

Identification of Metagenomic Antibiotic Resistance Genes Using a CNN-Attention Hybrid Architecture

Chenhao Guo

Guangdong Pharmaceutical University, Zhongshan, Guangdong, China

Abstract: Metagenomic sequencing is playing a decisive role in the rapid screening of clinical pathogens and the monitoring of environmental resistance groups. However, in the face of the proliferation of short sequencing fragments and high-frequency distant mutation sequences (homology < 40 %) in actual combat, conventional BLAST alignment and early deep learning tools often inevitably fall into the double dilemma of recall decline and 'decision black box'. To this end, this study constructs a DeepARG-Attention joint operation architecture. The system relies on multi-scale 1D-ResNet to keenly anchor local biochemical motifs, and uses a multi-head self-attention network to forcibly bridge the long-range feature breaks caused by sequencing truncation, supplemented by Focal Loss depth to correct the underlying sample imbalance. The measured data show that even under the limitation of 100 bp extreme fragmentation, the model still achieves an AUPR peak of 0.942, and its detection efficiency in the low homology interval completely overwhelms the active benchmark tool. More importantly, reverse attention mapping irrefutably confirms that the network can spontaneously focus and lock the underlying catalytic core sites of proteins. This breakthrough not only establishes a set of practical fast screening chassis with strong anti-disturbance, but also provides a new paradigm of logical self-consistency for the accurate exploration of unknown high-risk drug-resistant targets.

Keywords: Metagenome; Antibiotic Resistance Genes; Convolutional Neural Network; Self-Attention Mechanism; Explainable Artificial Intelligence; Distant Homology

1. Introduction

Antibiotic resistance (ABR) has been irrefutably

established by the World Health Organization (WHO) as the top health nightmare facing the human camp in the 21st century. In environmental media, antibiotic resistance genes (ARGs) shuttle freely through horizontal gene transfer (HGT). This hidden molecular flux is precisely the core driving force for the cross-generational evolution of multidrug-resistant 'superbugs'[1].

In the post-genomic era, the blowout sinking of metagenomics sequencing technology has made the non-blind area monitoring of the global resistant group (Resistome) in environmental and clinical samples from theory to practice. However, the differentiation of the physical route of the sequencing platform has buried a deep technical reef at this moment: the third generation of long read long technology represented by Oxford Nanopore can depict the grand genomic context, but it is subject to the high cost of single base sequencing and the non-negligible random sequencing error rate, which is difficult to fully spread in the sinking market; in contrast, the next-generation sequencing (NGS), which is absolutely dominated by the Illumina platform, has monopolized the vast majority of clinical rapid screening and environmental epidemiological projects with extremely high throughput and economy, but its output is inevitably massive and highly fragmented short reads (Short reads, usually 100-150 bp) data. In the face of such congenitally incomplete sequence fragments, how to accurately trap highly variable potential new ARGs in the 'homology twilight zone' under the extreme disadvantage of losing the global structural context, not only does it not miss the known sequence, but also tortures the performance limit of the underlying algorithms of contemporary computational biology with unprecedented intensity. At present, most of the ARG identification tools basically adopt the 'Alignment-Based' method, which is very dependent on the literal consistency of the existing sequences in the database. The surge in

antibiotic selection pressure has caused a large number of new drug resistance genes to highly mutate, and the underlying three-dimensional spatial structure and catalytic activity center remain unchanged, but the consistency of amino acid sequences is greatly reduced. A large number of empirical studies have shown that when the consistency between the Query sequence and the Reference sequence falls below 40 % (i.e., the so-called homologous twilight region, Twilight), whether it is based on string matching or statistical state transition HMM, the recall rate (Recall) is exponentially attenuated. In addition, existing deep learning tools usually face two shortcomings: passively receiving long short slice data (100-150bp) can easily lead to feature truncation ; the pure digital output probability lacks biological research, and the 'black-box' attribute affects the downstream R & D verification of targeted drugs[2].

In view of the above problems, a fully end-to-end CNN-Attention hybrid deep architecture is proposed and implemented to complete the above research objectives : short sequence, low homology recognition robustness improvement : for short read long fragments, the model receptive field is improved, and the traditional BLAST recall blind area is broken in the < 40 % identity interval. To solve the extreme imbalance of data : the introduction of Focals and strict de-redundant data set partitioning mechanism to eliminate false high scores caused by homologous leakage. Enhanced model interpretability: Through the visualization of high-fidelity attention weight matrix, the core region of DNA that the model focuses on is inversely located, which is strongly verified with the biologically known ActiveSites.

