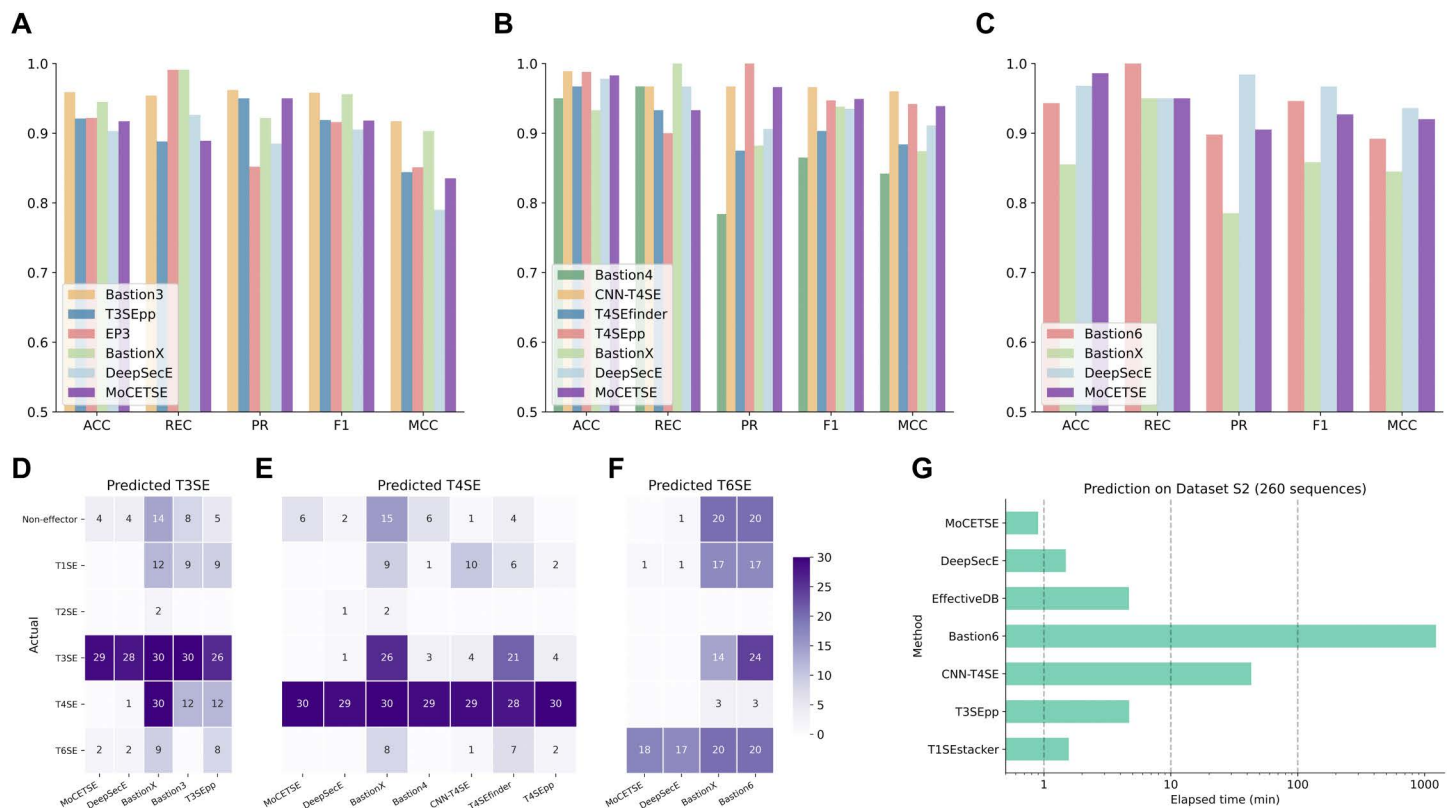


Fig 4. Benchmarking results of MoCETSE against mainstream secretion effector protein prediction models. (A–C) Performance comparison of MoCETSE with mainstream prediction models on T3SE (A), T4SE (B), and T6SE (C) prediction tasks. (D–F) Comparison of true positive and false positive predictions from different models on the independent test set (Dataset S2), which contains 150 non-effectors, 20 T1SEs, 10 T2SEs, 30 T3SEs, 30 T4SEs, and 20 T6SEs. The heatmaps illustrate the distribution of predicted samples across various effector protein types. (G) Evaluation of the prediction efficiency of each model on the independent test dataset (Dataset S2) containing 260 effector proteins.

<https://doi.org/10.1371/journal.pcbi.1013397.g004>

Table 2. Performance comparison of different model variants in ablation experiments under five-fold cross-validation (CV) and independent (Ind) test settings.

Preprocessing module	Attention mechanism	ACC		F1		AUPRC	
		CV	Ind	CV	Ind	CV	Ind
1D convolutional layer	Multi-head attention	0.873	0.883	0.843	0.836	0.885	0.882
1D convolutional layer	RPE-multi-head attention	0.873	0.872	0.837	0.831	0.877	0.884
Target preprocessing network	Multi-head attention	0.875	0.893	0.847	0.848	0.886	0.889
Target preprocessing network	RPE-multi-head attention	0.878	0.905	0.850	0.867	0.896	0.893

<https://doi.org/10.1371/journal.pcbi.1013397.t002>

from 0.873 to 0.875, and the F1 score improved from 0.843 to 0.847. On the Ind, the ACC rose from 0.883 to 0.893, and the F1 score increased from 0.836 to 0.848. These results indicate that TPN captures the underlying features of target sequences more precisely than 1D Conv, providing a more informative representational basis for subsequent feature fusion and interaction, and thus improving the overall decision-making ability of the model. Furthermore, a significant synergistic effect was observed between RPE-MHA and the different preprocessing modules. Under the 1D Conv-based configuration, the introduction of the RPE mechanism did not yield performance gains. However, when combined with TPN, model performance improved markedly; specifically, the ACC on the independent test set increased from 0.872 to

0.905, and the F1 score rose from 0.831 to 0.867. This comparison indicates that RPE effectively compensates for the lack of relative spatial information in the high-dimensional features extracted by TPN, thereby enabling the model to more accurately capture long-range dependencies within the sequences.

