# Moss: Adaptive Inductive Bias Attention for Efficient and Robust Sequence Processing

Lingkai Hu, Wenkai Huang, *Member, IEEE*, Feng Zhan and Weiming Gan

*Abstract*—Sequence processing is a fundamental research area in artificial intelligence (AI) that encompasses various tasks and applications. Existing models—such as recurrent neural networks (RNNs) and Transformers—have drawbacks such as slow computation, high complexity, and overfitting. In this paper, we propose Moss, a novel sequence processing model that leverages the adaptive inductive bias (AIB)-attention mechanism. Moss can capture positional inductive relations more robustly and efficiently than Transformers. Moreover, it enables parallel computation and can handle extremely long sequences with near-linear complexity. We conduct extensive experiments on multiple data sets and tasks and evaluate Moss using various metrics. The results demonstrate that Moss achieves state-of-the-art (SOTA) performance on several language modeling benchmarks and surpasses other models in terms of accuracy, speed, and memory usage.

*Index Terms*—Attention mechanism, complexity, inductive bias, language modelling, sequence processing

## I. INTRODUCTION

SEQUENCE processing is a fundamental research area in artificial intelligence (AI) that encompasses tasks such as sequence modeling [1]–[5], classification [6]–[8], and translation [1], [9]–[12]. These tasks have a wide range of applications in natural language processing (NLP), speech recognition [13], music synthesis [14], chatbots [15], machine translation, and more. Numerous influential models have been proposed to address these tasks, such as recurrent neural networks (RNNs) [16] and Transformers [12].

However, these models also have inherent drawbacks. RNNs operate sequentially along the sequence order, thereby resulting in slow computation and memory loss for long sequences [16], [17]. Transformers employ a self-attention mechanism that can handle long sequences, but its time and space complexity are quadratic functions of the input sequence length, thereby requiring a substantial amount of computational resources. These limitations hinder their performance on long sequence tasks [18]–[20]. Although numerous attempts have been made to overcome these limitations [21]–[32], the effect is marginal and the challenge of long sequences remains.

Attention mechanisms are essential for numerous sequence processing models [33]–[35], such as Transformers. However, unlike models that incorporate structural priors in the input, attention mechanisms have to learn the structural relations

Lingkai Hu, Wenkai Huang, Feng Zhan and Weiming Gan are with the School of Mechanical and Electrical Engineering, Guangzhou University, 510006 Guangzhou, China. E-mail: 2112107060@e.gzhu.edu.cn, smallkat@gzhu.edu.cn, 2007300008@e.gzhu.edu.cn, 2112107050@e.gzhu.edu.cn.

(Corresponding author: Wenkai Huang.)

from the data. This makes Transformer-like models prone to overfitting on small to medium-sized data sets [18], [20], [36], [37].

In this study, we propose Moss, a novel sequence processing model based on the adaptive inductive bias (AIB)-attention mechanism. Moss can effectively capture positional inductive relations with higher robustness and lower time and space complexity. In addition, It enables parallel computation and can build long-distance dependency relations with a complexity of $O(L \log_2 L)$.

The following are the main contributions of this paper:

- We introduce the AIB-attention mechanism, a novel technique that leverages positional inductive bias to enhance the learning ability of attention-based models.
- We present Moss, a new sequence processing model that incorporates the AIB-attention mechanism into its architecture. We analyze its advantages over existing models in terms of complexity, robustness, and generalization.
- We conduct extensive experiments on multiple data sets and tasks and evaluate Moss using multiple metrics, such as accuracy, calculation speed, and peak memory usage. In the character-level language modeling task, Moss easily achieved better results than GPTs of the same size in data set benchmarks, including text8 [39] and enwik8 [40]. In the enwik8 character-level language modeling task, this model exceeds the current SOTA. The results reveal that Moss outperforms other models on these metrics and can effecitvely handle long sequences.

## II. BACKGROUND

### A. Attention mechanism

Attention mechanism is a special structure that can be embedded into machine learning models to automatically learn and compute the importance of input data for output data [41], [42]. It draws inspiration from human visual attention, which selectively focuses on parts of images or texts [43], [44]. Mathematically, the attention mechanism can be considered a function that maps a query and a set of key-value pairs to an output, where the query, keys, values, and output are all vectors. Moreover, the attention mechanism can enhance the efficiency and performance of models in various fields, such as NLP and image processing [12], [45]–[47].