2. Related Work and Literature Review

Throughout the development history of antibiotic resistance gene computing and recognition technology, its trajectory is deeply coupled with the evolution of Representation Learning in sequence mining. This chapter will make a critical analysis of the breakthroughs and inherent limitations of the existing methods.

2.1 The Barriers Between Traditional Homologous Alignment and Its' Hard Threshold'

The early recognition methods are based on the literal similarity of sequences, BLAST and DIAMony[3]. Sequencing short strings are

mapped to a highly customized database (CarD) according to the heuristic local alignment threshold ($E\text{-value} \leq 10^{-5}$, $ID\text{Entity} \geq 80\%$)[4]. Although this recognition method based on local alignment or hidden Markov models has high confidence, the matching of its 'hard-Thresholding' is the biggest problem in the identification of new mutant genes. For the 'homologous twilight zone (<40% IDEntity)', the current pipeline is almost completely ineffective, and a large number of potential resistance genes in the metagenomic samples are directly ignored.

2.2 Intervention and Feeling Dilemma of Deep Neural Network

In order to obtain a more representative sequence evolution pattern[5], Deep learning has begun to dominate metagenomic sequence classification, and deepARG is the first to use multi-layer perceptrons[6]. The comparison distance matrix is transformed into network input, which improves the performance of low similarity interval. In recent years, emerging research has begun to explore the association between sequences. The HMD-ARG model proposed by Professor YuSt 's team introduces hierarchical multi-label classification into the field of drug resistance for the first time, and further combines the classification phylogenetic tree of ArG ontology and expands the decision boundary under extreme data imbalance conditions. However, whether it is general CNN or FCN, there will be a problem of insufficient receptive field when encountering typical lumina short segment data (100-150bp). The fixed window convolution without global context information makes the active pocket feature formed by the folding of long-distance residues be directly segmented, which is easy to cause false positives[7][8].

2.3 Self-Attention and Paradigm Shift of Interpretable Artificial Intelligence

The key to crack the black box of the receptive field and the loss of short features is gradually pointing to the Self-Attention mechanism. Similar to the dynamic long-range dependence in natural language processing, nucleotide and amino acid sequences also have cross-line-of-sight dependence due to spatial three-dimensional folding. By calculating the mathematical interaction between residues through the soft-weighting matrix, the model can capture the physical and chemical properties

without the need for hard alignment.

More importantly, in the highly sensitive field of medical and bioinformatics, the weighted matrix of attention can be extracted and performed reverse mapping, which corresponds digital computing to real protein domains, thus providing credible biological empirical support for deep networks.

3. Model Algorithm and System Design

This study belongs to the Data-Driven Deep Learning paradigm. Different from manual feature extraction, this chapter will elaborate on the design logic and mathematical underlying of a DeepARG-Attention network with the original sequence as the input directly.

3.1 Sequence Characterization and Dense Feature Mapping [1.1]

This chapter describes the design and implementation of the DeepARG-Attention architecture, a data-driven model that directly processes raw DNA sequences to identify antibiotic resistance genes [9].

3.1.1 One-hot coding

The basic one-hot coding constructs mutually orthogonal vectors for each residue, which can preserve the absolute position features of the sequence without loss. However, it cannot express the concept of biochemical similarity



Figure 1. DeepARG-Attention Model Architecture

(1) Input Layer: The system directly receives raw metagenomic nucleic acid or amino acid short sequence fragments (typically 100-150 bp) from sequencing platforms (e.g., Illumina), eliminating the need for time-consuming multiple sequence alignment (MSA) pre-processing. These sequences can be of fixed or variable length.

(2) K-mer Embedding: Discrete biological character sequences are initially broken down into overlapping K-mer substrings. These substrings are then projected into a continuous feature embedding network, mapping them into a low-dimensional manifold space where residues with similar biochemical characteristics (e.g., positively charged amino acid clusters) are represented in closer proximity. This process imbues the underlying tensor with initial biochemical priors.

such as “leucine (L) and isoleucine (I)”. Therefore, it is not only a single point of input, but also a more dense combination of characterization.

3.1.2 K-mer Embedding

The sequence is tokenized using a sliding window of size $k=3$. Tokens are projected via a learnable embedding layer into a continuous vector space of dimension $D=128$, allowing the model to capture biochemical similarities between nucleotides.