To further verify the robustness of each MoCETSE module across different effector protein classes, especially its generalization ability on sample-scarce categories, we performed a stratified ablation experiment on T1SE, T2SE, T3SE, T4SE, and T6SE using the independent test set (Dataset S2). Our results showed that the component improvements in MoCETSE led to consistent performance gains, both in the T1SE and T2SE categories with relatively small sample sizes and in other categories with more abundant samples (Figs 5A–C, S7). In particular, for T2SE—the category with the highest prediction difficulty—the baseline model (1D-Conv + MHA) achieved an F1 score of only 0.550, suggesting notable limitations in its ability to capture relevant sequence features. After integrating TPN and RPE-MHA modules, MoCETSE yielded performance improvement on this class, with ACC and F1 score reaching 0.961 and 0.636, respectively (S7 FigB). Taken together, these results demonstrate the effectiveness of the proposed improvements in capturing and representing the characteristic features of the target sequences.

Genome-wide prediction of secreted proteins

Genome-wide prediction of secreted proteins is a pivotal strategy for elucidating the pathogenicity, evolutionary dynamics, and distribution patterns within bacterial populations [49]. To evaluate the practical utility and generalization performance of MoCETSE in real-world scenarios, we integrated our model with Macsfinder [45], a specialized tool for structural identification of secretion systems. This integration established a high-throughput pipeline for the systematic recognition of effector proteins across the genomes of Gram-negative bacteria (detailed in the “Materials and methods” section).

To evaluate the genome-wide prediction capability of MoCETSE for secreted proteins, three representative Gram-negative bacterial strains were selected: *Pseudomonas syringae* pv. tomato str. DC3000 (T3SE) [National Center for Biotechnology Information (NCBI) accession number: NC_004578.1], *Legionella pneumophila* subsp. *pneumophila* str. Philadelphia 1 (T4SE) (NC_002942.5), and *Pseudomonas aeruginosa* PAO1 (T6SE) (NC_002516.2) (Fig 5D–F, S5 Table). For *P. syringae* pv. tomato str. DC3000 (T3SE), both MoCETSE and DeepSecE successfully detected 30 experimentally validated effectors, achieving a recall rate of 90.9%.

Notably, MoCETSE identified significantly fewer candidates (273 vs 330), leading to a higher precision (11.0% vs 9.1%) and F1 score (0.196 vs 0.165). These results, combined with a superior AUPRC (0.505 vs 0.412), demonstrates that MoCETSE can enhance prediction specificity without compromising sensitivity. For *L. pneumophila* subsp. *pneumophila* str. Philadelphia 1 (T4SE) and *P. aeruginosa* PAO1 (T6SE), MoCETSE achieved higher recall rates than DeepSecE (96.7% vs 86.3% and 90.0% vs 80.0%, respectively). Although MoCETSE produced more candidates in these two cases and yielded slightly lower precision, its F1 scores remained competitive (0.829 vs 0.827 for T4SE), indicating a robust balance between accuracy and discovery potential. Furthermore, MoCETSE exhibited a higher AUPRC value in the T4SE case, further confirming its superior discriminative power across different thresholds. In practical genomic screening, where minimizing false negatives is often prioritized to avoid missing functional effectors, MoCETSE offers a more comprehensive and reliable solution for identifying experimentally validated secreted proteins in Gram-negative bacteria.

Relative position attention identification sequence motifs

Signal peptides located at the N- or C-terminus of protein sequences are pivotal and widely recognized hallmarks of many secretory proteins, particularly those exported via classical secretion pathways [3]. To evaluate the universality and precision of MoCETSE in identifying functional secretion signals across diverse substrates, we selected representative effector proteins from T1SS, T2SS, T3SS, T4SS, and T6SS. We characterized their secretion motifs through sequence attention weight analysis, which allows for the direct visualization of the model’s focus during feature extraction (detailed in the “Materials and methods” section).