Encoder-decoder is a common framework that incorporates the attention mechanism [10], [48], [49]. It can transform one sequence (such as an audio signal) into another sequence (such as a text transcript). The encoder encodes the input

TABLE I
SEVERAL PERFORMANCE INDICATORS OF DIFFERENT TYPES OF NEURAL LAYERS.

| Layer Type | Complexity per Layer | Maximum Path Length | Receptive Field |
|---|---|---|---|
| Self-Attention | $O(L^2D)$ | $O(1)$ | $O(L)$ |
| AIB-Attention | $O(LD\log_2 L)$ | $O(1)$ | $O(L)$ |
| Fully Connected | $O(L^2D^2)$ | $O(1)$ | $O(L)$ |
| Convolutional | $O(KLD^2)$ | $O(\log_K L)$ | $O(K)$ |
| Recurrent | $O(LD^2)$ | $O(L)$ | $O(L)$ |

sequence into a semantic vector, and the decoder generates the output sequence based on this vector. The attention mechanism enables the decoder to dynamically adjust the attention weight to the encoder output, thereby improving the accuracy and fluency of speech recognition.

The self-attention mechanism is a classic implementation of the attention mechanism. It enables each element in the input sequence to interact with other elements, thereby capturing the long-distance dependency relationship within the sequence. The self-attention mechanism can also reduce the computational complexity and memory consumption compared to other attention mechanisms. The core component of the Transformer model is the self-attention mechanism. In the field of NLP, the Transformer model has achieved numerous breakthrough results in tasks such as machine translation [50], text summarization [51], question answering [52], and natural language generation [53].

*B. Transformers*

Transformers are a popular neural network architecture that aim to solve sequence-to-sequence tasks while handling long-range dependencies with ease and are now a state-of-the-art technique in the field of NLP. A Transformer neural network can take an input sentence in the form of a sequence of vectors and convert it into a vector called an encoding and then decode it back into another sequence. An important part of the Transformer model is the attention mechanism, which allows the processing of one input word to include relevant data from certain other words, while masking the words that do not convey relevant information. Transformer neural networks have achieved remarkable results in various NLP tasks.

In the original Transformer model, the multihead attention mechanism is calculated in the following manner:

$$Q_i = XW_i^Q \tag{1}$$

$$K_i = CW_i^K \tag{2}$$

$$V_i = CW_i^V \tag{3}$$

$$head_i = Softmax\left(\frac{Q_iK_i^T}{\sqrt{d_k}}\right)V_i \tag{4}$$

$$Y = Attention\left(X, C\right) = Concat\left(head_1, ..., head_h\right). \tag{5}$$

It contains two inputs: the query sequence $X \in \mathbb{R}^{L \times d}$ with length $L$ and the context sequence $C \in \mathbb{R}^{m \times d}$ with length $m$. The output sequence $head_i \in \mathbb{R}^{L \times d_k}$ with length $L$, which is the same as $X$. Here, $d$ represents the encoding dimension, $d_k$ represents the dimension of the key, $h$ represents the number

of heads, and generally $d_k = d/h$. Moreover, $W_Q$, $W_K$, $W_V \in \mathbb{R}^{d \times d_k}$ are three learnable parameters. The output $Y \in \mathbb{R}^{L \times d}$ of the attention mechanism is the concatenation of multiple heads in the last dimension. When calculating self-attention, $C = X$; when calculating cross-attention, $C$ and $X$ have different sources. Generally, $C$ is the output of the Transformer encoder and $X$ is the input of the decoder.

In the attention mechanism, each query vector needs to perform multiplication with all keys in the sequence, while each query vector in the attention matrix $Softmax\left(\frac{QK^T}{\sqrt{d_k}}\right) \in \mathbb{R}^{L \times m}$ also needs to perform multiplication with all values. This results in a time and space complexity of $O\left(LmD\right)$ for Transformers. In AIB, the complexity is reduced to $O\left(LDlog_2L\right)$, which significantly improves the computational efficiency of the neural network.