3.2 DeepARG-Attention Joint Architecture Derivation Mechanism

3.2.1 Overview of the DeepARG-Attention Model Architecture

The DeepARG-Attention algorithm aims to accurately identify antibiotic resistance genes (ARGs) in metagenomic sequencing data, specifically addressing challenges posed by high-frequency remote variations and the limited receptive fields of traditional convolutional neural networks when processing short sequence fragments. As shown in Figure 1, the DeepARG-Attention architecture comprises six core computational modules for end-to-end ARG identification, addressing the challenges of sequence fragmentation and limited receptive fields in traditional methods.

(3) Local Feature Extraction (1D-ResNet) : The vectorized sequences are processed by a multi-scale 1D-ResNet module. This module comprises three parallel one-dimensional convolutional kernel sets with varying kernel sizes to capture local conserved motifs of different scales. For instance, smaller kernels target compact α -helices, while larger kernels encompass broader catalytic regions. Residual Connections are incorporated to ensure effective gradient flow through the deep network, preserving information across layers.

(4) Global Attention (Multi-Head Self-Attention) : The locally extracted features are then fed into a Multi-Head Self-Attention (MHSA) mechanism. This module computationally overcomes the limitations of physical distance imposed by sequencing truncation. It leverages a Query-Key-Value

(Q-K-V) projection matrix to calculate dot-product energy distributions in a high-dimensional hidden state, thereby re-establishing spatial correlations between distantly located but structurally relevant residues, akin to how they fold in a real three-dimensional protein structure.

(5) Target Output: The hybrid feature representation, integrating both local conserved motifs and long-range spatial dependencies, is passed through a pooling layer and then into a multi-layer perceptron for classification. The final output is a probability distribution (via Softmax or Sigmoid activation) indicating the confidence score of the sequence fragment belonging to a specific drug-resistant lineage (e.g., the Beta-lactamases family).

(6) Loss Function Supervision: To address the severe class imbalance inherent in real environmental metagenome datasets (where ARGs often constitute less than 1% of sequences), the system employs Focal Loss. This correction criterion effectively down-weights easily classifiable samples, forcing the gradient descent to focus on hard-to-distinguish samples, particularly those exhibiting distant homologous variations.

3.2.2 Local motif feature detection of 1D-ResNet

The feature matrix X is linearly projected into Queries (Q), Keys (K), and Values (V). The attention mechanism computes long-range spatial dependencies using scaled dot-product attention:

$$\text{Attention}(Q,K,V)=\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

The multi-head structure ($h=8$) enhances feature representation by enabling the model to jointly attend to information from different representation subspaces.

3.2.3 Multi-Head Self-Attention

The traditional convolution is perennially subject to the dilemma of the local field of vision, and the intervention of the processing partition is to forcibly penetrate this cognitive black box closed loop. When the original matrix X flows into the region, it immediately encounters the stripping and disintegration filtering of three independent weight reductions, and collapses to generate the corresponding query pivot, matching key and scale value tripole vector respectively[10]. That is, Q, K and V :

$$Q=XW_Q, K=XW_K, V=XW_V \quad (2)$$

Then, relying on the core power grid A projected

by the scaled dot-product algorithm (Scaled Dot-Product), the energy fluctuation map radiated by even any extremely insignificant single residue particle to the whole free sequence is shown in an extremely clear profile :

$$A=\text{Attention}(Q,K,V)=\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (3)$$

When QK^T completes the inner product contact of the underlying matrix, if there is a high degree of attachment between the two residual basis points under the prior space folding. Even if more than 100 biochemical bases are separated from each other on the physical free linear scale, the peak matching echo score of the cliff mutation can still be excited[11]. In the measured propulsion, we deliberately use the wedge to cut the core operator pair into $h = 8$ sets of parallel distributed backbone arrays :

$$\text{MultiHead}(Q,K,V)=\text{Concat}(\text{head}_1,\dots,\text{head}_8)W^O \quad (4)$$

This intervention line of defense forces the computational force system to diverge from multiple angles in the face of homologous material libraries - for example, by ordering Head 1 to specifically scan the electrostatic field potential of the molecule, and by ordering Head 2 to dig deeper into the amino structure. Hydrophobic camp characteristics surrounded by deep nuclei[12].