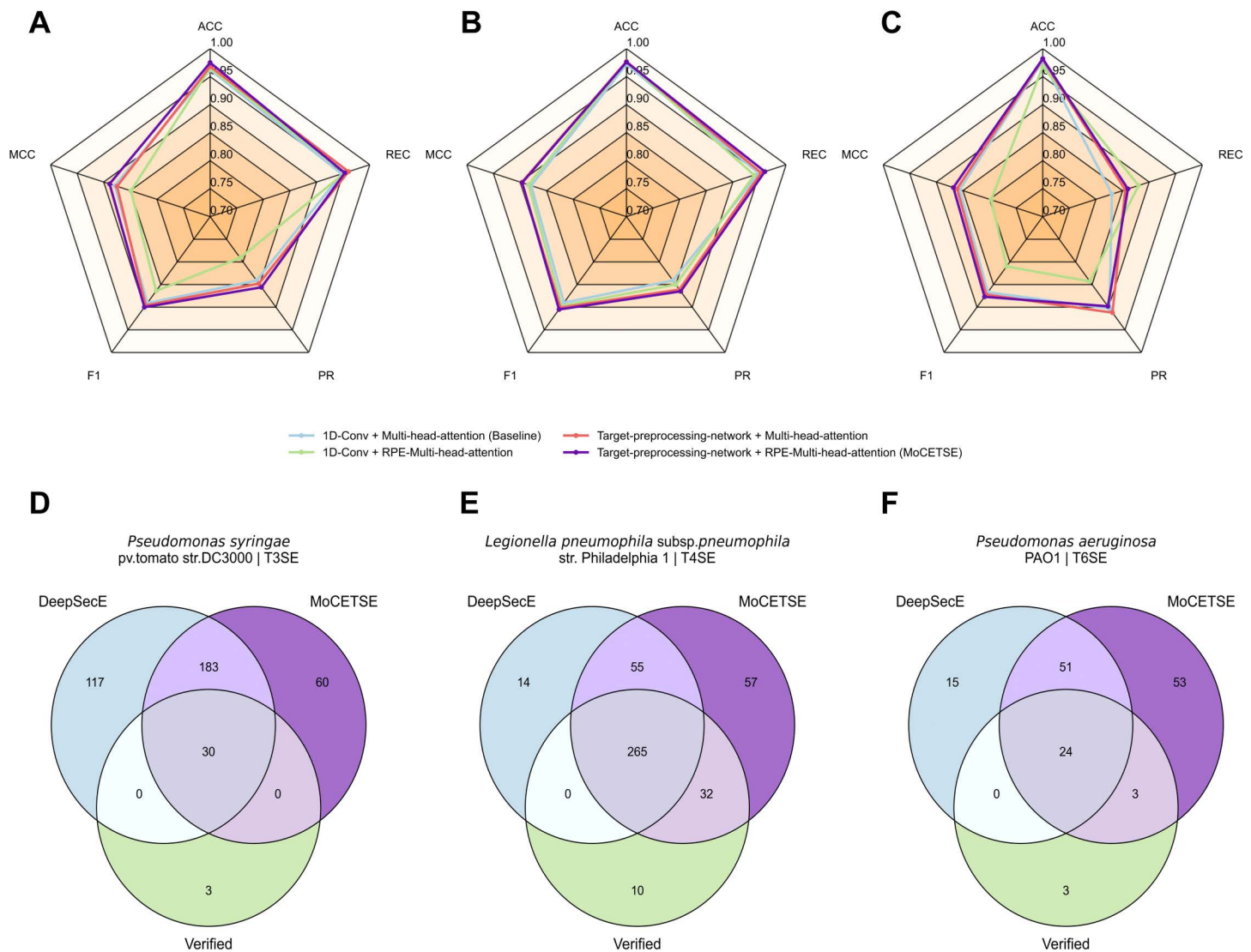


Fig 5. Stratified ablation analysis and genome-wide effector protein prediction. (A–C) Performance comparison of four model configurations across different effector protein classes. The ablation results for (A) T3SEs, (B) T4SEs, and (C) T6SEs on the independent test set are shown. (D–F) Comparison of genome-wide effector protein predictions between MoCETSE and DeepSecE across three representative Gram-negative bacterial strains: (D) *Pseudomonas syringae* pv. tomato str. DC3000 (T3SE); (E) *Legionella pneumophila* subsp. *pneumophila* str. Philadelphia 1 (T4SE); and (F) *Pseudomonas aeruginosa* PAO1 (T6SE). The Venn diagrams illustrate the overlap between the model predictions and experimentally verified effector proteins (Verified).

<https://doi.org/10.1371/journal.pcbi.1013397.g005>

In the analysis of the *Escherichia coli* T1SS effector HlyA (UniProt accession P08715), MoCETSE not only precisely localized the critical non-cleavable C-terminal secretion signal within the last 20 residues (S8 FigB) but also accurately identified multiple RTX (Repeat in ToXins) motifs characterized by conserved aspartic acid (D) residues distributed between positions 700 and 800 (S8 FigA) [50]. For the *Vibrio cholerae* serotype O1 T2SS effector CtxA (UniProt accession P11439), the model captured the Sec-dependent signal peptide at the N-terminus (residues 1–18) (S8 FigC); this finding is highly consistent with the established biological feature that CtxA, as a T2SS substrate, must first be translocated to the periplasm via the Sec pathway [51]. Similarly, MoCETSE precisely identified key functional regions of the *Salmonella typhimurium*

T3SS secretory protein SptP (UniProt accession P74873), including its N-terminal secretion signal (residues 1–15) and the chaperone-binding domain (CBD, residues 15–35) associated with specific chaperone binding (Fig 6A) [52]. In the characterization of the *Legionella pneumophila* subsp. *pneumophila* str. Philadelphia 1 T4SS effector VpdB (UniProt accession Q9KNE5), the model localized the N-terminal auxiliary initiation sequence (Fig 6B) and further pinpointed the core C-terminal secretion signal at residues 580–600 (Fig 6C). This result provides robust support for the biological mechanism reported in the literature regarding the recognition of C-tail features by the Dot/Icm system for translocation [53]. Regarding the *Vibrio cholerae* serotype O1 T6SS effector VasX (UniProt accession Q9KNE5), MoCETSE successfully captured sparse conservation patterns within the first 100 N-terminal residues and precisely focused on key sites highly consistent with the MIX (Marker for Type VI effectors) motif, assigning particularly high attention weights to residues D78, Y83, and K89 (S8 FigD). This identification aligns perfectly with the known mechanism of T6SS-mediated protein translocation via specific MIX signals [54]. Collectively, these experimental results demonstrate that MoCETSE can transcend the limitations of sequence position, accurately identifying characteristic motifs of various secretion systems across the entire sequence scale.