Table I summarizes the complexity, maximum path length, and receptive field of various commonly used neural network layers. From this, the performance indicators of different types of layers can be clearly compared.

*C. Inductive Bias*

Inductive bias refers to the necessary assumptions that a machine learning algorithm makes about the target function of a learning problem [54], [55]. In machine learning, we want to construct algorithms that can predict the output based on the given input. However, for inputs that we have not encountered before, we cannot determine the output value without any additional assumptions. Therefore, we need to make a few assumptions about the properties of the target function, and these assumptions constitute the inductive bias. Inductive bias originates from the no free lunch theorem, which states that no learning algorithm can perform well on all possible problems without any prior knowledge [56].

Inductive bias can help deep learning algorithms converge faster and generalize better, while reducing computational and parameter costs. Inductive biases can be categorized into two different groups—explicit and implicit. The former represents the assumptions that are explicitly encoded in the model architecture or the learning algorithm, while the latter represents the assumptions that are implicitly induced by factors such as initialization, optimization, and regularization [55], [57].

In deep learning, many neural network structures are designed based on explicit inductive biases. For example, convolutional neural networks (CNNs) assume local connectivity and translation invariance, recurrent neural networks (RNNs) assume temporal dependency and sequential order, and graph neural networks (GNNs) assume relational structure

and permutation invariance. These structures mainly model local structural relationships, while neglecting the modeling of global relational patterns. On the other hand, neural network structures with less explicit inductive biases, such as fully connected and Transformer layers, can capture more contextual information but may have higher complexity and be more prone to overfitting [56].

In this paper, we introduce the main idea and structure of the Moss model, which is a novel sequence processing model that combines explicit and implicit inductive biases to achieve superior robustness and fitting ability. We compare the Moss model with the Transformer model theoretically and empirically, and show that the Moss model has a lower complexity and faster computation speed than that of the Transformer model, while also being able to perform parallel computation. We also apply the Moss model to several NLP tasks, and demonstrate its effectiveness and superiority over the Transformer model. We hope that our study can inspire more research on improving attention mechanism and sequence processing models.

## III. Moss

In this section, we introduce our proposed Moss model. The core structure of Moss is AIB-Attention. AIB-Attention applies adaptive information induction to attention computation, thereby allowing the model to dynamically establish inductive bias relationships in the data. A notable feature of AIB-Attention is its use of a weighted attention mechanism to fuse contextual information. We demonstrate that by applying this structure, we can achieve lower time and space complexity while maintaining excellent robustness and fitting ability.

### A. AIB-Attention

We refer to the attention computation in Moss as AIB-attention, which involves two steps: 1) Calculate the attention matrix and values matrix through linear transformation. 2) Combine the attention matrix and the value matrix using a special weighted prefix sum algorithm to produce the output.

This paragraph introduces the mechanism of the AIB-attention computation. Given that a sequence $X \in \mathbb{R}^{L \times d}$ is the input to AIB-attention. Then, we compute the attention matrix $A \in \mathbb{R}^{L \times d}$ and the value matrix $V \in \mathbb{R}^{L \times d}$ in the following manner:

$$A = XW^A + b_1 \tag{6}$$
$$V = XW^V + b_2, \tag{7}$$

where $W^A \in \mathbb{R}^{d \times d}$, $W^V \in \mathbb{R}^{d \times d}$, $b_1 \in \mathbb{R}^d$ and $b_2 \in \mathbb{R}^d$ are learnable parameters. Then, we calculate the attention output for each time step based on $A$ and $V$. In the decoder, the $i$-th output can only access information from the first $i$ inputs. Therefore,

$$O_i = Sum(Softmax(A[:i]) \bullet V[:i]) \tag{8}$$
$$O = [O_1, ..., O_L] \tag{9}$$
$$Y = OW^O + b_3, \tag{10}$$