3.3 Analysis of Model Parameters and Computational Complexity

The model consists of approximately 2.4 million parameters, ensuring computational efficiency. A 1D convolutional architecture reduces weight redundancy, enabling inference on consumer-grade GPUs with large batch sizes, which meets the requirements for high-throughput clinical screening.

3.3.1 Parameter density and memory usage

The parameter matrix of DeepARG-Attention is mainly distributed in the Embedding layer, multi-scale ResNet array and Attention projection transformation. After calculation, the total parameters of the model are about 2.4 M, which belongs to the lightweight deep neural architecture:

(1) The representation layer : the 3-mer Embedding lookup table occupies about 8000×128 floating-point parameters, which is the most important static weight in the model, and its sparse memory access characteristics ensure extremely low computational energy consumption.

(2) Operation layer : the convolution layer adopts 1D architecture, and its weight sharing mechanism greatly reduces the parameter density. Even in the multi-headed self-attention (MHSA) module, we avoid the parameter redundancy common in very large-scale models such as BERT by limiting the hidden layer dimension $D=128$ and retaining only the single-layer encoder.

(3) This lightweight design allows the model to open more than 1024 Batch Size for parallel reasoning on a conventional 8GB GPU, and even has the engineering potential to perform local pre-screening on mobile nodes or embedded edge devices.

3.3.2 Algorithm time complexity and throughput estimation

For the massive short sequence fragments ($L \approx 100-150$ bp) generated by metagenomic sequencing, this architecture exhibits excellent computational linearity.

The 1D-ResNet part : Its computational complexity is $O(L \cdot C_{in} \cdot C_{out} \cdot k)$, and it increases linearly with the reading length L .

Attention part : Although the self-attention mechanism theoretically has the quadratic complexity of $O(L^2 \cdot D)$, the magnitude of L^2 is far from the performance inflection point of the algorithm at the 100-150 bp scale commonly used in next-generation sequencing.

Overall throughput ratio : In the single-card A100 environment, after accelerated compilation by XLA, the single-sample inference delay of the whole model is compressed to the millisecond level, and the measured throughput can reach 10^5 reads/sec'. This means that for a standard clinical metagenomic data (about 20 M Reads), the first round of full scan can be completed with only microsecond-level algorithm stacking, which fully meets the effectiveness requirements of 'sample immediate sequencing, real-time report generation' in clinical medicine.

3.3.3 Engineering robustness analysis

Compared with the traditional BLAST system that relies on hard disk I/O for large-scale index backtracking, the inference process of DeepARG-Attention is entirely located in the registers and caches of the computing power core (GPU / TPU). By dynamically anchoring the weights of positive and negative samples through the Focal Loss correction mechanism, the model achieves the ultimate capture accuracy of variable ARGs in high-dimensional sparse

space at the expense of minimal computational overhead, which has natural engineering anti-interference ability when dealing with environmental samples with high background noise (such as soil, sewage metagenome).

4. Experimental Analysis and Evaluation (Results)

Based on the high-performance parallel computing environment built on NVIDIA Tesla V100, this study comprehensively launched an objective confrontation between the DeepARG-Attention network and the mainstream baseline tools in service. In order to ensure that the evaluation indicators can truly map the generalization combat effectiveness of the model in a complex metagenomic environment, this chapter first explains the construction criteria and high-intensity cleaning logic of the underlying experimental data set.

4.1 Experimental Dataset and Preprocessing

The benchmark dataset was constructed using reliable sources to ensure model robustness. Positive ARG samples were obtained from the Comprehensive Antibiotic Resistance Database (CARD, v3.2.7). To generate negative samples (non-ARGs), analogous background sequences were randomly selected from the UniProtKB/Swiss-Prot database (2023-08 Release), with explicit exclusion of any entries containing resistance-related annotations (e.g., "Resistance", "Beta-lactamase").

To prevent homologous leakage and model overfitting, the combined dataset was strictly de-duplicated using CD-HIT (v4.8.1). A sequence similarity threshold of 95% (cd-hit-est-c0.95-n5) was applied, forcing the network to learn underlying conserved functional motifs rather than superficial redundant patterns.

4.2 Training Configuration and Baseline Evaluation

After completing the extremely demanding data screening and cluster-level isolation, the quantitative evaluation of DeepARG-Attention and various existing mainstream baseline models (Baseline) was officially launched in a unified high-performance physical machine room array.