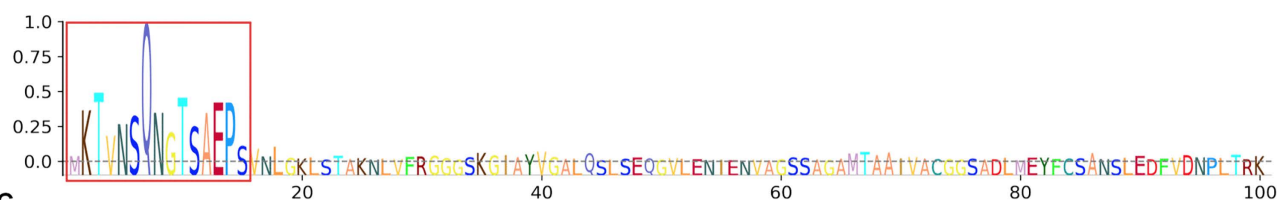
Discussion

The advancement of machine learning technologies has provided crucial momentum for the identification of virulence factors in bacterial pathogens, particularly in the study of secreted effector proteins [49]. In particular, the rapid development of PLMs has endowed effector prediction tools with enhanced feature representation capabilities, significantly deepening

A *Salmonella typhimurium* T3SS effector SptP



B *Legionella pneumophila* subsp. *pneumophila* str. Philadelphia 1 T4SS effector VpdB



C



Fig 6. Visualization of sequence motifs and attention weight distributions for T3SS and T4SS effector proteins. (A) Attention weight distribution for the *Salmonella typhimurium* T3SS effector protein SptP (P74873). The sequence map highlights the high attention weights assigned to the N-terminal secretion signal and the chaperone-binding domain (CBD). (B–C) Attention weight distribution for the *Legionella pneumophila* subsp. *pneumophila* str. Philadelphia 1 T4SS effector protein VpdB (Q5ZW60). (B) Highlights the N-terminal initiating auxiliary sequence recognized by the model. (C) Displays the core secretion signal concentrated in the C-terminal region.

<https://doi.org/10.1371/journal.pcbi.1013397.g006>

the models' understanding of complex sequence semantics and thereby improving prediction accuracy [21,27]. To further optimize the performance of effector protein prediction models, we developed MoCETSE, an end-to-end framework that refines the raw, general protein representations generated by PLMs to focus on key features associated with secretion functions and directly performs classification, thus integrating feature processing and prediction into a unified workflow. Benchmarking results demonstrated that MoCETSE achieved the best performance in T1SE and T6SE classification tasks and remained highly competitive in T3SE and T4SE prediction (Fig 4A–C). Moreover, MoCETSE exhibited lower false-positive rates and higher specificity in mixed-effector scenarios (Fig 4D–F), a characteristic that could substantially reduce the burden of subsequent experimental validation in practical applications.

Further analyses revealed that TPN and the secretion-specific transformer module in MoCETSE synergistically capture and enhance secretion-relevant features, thereby significantly boosting overall model performance (Fig 3A–D, Table 1). ROC curve and confusion matrix analyses indicated that MoCETSE achieved relatively balanced predictions across all effector classes, demonstrating strong multi-class classification capability with minimal bias (Fig 3E–H). Stratified analysis based on sequence similarity in the independent test set further showed that MoCETSE maintained stable and high performance across different sequence identity ranges (S5 Fig). Collectively, these findings suggest that MoCETSE is robust to sequence heterogeneity and exhibits high reliability when confronted with uncharacterized or distantly related effector proteins.

Leveraging the RPE-MHA mechanism, MoCETSE accurately captures core signal features of secreted proteins, including key signal peptide regions located at the N-terminus, C-terminus, or within the internal sequence (Figs 6, S8). This not only improves the model's predictive performance but also provides robust biological interpretability by aligning with established biological mechanisms. Such alignment reinforces the credibility of deep learning-based approaches in effector protein prediction.

Genome-scale prediction results further validated the robustness and reliability of MoCETSE in handling real-world biological sequence data. In analyses of representative pathogenic species such as *P. syringae*, *L. pneumophila*, and *P. aeruginosa*, MoCETSE not only outperformed existing tools in recall but also effectively reduced redundancy among candidate proteins while maintaining high sensitivity (Fig 5D–F, S5 Table). This high specificity is particularly valuable for large-scale omics screening, as it can significantly narrow down the range of candidates requiring experimental validation, thereby reducing research costs. The ability of MoCETSE to accurately identify effectors at the genome-wide scale provides a powerful tool for investigating bacterial pathogenic mechanisms and host–pathogen interactions.

Despite its strong performance in predicting bacterial secretion system effectors, MoCETSE still has room for improvement. First, constrained by the size of experimentally validated effector datasets, the model may not fully capture the complete feature space of some rare substrates. Second, the uneven distribution of data across different secretion system classes (e.g., the relative scarcity of T2SE samples) may affect the model's generalization ability for minority classes. Future work will focus on integrating data from a broader range of species and newly discovered effector substrates to further enhance model robustness. Additionally, adopting larger or more advanced PLMs (such as ESM2 or SaProt) [21,55] may enable the extraction of more biologically meaningful features, thereby further improving prediction accuracy. However, balancing computational efficiency and model interpretability while pursuing performance gains remains an important consideration for future research.

Conclusion

MoCETSE achieves significant improvements in prediction accuracy and generalization by integrating evolutionary and structural information encoded by the pre-trained protein language model ESM-1b, secretion-specific core features refined by the target preprocessing network, and long-range sequence dependencies captured by the relative position-enhanced transformer module. Beyond providing an efficient tool for the rapid identification of effector proteins, this model enhances the understanding of bacterial pathogenic mechanisms through the precise localization of functional motifs. MoCETSE offers a valid paradigm for the integration of large-scale protein language models with task-specific biological modules,

thereby establishing a more robust research framework for pathogen-associated proteins. With the continuous advancement of large language models, future studies will focus on adopting more advanced and larger-scale protein language models to further strengthen the robustness and practicality of the model, enabling it to play a pivotal role in addressing increasingly complex biological tasks.