where $W^O \in \mathbb{R}^{d \times d}$ and $b_3 \in \mathbb{R}^d$ are learnable parameters, $O_i \in \mathbb{R}^d$ denotes the $i$-th vector of the output matrix, $O \in \mathbb{R}^{L \times d}$, $A[1:i] \in \mathbb{R}^{i \times n}$ denotes the matrix formed by the first $i$ vectors of $A$, and $V[1:i] \in \mathbb{R}^{i \times n}$ denotes the matrix formed by the first $i$ vectors of $V$. In the encoder, the output must capture the complete context. Therefore,

$$\overrightarrow{O_i} = Sum(Softmax(A[:i,:d|2]) \bullet V[:i,:d|2]) \tag{11}$$
$$\overleftarrow{O_i} = Sum(Softmax(A[i:,d|2:]) \bullet V[i:,d|2:]) \tag{12}$$
$$\overrightarrow{O} = \left[\overrightarrow{O_1}, ..., \overrightarrow{O_L}\right] \tag{13}$$
$$\overleftarrow{O} = \left[\overleftarrow{O_1}, ..., \overleftarrow{O_L}\right] \tag{14}$$
$$O = Concat(\overrightarrow{O}, \overleftarrow{O}) \tag{15}$$
$$Y = OW^O + b_3, \tag{16}$$

where $\overrightarrow{O} \in \mathbb{R}^{L \times (d|2)}$ and $\overleftarrow{O} \in \mathbb{R}^{L \times (d|2)}$ are tensors of forward and backward attention calculation results, respectively; $\overrightarrow{O_i} \in \mathbb{R}^{d|2}$ is the $i$-th vector of $\overrightarrow{O}$ and only contains information at or before the $j$-th time step; and $\overleftarrow{O_i} \in \mathbb{R}^{d|2}$ is the $i$-th vector of $\overleftarrow{O}$ and only contains information at or after the $j$-th time step. The output matrix $O \in \mathbb{R}^{L \times d}$ is the concatenation of $\overrightarrow{O}$ and $\overleftarrow{O}$. In this manner, the output vector of each time step, $O_i$, can capture the complete contextual information.

In the above computation, the softmax operation is performed at each time step, which would require too much computation. Therefore, we use the following method to replace the calculation of $O$ above. In decoder,

$$E = Exp(A) \tag{17}$$
$$H_i = \sum_{j=1}^{i} (E_j \bullet V_j) \tag{18}$$
$$U_i = \sum_{j=1}^{i} E_j \tag{19}$$
$$H = [H_1, ..., H_L] \tag{20}$$
$$U = [U_1, ..., U_L] \tag{21}$$
$$O = \frac{H}{U}. \tag{22}$$

Let $H \in \mathbb{R}^{L \times d}$ and $U \in \mathbb{R}^{L \times d}$ be temporary variables in the computation. $H_i$ denotes the weighted sum of the first $i$ vectors of $V$, and $U_i$ denotes the denominator of the $i$-th Softmax computation. In this computation, $H$ can be expressed as the prefix sum of $EV$ and $U$ as the prefix sum of $E$. To enhance the parallelism of computation, we propose the following algorithm for prefix sum computation:

$$H = E \bullet V \tag{23}$$
$$U = E \tag{24}$$
$$For \ j = 0 : \lceil log_2 L \rceil \tag{25}$$
$$H[2^j :] + = H[: -2^j] \tag{26}$$
$$U[2^j :] + = U[: -2^j] \tag{27}$$
$$End \ For \tag{28}$$

$$O = \frac{H}{U}. \tag{29}$$

In each iteration of the loop operation, we update the data with information from the original data, which are separated by $2^j$ elements. To enable the neural network to learn the relationship between the data at different distances, we introduce learnable weights in this step. We call this technique AIB. The algorithm implementation is given below:

$$w = \left[ w_1, ..., w_{\lceil log_2 Max(L) \rceil} \right] \tag{30}$$

$$W = Exp\left( Prefix\ sum\left( w \right) \right) \tag{31}$$

$$For\ j = 0 : \lceil log_2 L \rceil \tag{32}$$

$$H[2^j :] += H[: -2^j] \bullet W_{j+1} \tag{33}$$

$$U[2^j :] += U[: -2^j] \bullet W_{j+1} \tag{34}$$

$$End\ For \tag{35}$$

$$O = \frac{H}{U}. \tag{36}$$

In the above formula, $w \in \mathbb{R}^{\lceil log_2 Max(L) \rceil \times d}$ is a trainable parameter, where each sub-vector is independent in the calculation. By this method, in each iteration, $H_i$ and $U_i$ are weighted and added with the information of their current $2^j$ time step. In this manner, the attention mechanism can establish AIB, thereby scaling the reception of distant information. The algorithm we use to calculate $Y$ is represented as Algorithm 1.