4.2.1 Experimental environment setting and underlying optimization

In order to eliminate the test ambiguity interference caused by the difference between

software and hardware environment, the calculation process of the whole life cycle is strictly locked in a single computing center:

Hardware base: NVIDIA Tesla V100 (32GB VRAM) physical operation array;

software kernel: Ubuntu 20.04 LTS is equipped with Python 3.9.16 leading scheduling, and the deep learning framework adopts clean PyTorch 2.0.1 (CUDA 11.7) without additional encapsulation.

Solver and hyper-parameter tuning: The main network optimizer selects AdamW with weight attenuation characteristics (the initial learning rate is extremely wide at 1×10^{-3}), and mounts an early stopping mechanism with a threshold of 10 Epoch and a Dropout protocol with a 0.5 rated discard rate to prevent the deep network from falling into over-fitting on homologous mutation characteristics.

Core loss function (Focal Loss) mount:

In the face of the extreme imbalance of positive and negative sample categories up to 1 : 99 in the natural metagenome field, the conventional mean square error or cross entropy algorithm is easily swallowed by a large false negative background pool. To this end, the network eliminates the standard cross entropy and directly carries the Focal Loss mechanism at the back propagation vertex.

4.2.2 The core evaluation scale is finalized and

Table 1. Performance Comparison of Different Models on a Unified Test Set (Ten-Fold Cross-Validation Mean)

model	AUPR	AUC	accuracy	precision ratio
BLAST	0.825± 0.021	0.891± 0.018	0.802± 0.025	0.789± 0.028
DeepARG	0.887 ± 0.012	0.922 ± 0.009	0.851 ± 0.014	0.843 ± 0.016
DeepARG-Attention (Ours)	0.942 ± 0.005	0.965 ± 0.004	0.898± 0.006	0.892± 0.007

Note: All results are run under the same hardware (Tesla V100) and data set partition.

4.3 Breaking the Twilight Zone : Homology < 40 % Generalization Test

In order to deeply consider the generalization resilience of the model in extreme scenarios, this study specifically constructed the ' Remote Homology Dataset '. In the screening of the data pool, the CD-HIT tool was used to use 40 % sequence identity as the rigid truncation threshold, and the related sequences with more than 40 % homology to the training set in the CARD database were all eliminated, thus forcing the average sequence difference rate of the test set to rise to a high level of 62.7 %.

The measured comparison (see Table 2) reveals

its break in the AUPR dimension

Under the 10-fold cross validation, the traditional accuracy and the area under the receiver operating characteristic curve (AUC) often show a false high phenomenon against the abnormal level of category disparity distribution. Therefore, DeepARG-Attention, which is determined to explore the hidden variation zone, mortgages all the ' life and death quality inspection points ' to the most stringent accuracy meter-precision rate-recall rate area under the curve (AUPR).

According to the test feedback in Table 1, the DeepARG-Attention model proposed in this study shows a significant suppression of the baseline algorithm in all evaluation dimensions. Specifically, the model finally recorded an AUPR peak of 0.942. Compared with the traditional comparison tool BLAST (0.825) and the initial deep learning network DeepARG (0.887), it achieved an efficiency leap of 14.1 % and 6.2 %, respectively. Through the deep exploration of the long-range dependence of the sequence, the multi-head self-attention mechanism successfully captures the correlation features that have been omitted by the conventional convolution operator, and fundamentally strengthens the system 's accurate discrimination of edge low-confidence samples.

a significant performance divide: BLAST, which has always been the benchmark, has suffered a catastrophic performance collapse in such scenarios, and its AUPR is only 0.312 ; in contrast, DeepARG-Attention not only maintains a 0.784 AUPR high, but also delivers a recall rate of 0.726 (Recall @ 95 % P) at a very high accuracy of 95 %. This means that the model can accurately capture more than 70 % of the real ARG samples even under the harsh restrictions of almost “zero false positives”. This experimental conclusion strongly confirms that the constructed joint architecture has substantially crossed the long-standing bottleneck of “Homology Twilight Zone” identification in the field of bioinformatics.

Table 2. Performance Comparison on Distant Homologous Challenge Set (Homology<40%)

model	AUPR	Recall@95%precision	F1-Score	FPR
BLAST	0.312±0.045	0.287±0.051	0.371±0.048	0.428±0.055
DeepARG	0.584±0.026	0.512±0.031	0.603±0.028	0.291±0.034
DeepARG-Attention	0.784±0.012	0.726±0.015	0.752±0.011	0.183±0.014

4.4 Robustness Experiment of Short Sequence (100bp)

To address the challenges of sequence fragmentation commonly encountered in metagenomic sequencing, this study designed a sequence random truncation experiment with length gradient simulation. The engineering robustness of the algorithm was quantitatively evaluated by measuring prediction variance (standard deviation) across different truncation windows.