Supporting information

S1 Table. Summary of the training, testing, and benchmarking datasets.

(PDF)

S2 Table. Web server links of the comparative models used in this study.

(PDF)

S3 Table. Performance comparison of MoCETSE based on ESM-1b and ESM-2 across five-fold cross-validation (CV) and independent (Ind) tests.

(PDF)

S4 Table. Performance comparison between MoCETSE and existing popular methods on benchmark datasets for T1SE, T2SE, T3SE, T4SE, and T6SE prediction.

(PDF)

S5 Table. Comparison of genome-wide effector protein prediction results for three representative Gram-negative bacterial strains across different models.

(PDF)

S1 Fig. Distribution of protein sequence lengths in the training and test datasets. (A–B) Histograms illustrating the sequence length composition of the training dataset (A) and the test dataset (B). The x-axis represents the protein sequence length (number of amino acids), while the y-axis represents the frequency of proteins within each interval.

(PDF)

S2 Fig. Composition distribution of the secreted substrate dataset. (A–F) Proportional distribution of each protein category across the training, independent test, and benchmark datasets employed in this study.

(PDF)

S3 Fig. UMAP visualization of effector and non-effector protein embeddings. (A) UMAP projections of secretion embeddings generated by DeepSecE. (B) UMAP projections of secretion embeddings generated by MoCETSE.

(PDF)

S4 Fig. Performance comparison of various models across cross-validation and independent testing. (A,B) Average accuracy, F1 score, and AUPRC of different models evaluated by five-fold cross-validation (A) and independent testing (B). Data are presented as mean values with 95% confidence intervals. MoCETSE consistently achieved the highest overall predictive performance across all metrics.

(PDF)

S5 Fig. Evaluation of model generalization across sequence identity gradients. Generalization performance of effector prediction models assessed by (A) Accuracy and (B) F1 score across different sequence identity intervals. The independent test set (Dataset S2, $n = 110$) was categorized into six gradients based on sequence identity relative to the training set (Dataset S1): $< 20\%$, $20\text{--}30\%$, $30\text{--}40\%$, $40\text{--}50\%$, $50\text{--}60\%$, and $\geq 60\%$. Shaded areas represent the 95% confidence intervals (CI).

(PDF)

S6 Fig. Performance comparison of MoCETSE and mainstream models in type I (T1SE) and type II (T2SE) secretion effector protein prediction tasks. (A–B) Model performance evaluated on the benchmark test set (Dataset S3) and a dataset containing 150 non-effector proteins and 10 T6SEs, using metrics including accuracy (ACC), recall (REC), precision (PR), F1 score (F1), and Matthews correlation coefficient (MCC). (C–D) Comparison of true positive and false positive predictions for T1SE and T2SE across different models. Heatmaps show the distribution of samples predicted as T1SE, T2SE, and other effector protein types.
(PDF)

S7 Fig. Stratified ablation analysis for sample-scarce effector classes. (A) and (B) show the ablation performance comparison of different model configurations for the T1SE and T2SE classes, respectively.
(PDF)

S8 Fig. Interpretability analysis of MoCETSE attention weights across T1SS, T2SS, and T6SS effectors. (A–B) Motif identification of *Escherichia coli* T1SS effector HlyA (P08715). (A) shows the RTX (Repeat in ToXins) motifs distributed between residues 700–800. (B) shows the C-terminal secretion signal located within the last 20 residues of the sequence. (C) Identification of the Sec-dependent signal peptide in *Vibrio cholerae* serotype O1 T2SS effector CtxA (P11439). High attention weights are concentrated in the N-terminal region (residues 1–18), which is consistent with the translocation characteristics of the Sec pathway. (D) Characterization of the MIX motif in *Vibrio cholerae* T6SS effector VasX (Q9KNE5). The model precisely focuses on key sites (D78, Y83 and K89) within the first 100 N-terminal residues.
(PDF)

Acknowledgments

We thank Jiani Chen and Zhouying Li for discussion on related topics.

Author contributions

Conceptualization: Hua Shi, Quan Zou.

Data curation: Yihang Lin, Dachen Liu.

Formal analysis: Yihang Lin, Dachen Liu.

Investigation: Yihang Lin.

Methodology: Hua Shi, Yihang Lin.

Project administration: Hua Shi, Quan Zou.

Software: Yihang Lin.

Supervision: Hua Shi, Quan Zou.

Validation: Dachen Liu.

Writing – original draft: Hua Shi, Yihang Lin.

Writing – review & editing: Hua Shi, Yihang Lin, Dachen Liu, Quan Zou.

References

1. Costa TRD, Felisberto-Rodrigues C, Meir A, Prevost MS, Redzej A, Trokter M, et al. Secretion systems in Gram-negative bacteria: structural and mechanistic insights. *Nat Rev Microbiol*. 2015;13(6):343–59. <https://doi.org/10.1038/nrmicro3456> PMID: 25978706
2. Green ER, Meccas J. Bacterial Secretion Systems: An Overview. *Microbiol Spectr*. 2016;4(1):10.1128/microbiolspec.VMBF-0012–2015. <https://doi.org/10.1128/microbiolspec.VMBF-0012-2015> PMID: 26999395
3. Hui X, Chen Z, Zhang J, Lu M, Cai X, Deng Y, et al. Computational prediction of secreted proteins in gram-negative bacteria. *Comput Struct Biotechnol J*. 2021;19:1806–28. <https://doi.org/10.1016/j.csbj.2021.03.019> PMID: 33897982