---

**Algorithm 1** AIB-Attention in decoder

---

**input** : $X \in \mathbb{R}^{L \times d}$
**output:** $Y \in \mathbb{R}^{L \times d}$
$A, V = XW^A + b_1, XW^V + b_2$
$U = E = \exp(A)$
$H = E \bullet V$
$W = \exp(Profix\ sum(w))$
**for** $j \leftarrow 0$ **to** $\lceil log_2 L \rceil$ **do**
  $H[2^j :] += H[: -2^j] \bullet W_{j+1}$
  $U[2^j :] += U[: -2^j] \bullet W_{j+1}$
**end**
$O = \frac{H}{U}$
$Y = OW^O$

---

Similarly, in the encoder,

$$E = Exp\left( A \right) \tag{37}$$

$$\overrightarrow{H} = E\left[ :, : d|2 \right] \bullet V\left[ :, : d|2 \right] \tag{38}$$

$$\overleftarrow{H} = E[:, d|2 :] \bullet V[:, d|2 :] \tag{39}$$

$$\overrightarrow{U} = E\left[ :, : d|2 \right] \tag{40}$$

$$\overleftarrow{U} = E[:, d|2 :] \tag{41}$$

$$\overrightarrow{w} = \left[ \overrightarrow{w}_1, ..., \overrightarrow{w}_{\lceil log_2 Max(L) \rceil} \right] \tag{42}$$

$$\overleftarrow{w} = \left[ \overleftarrow{w}_1, ..., \overleftarrow{w}_{\lceil log_2 Max(L) \rceil} \right] \tag{43}$$

$$\overrightarrow{W} = Exp\left( Prefix\ sum\left( \overrightarrow{w} \right) \right) \tag{44}$$

$$\overleftarrow{W} = Exp\left( Prefix\ sum\left( \overleftarrow{w} \right) \right) \tag{45}$$

$$For\ j = 0 : \lceil log_2 L \rceil \tag{46}$$

$$H[2^j :] += H[: -2^j] \bullet W_{j+1} \tag{47}$$

$$H[: -2^j] += H[2^j :] \bullet W_{j+1} \tag{48}$$

$$U[2^j :] += U[: -2^j] \bullet W_{j+1} \tag{49}$$

$$U[: -2^j] += U[2^j :] \bullet W_{j+1} \tag{50}$$

$$End\ For \tag{51}$$

$$O = Concat(\frac{\overrightarrow{H}}{\overrightarrow{U}}, \frac{\overleftarrow{H}}{\overleftarrow{U}}). \tag{52}$$

The algorithm we use to calculate $Y$ in encoder is represented as Algorithm 2:

---

**Algorithm 2** AIB-Attention in encoder

---

**input** : $X \in \mathbb{R}^{L \times d}$
**output:** $Y \in \mathbb{R}^{L \times d}$
$A, V = XW^A + b_1, XW^V + b_2$
$E = \exp(A)$
$\overrightarrow{U} = E[:, : d|2]$
$\overleftarrow{U} = E[:, d|2 :]$
$\overrightarrow{H} = E[:, : d|2] \bullet V[:, : d|2]$
$\overleftarrow{H} = E[:, d|2 :] \bullet V[:, d|2 :]$
$\overrightarrow{W} = \exp(Profix\ sum(\overrightarrow{w}))$
$\overleftarrow{W} = \exp(Profix\ sum(\overleftarrow{w}))$
**for** $j \leftarrow 0$ **to** $\lceil log_2 L \rceil$ **do**
  $\overrightarrow{H}[2^j :] += \overrightarrow{H}[: -2^j] \bullet \overrightarrow{W}_j$
  $\overleftarrow{H}[: -2^j] += \overleftarrow{H}[2^j :] \bullet \overleftarrow{W}_j$
  $\overrightarrow{U}[2^j :] += \overrightarrow{U}[: -2^j] \bullet \overrightarrow{W}_j$
  $\overleftarrow{U}[: -2^j] += \overleftarrow{U}[2^j :] \bullet \overleftarrow{W}_j$
**end**
$O = concat(\frac{\overrightarrow{H}}{\overrightarrow{U}}, \frac{\overleftarrow{H}}{\overleftarrow{U}})$
$Y = OW^O$

---

Using the above method, we achieved efficient parallel computation of $Y$.