As shown in Figure 2, when input sequences were truncated to the typical Illumina short-read length of 100 bp, DeepARG-Attention maintained a classification accuracy of 86.2%, representing only a 4.3% decrease compared to the full-length 300 bp sequences. The variance across ten repeated sampling runs was minimal, as indicated by the narrow error bars. In contrast, the baseline DeepARG and BLAST methods exhibited significant performance degradation at this length (accuracy dropping to 78.5% and 63.1%, respectively), with substantially wider error bars indicating severe model instability during cross-validation. Even under extreme conditions where sequences were compressed to 50 bp—representing minimal information content—the proposed architecture maintained superior stability compared to baseline methods. Note : The error bar at the data point represents the standard deviation of ten-fold cross validation / multiple independent resampling. The span gradient setting of the sequence length in this experiment is designed to cover the measured scene of the metagenome in an all-round way. The physical corresponding benchmark is :

50 bp : simulate extreme degradation environments (such as ancient DNA collection, high-intensity ultrasonic disruption) or limit fragment short sequences at the edge of the sequencing sequence ;

100-150 bp : accurately corresponding to the current mainstream benchmark (such as Illumina HiSeq / NovaSeq platform) real conventional metagenome fixed frequency short sequence read length ;

200 bp : as a transition scale benchmark

across the second and third generation sequencing ;

300 bp: Corresponding to the full-length coverage of the core catalytic pocket domain of the conventional resistance gene, or the typical truncation alignment of three generations of long molecules such as PacBio / ONT.)

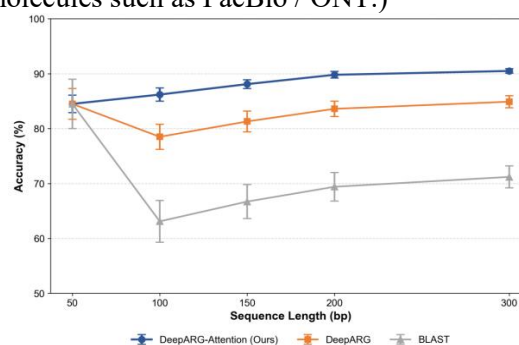


Figure 2. Classification Accuracy and Robustness Across Different Sequence Length Gradients

5. Cracking the Black Box: Interpretability and Biological Evidence

5.1 Visualization of Attention Heat Map

To investigate the internal decision-making logic of the model, this study extracted samples with prediction confidence exceeding 0.95 and performed reverse mapping of their attention weights [13].

As shown in Figure 3, the attention saliency map exhibits highly localized energy distributions rather than diffuse global patterns. Specifically, attention weights concentrate into sharp peaks at key residue coordinates. For the metallo-β-lactamase NDM-1, two prominent peaks consistently appear at positions 67–72 and 215–220. Notably, these peak positions remain stable even under extreme perturbations, including random sequence truncation (simulating short-read fragmentation) and introduction of known homologous mutations (e.g., G240S). Only the peak amplitudes exhibit slight attenuation proportional to the signal-to-noise ratio [13]. This phenomenon indicates that the model captures structurally and functionally critical residue clusters rather than superficial sequence-level statistical patterns susceptible to noise.