4. Sanchez-Garrido J, Ruano-Gallego D, Choudhary JS, Frankel G. The type III secretion system effector network hypothesis. *Trends Microbiol.* 2022;30(6):524–33. <https://doi.org/10.1016/j.tim.2021.10.007> PMID: [34840074](https://pubmed.ncbi.nlm.nih.gov/34840074/)
5. Singh RP, Kumari K. Bacterial type VI secretion system (T6SS): an evolved molecular weapon with diverse functionality. *Biotechnol Lett.* 2023;45(3):309–31. <https://doi.org/10.1007/s10529-023-03354-2> PMID: [36683130](https://pubmed.ncbi.nlm.nih.gov/36683130/)
6. Ruano-Gallego D, Sanchez-Garrido J, Kozik Z, Núñez-Berruoco E, Cepeda-Molero M, Mullineaux-Sanders C, et al. Type III secretion system effectors form robust and flexible intracellular virulence networks. *Science.* 2021;371(6534):eabc9531. <https://doi.org/10.1126/science.abc9531> PMID: [33707240](https://pubmed.ncbi.nlm.nih.gov/33707240/)
7. Böck D, Hüsler D, Steiner B, Medeiros JM, Welin A, Radomska KA, et al. The Polar Legionella Icm/Dot T4SS Establishes Distinct Contact Sites with the Pathogen Vacuole Membrane. *mBio.* 2021;12(5):e0218021. <https://doi.org/10.1128/mBio.02180-21> PMID: [34634944](https://pubmed.ncbi.nlm.nih.gov/34634944/)
8. Colautti J, Kelly SD, Whitney JC. Specialized killing across the domains of life by the type VI secretion systems of *Pseudomonas aeruginosa*. *Biochem J.* 2025;482(1):1–15. <https://doi.org/10.1042/BCJ20230240> PMID: [39774785](https://pubmed.ncbi.nlm.nih.gov/39774785/)
9. Chen Z, Zhao Z, Hui X, Zhang J, Hu Y, Chen R, et al. T1SEstacker: A Tri-Layer Stacking Model Effectively Predicts Bacterial Type 1 Secreted Proteins Based on C-Terminal Non-repeats-in-Toxin-Motif Sequence Features. *Front Microbiol.* 2022;12:813094. <https://doi.org/10.3389/fmicb.2021.813094> PMID: [35211101](https://pubmed.ncbi.nlm.nih.gov/35211101/)
10. Wang J, Li J, Yang B, Xie R, Marquez-Lago TT, Leier A, et al. Bastion3: a two-layer ensemble predictor of type III secreted effectors. *Bioinformatics.* 2019;35(12):2017–28. <https://doi.org/10.1093/bioinformatics/bty914> PMID: [30388198](https://pubmed.ncbi.nlm.nih.gov/30388198/)
11. Hui X, Chen Z, Lin M, Zhang J, Hu Y, Zeng Y, et al. T3SEpp: an Integrated Prediction Pipeline for Bacterial Type III Secreted Effectors. *mSystems.* 2020;5(4):e00288–20. <https://doi.org/10.1128/mSystems.00288-20> PMID: [32753503](https://pubmed.ncbi.nlm.nih.gov/32753503/)
12. Li J, Wei L, Guo F, Zou Q. EP3: an ensemble predictor that accurately identifies type III secreted effectors. *Brief Bioinform.* 2021;22(2):1918–28. <https://doi.org/10.1093/bib/bbaa008> PMID: [32043137](https://pubmed.ncbi.nlm.nih.gov/32043137/)
13. Wang J, Yang B, An Y, Marquez-Lago T, Leier A, Wilksch J, et al. Systematic analysis and prediction of type IV secreted effector proteins by machine learning approaches. *Brief Bioinform.* 2019;20(3):931–51. <https://doi.org/10.1093/bib/bbx164> PMID: [29186295](https://pubmed.ncbi.nlm.nih.gov/29186295/)
14. Hong J, Luo Y, Mou M, Fu J, Zhang Y, Xue W, et al. Convolutional neural network-based annotation of bacterial type IV secretion system effectors with enhanced accuracy and reduced false discovery. *Brief Bioinform.* 2020;21(5):1825–36. <https://doi.org/10.1093/bib/bbz120> PMID: [31860715](https://pubmed.ncbi.nlm.nih.gov/31860715/)
15. Zhang Y, Zhang Y, Xiong Y, Wang H, Deng Z, Song J, et al. T4SEfinder: a bioinformatics tool for genome-scale prediction of bacterial type IV secreted effectors using pre-trained protein language model. *Brief Bioinform.* 2022;23(1):bbab420. <https://doi.org/10.1093/bib/bbab420> PMID: [34657153](https://pubmed.ncbi.nlm.nih.gov/34657153/)