### B. Feed-Forward Networks

We implemented a feed-forward layer after the self-attention layer. This layer comprises two linear transformations with a Leaky-ReLU activation [58] in between.

$$Leaky\ ReLU\left( X \right) = Max\left( 0.01X,\ X \right) \tag{53}$$

$$FFN\left( X \right) = Leaky\ ReLU\left( XW_1 + b_1 \right) W_2 + b_2. \tag{54}$$

We selected Leaky-ReLU as the activation function for this structure because it maintains a non-zero gradient in the negative region, which prevents the issue of dying neurons. The dimensionality of input and output is $d$, and the dimensionality of the inner layer is $d_{ff}$.

### C. Model Architecture

Neural sequence transduction models typically use an encoder–decoder structure. Moss also adopts this architecture, with stacked self-attention and point-wise fully connected
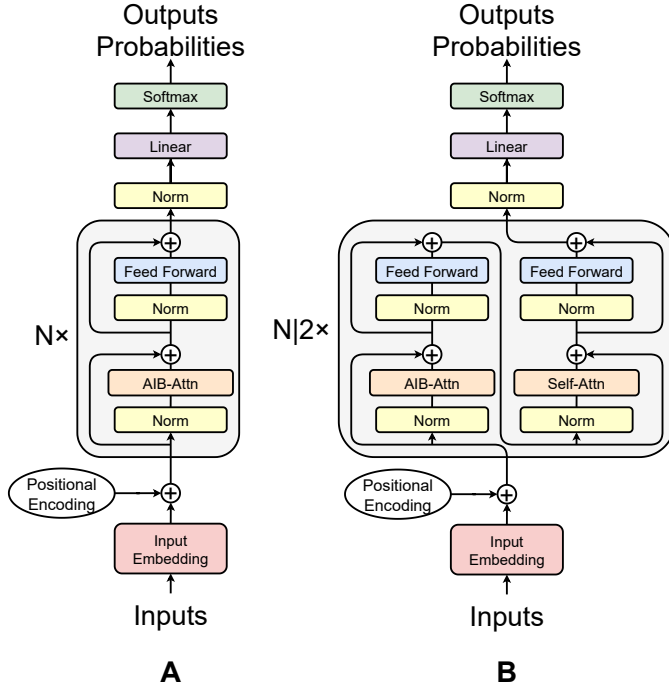
Fig. 1. The Architecture of the Moss Model.

TABLE II
HYPERPARAMETER OF THE NEURAL NETWORK IN THE LRA
EXPERIMENT.

| Task | $N$ | $d_{model}$ | $d_{ff}$ | Batch size | Iterations | Epochs |
|---|---|---|---|---|---|---|
| ListOps | 4 | 512 | 1024 | 32 | 5000 | \ |
| Text | 4 | 256 | 1024 | 32 | 20000 | \ |
| Retrieval | 4 | 128 | 512 | 32 | 5000 | \ |
| Image | 1 | 32 | 64 | 256 | \ | 200 |
| Pathfinder | 1 | 32 | 64 | 512 | \ | 200 |
| Path-X | 1 | 32 | 64 | 64 | \ | 200 |

layers in both the encoder and decoder. Fig. 1 shows its encoder and decoder structures, where Fig. 1A is the encoder and Fig. 1B is the decoder:

**Encoder:** We use randomly initialized parameters as positional encoding. The encoder consists of N layers, with each layer containing two sub-layers. The first sub-layer employs an AIB-attention encoding form, while the second sub-layer implements a feed-forward neural network. Both sub-layers adopt normalization and residual connections. Moreover, the model preserves the same embedding dimension across all layers to ensure consistent residual connections.