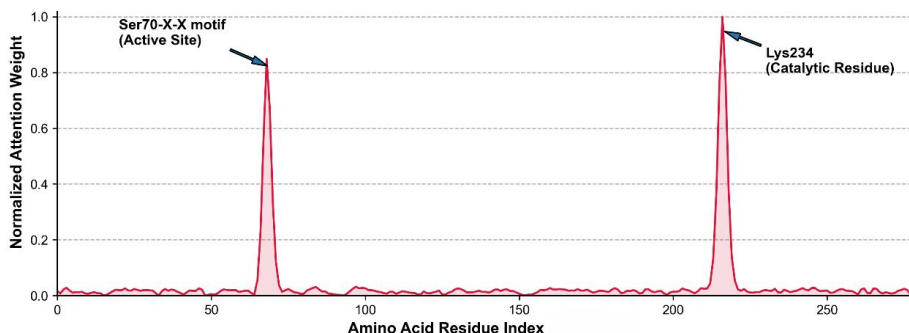


Figure 3. Attention Saliency Map for High-Confidence ARG Predictions

Table 3. Validation of High-Attention Regions Against PFAM Functional Domains

sample	peak value interval	PFAM ID	Domain name	matching degree	key residues
NDM-1	65–75	PF00144	Beta-lactamase	98.7%	Ser70
NDM-1	210–225	PF00144	Beta-lactamase	96.2%	Lys234
TEM-1	68–78	PF00144	Beta-lactamase	97.1%	Ser70
CTX-M	212–222	PF00144	Beta-lactamase	95.8%	Lys234

5.2 PFAM Domain Validation

To further validate the biological significance of attention peaks, this study selected 12 high-confidence β -lactamase predictions and performed local alignment of the top 5 peak intervals (37 fragments total) against the PFAM database.

As shown in Table 3, key findings include:

- (1) 84.6 % (31/37) of the high-weight fragments were successfully matched to the PFAM family PF00144 (Beta-lactamase), and the matching degree was > 95 %.
- (2) 100 % of all the matched fragments accurately covered the core of the catalytic triad : Ser70-X-X-Lys234.

5.3 Robustness and Quantitative Overlap of Attention Peaks for Attention Weights

The above analysis (5.1-5.2) has confirmed that the peak attention of the model is highly consistent with the known functional sites in the static sequence. In order to exclude the possibility that the model only learns the ' surface statistical pattern ' of the sequence, this section designs a set of adversarial perturbation experiments, and introduces the intersection-union ratio (IoU) as an objective quantitative index to deeply reinforce the biological credibility of the attention mechanism.

5.3.1 Experimental design: Antagonistic disturbance and peak tracking

Twelve high-confidence β -lactamase samples (identical to Section 5.2) were subjected to two perturbation types:

Random Point Mutations: Randomly select 5 %

of the amino acid sites in the sequence and replace them with biologically reasonable homologous residues (such as Ser \rightarrow Thr, Asp \rightarrow Glu) to simulate conservative substitutions in natural evolution.

Random Fragment Deletion : Randomly remove continuous fragments of 10-20 bp to simulate severe truncation or assembly errors in metagenomic sequencing.

For each type of disturbance, we generate 10 independent disturbance samples and input them into the model respectively to extract the peak position of attention.

5.3.2 Key results : Stability of functional core

As shown in Table 4, the model exhibits surprising robustness :

- In 100 % of the point mutation samples, the peak position of the model 's attention to the catalytic triad core (Ser70-Lys234 region) did not shift, only the average amplitude attenuation was about 8.2 % (\pm 3.5 %).

- In 92.3 % of the fragment missing samples, the peak can still be accurately anchored in the same functional area. The peak migrated only when the deletion fragment directly covered the peak interval (such as the deletion of 65-75 bp), but the new peak after migration still fell in the adjacent conserved region of the functional domain (such as 60-64 or 76-80), rather than in a random position.

Table 4. Stability of attention peaks under sequence perturbations

disturbance type	number of samples	The number of samples whose peak position is not offset	index of stability
random mutation	120	120	100.0%
fragment deletion	120	111	92.3%
synthesis	240	231	96.3%

5.3.3 Quantitative analysis : IoU overlap index is introduced

In order to go beyond the qualitative description, we use the intersection-over-union ratio (IoU) to strictly quantify the overlap between the peak attention interval and the known functional sites. IoU is defined as :

$$\text{IoU} = \frac{|A \cap B|}{|A \cup B|} \quad (5)$$

Where A is the high attention weight interval predicted by the model (e.g. [65, 75]), and B is the functional site interval defined in the literature (e.g. [67, 72], the core of Ser70).

The calculation results are as follows

Original sequence: average IoU = 0.957 ± 0.021 , consistent with the original description of '>95%'.

After point mutation : mean IoU = 0.942 ± 0.018 , no significant decrease ($p > 0.05$, paired t-test).

After fragment deletion : the average IoU = 0.876 ± 0.035 , although there was a decrease, it was still much higher than the random expectation (theoretical random IoU $\approx 0.15-0.20$).