16. Hu Y, Wang Y, Hu X, Chao H, Li S, Ni Q, et al. T4SEpp: A pipeline integrating protein language models to predict bacterial type IV secreted effectors. *Comput Struct Biotechnol J.* 2024;23:801–12. <https://doi.org/10.1016/j.csbj.2024.01.015> PMID: [38328004](https://pubmed.ncbi.nlm.nih.gov/38328004/)
17. Li J, He S, Zhang J, Zhang F, Zou Q, Ni F. T4Seeker: a hybrid model for type IV secretion effectors identification. *BMC Biol.* 2024;22(1):259. <https://doi.org/10.1186/s12915-024-02064-z> PMID: [39543674](https://pubmed.ncbi.nlm.nih.gov/39543674/)
18. Wang J, Yang B, Leier A, Marquez-Lago TT, Hayashida M, Rocker A, et al. Bastion6: a bioinformatics approach for accurate prediction of type VI secreted effectors. *Bioinformatics.* 2018;34(15):2546–55. <https://doi.org/10.1093/bioinformatics/bty155> PMID: [29547915](https://pubmed.ncbi.nlm.nih.gov/29547915/)
19. Sen R, Nayak L, De RK. PyPredT6: A python-based prediction tool for identification of Type VI effector proteins. *J Bioinform Comput Biol.* 2019;17(3):1950019. <https://doi.org/10.1142/S0219720019500197> PMID: [31288641](https://pubmed.ncbi.nlm.nih.gov/31288641/)
20. Zeng C, Zou L. An account of in silico identification tools of secreted effector proteins in bacteria and future challenges. *Brief Bioinform.* 2019;20(1):110–29. <https://doi.org/10.1093/bib/bbx078> PMID: [28981574](https://pubmed.ncbi.nlm.nih.gov/28981574/)
21. Lin Z, Akin H, Rao R, Hie B, Zhu Z, Lu W, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science.* 2023;379(6637):1123–30. <https://doi.org/10.1126/science.ade2574> PMID: [36927031](https://pubmed.ncbi.nlm.nih.gov/36927031/)
22. Rives A, Meier J, Sercu T, Goyal S, Lin Z, Liu J, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proc Natl Acad Sci U S A.* 2021;118(15):e2016239118. <https://doi.org/10.1073/pnas.2016239118> PMID: [33876751](https://pubmed.ncbi.nlm.nih.gov/33876751/)
23. Rao R, Meier J, Sercu T, Ovchinnikov S, Rives A, Jones P, et al. Transformer protein language models are unsupervised structure learners. In: *Proceedings of the 9th International Conference on Learning Representations; 2021 May 3–7; Vienna, Austria.* Vienna (Austria): ICLR Press; 2021. p. 123–30.
24. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez A, et al. Attention is all you need. In: *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS 2017); 2017 Dec 4–9; Long Beach, CA.* New York: Curran Associates; 2017. p. 5998–6008.
25. Brandes N, Ofer D, Peleg Y, Rappoport N, Linial M. ProteinBERT: a universal deep-learning model of protein sequence and function. *Bioinformatics.* 2022;38(8):2102–10. <https://doi.org/10.1093/bioinformatics/btac020> PMID: [35020807](https://pubmed.ncbi.nlm.nih.gov/35020807/)
26. Chen J-Y, Wang J-F, Hu Y, Li X-H, Qian Y-R, Song C-L. Evaluating the advancements in protein language models for encoding strategies in protein function prediction: a comprehensive review. *Front Bioeng Biotechnol.* 2025;13:1506508. <https://doi.org/10.3389/fbioe.2025.1506508> PMID: [39906415](https://pubmed.ncbi.nlm.nih.gov/39906415/)
27. Zhang Y, Guan J, Li C, Wang Z, Deng Z, Gasser RB, et al. DeepSecE: A Deep-Learning-Based Framework for Multiclass Prediction of Secreted Proteins in Gram-Negative Bacteria. *Research (Wash D C).* 2023;6:0258. <https://doi.org/10.34133/research.0258> PMID: [37886621](https://pubmed.ncbi.nlm.nih.gov/37886621/)
28. Liao XJ, He TT, Liu LY, Jiang XL, Sun SS, Deng YH, et al. Unraveling and characterization of novel T3SS effectors in *Edwardsiella piscicida*. *mSphere.* 2023;8(5):e0034623. <https://doi.org/10.1128/msphere.00346-23> PMID: [37642418](https://pubmed.ncbi.nlm.nih.gov/37642418/)

29. Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics*. 2012;28(23):3150–2. <https://doi.org/10.1093/bioinformatics/bts565> PMID: 23060610
30. Hendrycks D, Gimpel K. Gaussian error linear units (GELUs) [Preprint]. arXiv:1606.08415 [cs.LG]. 2016 Jun 27 [revised 2023 Jun 6; version 5]. Available from: <https://arxiv.org/abs/1606.08415>
31. Madani A, Krause B, Greene ER, Subramanian S, Mohr BP, Holton JM, et al. Large language models generate functional protein sequences across diverse families. *Nat Biotechnol*. 2023;41(8):1099–106. <https://doi.org/10.1038/s41587-022-01618-2> PMID: 36702895
32. Teufel F, Almagro Armenteros JJ, Johansen AR, Gíslason MH, Pihl SI, Tsirigos KD, et al. SignalP 6.0 predicts all five types of signal peptides using protein language models. *Nat Biotechnol*. 2022;40(7):1023–5. <https://doi.org/10.1038/s41587-021-01156-3> PMID: 34980915
33. Ahmad S, Jose da Costa Gonzales L, Bowler-Barnett EH, Rice DL, Kim M, Wijerathne S, et al. The UniProt website API: facilitating programmatic access to protein knowledge. *Nucleic Acids Res*. 2025;53(W1):W547–53. <https://doi.org/10.1093/nar/gkaf394> PMID: 40331428
34. Jacobs RA, Jordan MI, Nowlan SJ, Hinton GE. Adaptive mixtures of local experts. *Neural Comput*. 1991;3(1):79–87. <https://doi.org/10.1162/neco.1991.3.1.79> PMID: 31141872
35. Dai D, Deng C, Zhao C, Xu RX, Gao H, Chen D, et al. DeepSeekMoE: towards ultimate expert specialization in mixture-of-experts language models [Preprint]. arXiv:2401.05639 [cs.CL]. 2024 Jan 11. Available from: <https://arxiv.org/abs/2401.05639>
36. Huang Q, An Z, Zhuang N, Tao M, Zhang C, Jin Y, et al. Harder tasks need more experts: dynamic routing in MoE models [Preprint]. arXiv:2403.07679 [cs.LG]. 2024 Mar 12. Available from: <https://arxiv.org/abs/2403.07679>
37. Liu H, Xia M, Gao T, Wang R, Chen D. Gating Dropout: Communication-efficient Regularization for Sparsely Activated Transformers. arXiv [Preprint]. 2022 May 27. Available from: <https://arxiv.org/abs/2205.14336>
38. Zhang Y, Cai R, Chen T, Zhang G, Zhang H, Chen P, et al. Robust mixture-of-expert training for convolutional neural networks [Preprint]. arXiv:2308.09751 [cs.CV]. 2023. Available from: <https://arxiv.org/abs/2308.09751>
39. Cho K, van Merriënboer B, Gulcehre C, Bahdanau D, Bougares F, Schwenk H, et al. Learning phrase representations using RNN encoder-decoder for statistical machine translation. arXiv. 2014. <https://doi.org/10.48550/arXiv.1406.1078>
40. Gehring J, Auli M, Grangier D, Yarats D, Dauphin YN. Convolutional sequence to sequence learning [Preprint]. arXiv:1705.03122 [cs.CL]. 2017 May 8 [revised 2017 Jul 25; version 3]. Available from: <https://arxiv.org/abs/1705.03122>
41. Zhou Y, Zhao X, Guo X, Li J, Liu S. Design of a modified transformer architecture based on relative position coding. *Int J Comput Intell Syst*. 2021;14(1):89–99. <https://doi.org/10.1007/s44196-023-00345-z>
42. Shaw P, Uszkoreit J, Vaswani A. Self-attention with relative position representations. In: *Proceedings of the 2018 Conference on Neural Information Processing Systems (NeurIPS 2018)*; 2018 Dec 2–8; Montréal, Canada. Curran Associates; 2018. p. 4644–52. <https://doi.org/10.18653/v1/N18-2074>
43. Wang J, Li J, Hou Y, Dai W, Xie R, Marquez-Lago TT, et al. BastionHub: a universal platform for integrating and analyzing substrates secreted by Gram-negative bacteria. *Nucleic Acids Res*. 2021;49(D1):D651–9. <https://doi.org/10.1093/nar/gkaa899> PMID: 33084862
44. Healy J, McInnes L. Uniform manifold approximation and projection. *Nature Reviews Methods Primers*. 2024;4(1):82. <https://doi.org/10.1038/s43586-024-00369-5>
45. Abby SS, Denise R, Rocha EPC. Identification of Protein Secretion Systems in Bacterial Genomes Using MacSyFinder Version 2. *Methods Mol Biol*. 2024;2715:1–25. https://doi.org/10.1007/978-1-0716-3445-5_1 PMID: 37930518
46. Tareen A, Kinney JB. Logomaker: beautiful sequence logos in Python. *Bioinformatics*. 2020;36(7):2272–4. <https://doi.org/10.1093/bioinformatics/btz921> PMID: 31821414
47. An Y, Wang J, Li C, Leier A, Marquez-Lago T, Wilksch J, et al. Comprehensive assessment and performance improvement of effector protein predictors for bacterial secretion systems III, IV and VI. *Brief Bioinform*. 2018;19(1):148–61. <https://doi.org/10.1093/bib/bbw100> PMID: 27777222
48. Rao R, Bhattacharya N, Thomas N, Duan Y, Chen X, Canny J, et al. Evaluating Protein Transfer Learning with TAPE. *Adv Neural Inf Process Syst*. 2019;32:9689–701. PMID: 33390682
49. Zhao Z, Hu Y, Hu Y, White AP, Wang Y. Features and algorithms: facilitating investigation of secreted effectors in Gram-negative bacteria. *Trends Microbiol*. 2023;31(11):1162–78. <https://doi.org/10.1016/j.tim.2023.05.011> PMID: 37349207
50. Spitz O, Erenburg IN, Beer T, Kanonenberg K, Holland IB, Schmitt L. Type I Secretion Systems—One Mechanism for All? *Microbiol Spectr*. 2019;7(2):10.1128/microbiolspec.psib-0003–2018. <https://doi.org/10.1128/microbiolspec.PSIB-0003-2018> PMID: 30848237
51. Tsirigotaki A, De Geyter J, Šoštarić N, Economou A, Karamanou S. Protein export through the bacterial Sec pathway. *Nat Rev Microbiol*. 2017;15(1):21–36. <https://doi.org/10.1038/nrmicro.2016.161> PMID: 27890920
52. Jia L, Zhu L. The Bacterial Type III Secretion System as a Broadly Applied Protein Delivery Tool in Biological Sciences. *Microorganisms*. 2025;13(1):75. <https://doi.org/10.3390/microorganisms13010075> PMID: 39858842
53. Kim H, Kubori T, Yamazaki K, Kwak M-J, Park S-Y, Nagai H, et al. Structural basis for effector protein recognition by the Dot/Icm Type IVB coupling protein complex. *Nat Commun*. 2020;11(1):2623. <https://doi.org/10.1038/s41467-020-16397-0> PMID: 32457311
54. Monjarás Feria J, Valvano MA. An Overview of Anti-Eukaryotic T6SS Effectors. *Front Cell Infect Microbiol*. 2020;10:584751. <https://doi.org/10.3389/fcimb.2020.584751> PMID: 33194822
55. Su J, Han C, Zhou Y, Shan J, Zhou X, Yuan F. SaProt: Protein language modeling with structure-aware vocabulary. In: *Proceedings of the International Conference on Learning Representations (ICLR 2024)*; 2024 May; Kigali, Rwanda. ICLR Press; 2024. Poster Spotlight.