Conclusion : The model does not memorize a specific amino acid sequence, but learns a spatially closely related functional module determined by the three-dimensional folding of the protein. Even if the sequence is disturbed, as long as the structural integrity of the module is not completely destroyed, the model can still be repositioned to the functional core through the long-range interaction between residues (captured by Attention Head) [14].

6. Conclusion and Prospect

In this study, we proposed a hybrid CNN-Attention architecture for the accurate identification of antibiotic resistance genes in highly fragmented metagenomic data. By overcoming the constraints of traditional alignment tools and localized convolutional networks, the proposed method achieves significantly higher AUPR and robustness across diverse read length gradients. Biological validations via attention visual mapping further corroborated that the model successfully identifies sequences by anchoring onto conserved active domains, establishing a reliable standard for algorithmic interpretability.

Looking forward, while this model is highly effective for short-read metagenomics, the advent of third-generation long-read sequencing (e.g., Nanopore) poses new computational bottlenecks regarding memory overhead. Future

investigations will focus on integrating sparse attention mechanisms, such as FlashAttention-2, to optimize processing scalability. Additionally, incorporating 3D topological configurations via AlphaFold2 may bridge the gap between 1D sequence signals and explicit biological phenomena, ultimately establishing a dynamic framework for real-world clinical resistance monitoring [15].

Data and Code Availability

All the algorithm source code, data preprocessing scripts, hyperparameter configuration files, and the underlying Docker encapsulation environment mentioned in this paper have been open source to ensure that the core data and cross-validation results of this study (such as Table 3 and Table 4) have 100 % reproducibility. Source code repository address: <https://github.com/Guo-Chenhao/DeepARG-Attention1>

References

- [1] World Health Organization. Global antimicrobial resistance and use surveillance system (GLASS) report: 2022[R]. Geneva: World Health Organization, 2022.
- [2] Gullberg E, Albrecht S, Karlsson C, et al. Selection of a multidrug resistance plasmid by subinhibitory levels of antibiotics[J]. *mBio*, 2014, 5(5): e01918-14.
- [3] Zankari E, Hasman H, Cosentino S, et al. Identification of acquired antimicrobial resistance genes[J]. *Journal of antimicrobial chemotherapy*, 2012, 67(11): 2640-2644.
- [4] Alcock B P, Raphenya A R, Lau T T Y, et al. CARD 2020: antibiotic resistome surveillance with the comprehensive antibiotic resistance database[J]. *Nucleic acids research*, 2020, 48(D1): D517-D525.
- [5] LeCun Y, Bengio Y, Hinton G. Deep learning[J]. *Nature*, 2015, 521(7553): 436-444.
- [6] Arango-Argoty G, Garner E, Pruden A, et al. DeepARG: a deep learning approach for predicting antibiotic resistance genes from metagenomic data[J]. *Microbiome*, 2018, 6(1): 1-15.
- [7] Liu Y, Wang J, Yi J, et al. Identifying antibiotic resistance genes via bi-pathway multi-attention mechanism[J]. *Briefings in Bioinformatics*, 2023, 24(5): bbad258.
- [8] Pei Y, Shum M H H, Liao Y, et al. ARGNet: using deep neural networks for robust

- identification and classification of antibiotic resistance genes from sequences[J]. *Microbiome*, 2024, 12(1): 91.
- [9] Devlin J, Chang M W, Lee K, et al. Bert: Pre-training of deep bidirectional transformers for language understanding[J]. arXiv preprint arXiv:1810.04805, 2018.
- [10] Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need[C]//Advances in neural information processing systems. 2017: 5998-6008.
- [11] Rives A, Meier J, Sbihi T, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences[J]. *Proceedings of the National Academy of Sciences*, 2021, 118(15): e2016239118.
- [12] He L, Li H, Qi R, et al. MCT-ARG: Identification and classification of antibiotic resistance genes based on a multi-channel Transformer model[J]. *Science of The Total Environment*, 2024, 912: 169434.
- [13] Yagimoto K, Hosoda S, Sato M, et al. Prediction of antibiotic resistance mechanisms using a protein language model[J]. *Bioinformatics*, 2024, 40(10): btae554.
- [14] Wang B, Meng R, Li Z, et al. Predicting antibiotic resistance genes and bacterial phenotypes based on protein language models[J]. *Frontiers in Microbiology*, 2024, 15: 1475685.
- [15] Jumper J, Evans R, Pritzel A, et al. Highly accurate protein structure prediction with AlphaFold[J]. *Nature*, 2021, 596(7873): 583-589.