**Decoder:** The decoder also consists of N layers, each containing two sublayers. In the odd-numbered layers, the first sublayer is a unidirectional AIB-attention that models only the forward context. The second sublayer is a feedforward network that applies a nonlinear transformation to the output of the self-attention layer.

The decoder interleaves the AIB-attention and the self-attention layer because better accuracy is achieved by this strategy than using only the AIB-attention or only the self-attention layer in our experiments. It is conjectured that the AIB-attention excels at modeling relative-position dependency

whereas the self-attention layer excels at computing long-term relevance. This mixed strategy leverages the strengths of both structures.

### D. Complexity Analysis

The main source of the computational cost of AIB-attention is its iterative calculation process. Its complexity can be expressed in the following manner:

$$O\left(\int_0^{\lceil log_2 L \rceil} d\left(L - 2^j\right) dj\right) = O\left(\frac{d}{ln2}\left(LlnL - L + 1\right)\right) \quad (55)$$

$$\approx O\left(Ldlog_2 L\right). \quad (56)$$

Hence, AIB-attention has lower computational complexity than self-attention when the sequence length $L$ exceeds $4$. Its complexity increases as $O\left(LDlog_2 L\right)$, which is significantly lower than the quadratic complexity of self-attention.

### E. Initialization of Parameters

We initialize all parameters in the neural network to ensure gradient stability during training. All biases are set to zero and all weights are sampled from normal distributions, with zero mean and specific standard deviations. For position encoding, word embedding parameters, $W^A$, $W^V$, and $W_1$, we use $\frac{1}{\sqrt{d}}$. For $w$, $W^O$, and $W_2$, we use $1$, $\sqrt{\frac{1-\frac{2}{d}}{2Nd}}$, and $1.7047\sqrt{\frac{1-\frac{2}{d}}{2Nd_{ff}}}$, respectively.

## IV. EXPERIMENTS

In the previous section, we described the encoder and decoder architectures based on AIB-attention. In this section, we assess the performance of these architectures on several long-sequence tasks using the Long Range Arena (LRA) [36], wikitext-103 [59], text8 [39], and enwik8 [40] data sets. We implemented our neural network models in PyTorch and used Adam [60] as the optimizer. In addition, we trained and tested our models on a single NVIDIA GeForce RTX 3090 GPU.

### A. Long-range Modeling on LRA

We evaluated the encoder performance of our model on the LRA data set, a systematic benchmark for efficient sequence processing models. The LRA data set measures model quality under long-sequence scenarios across six tasks and various data types, such as text, natural images, synthetic images, and mathematical expressions. The sequence lengths in the LRA data set vary from 1K to 16K, requiring different reasoning abilities, such as similarity, structural, and visual-spatial. The LRA data set can help comprehensively examine the generalization ability, computational efficiency, memory usage, and other aspects of sequence processing models.

This model's attention calculation transmits information in a directional manner. Therefore, the classifier built by this model cannot use the common [CLS] method [61]. Instead, it uses the average of the output layer of the neural network as the output.

TABLE III
EXPERIMENTAL RESULTS OF LONG-RANGE ARENA BENCHMARK.

| Model | ListOps | Text | Retrieval | Image | Pathfinder | Path-X | Avg |
|---|---|---|---|---|---|---|---|
| Transformer | 36.37 | 64.27 | 57.46 | 42.44 | 71.40 | FAIL | 54.39 |
| Local Attention [23] | 15.82 | 52.98 | 53.39 | 41.46 | 66.63 | FAIL | 46.06 |
| Sparse Trans [24] | 17.07 | 63.58 | 59.59 | 44.24 | 71.71 | FAIL | 51.24 |
| Longformer [25] | 35.63 | 62.85 | 56.89 | 42.22 | 69.71 | FAIL | 53.46 |
| Linformer [26] | 35.70 | 53.94 | 52.27 | 38.56 | 76.34 | FAIL | 51.36 |
| Reformer [27] | 37.27 | 56.10 | 53.40 | 38.07 | 68.50 | FAIL | 50.67 |
| Sinkhorn Trans [28] | 33.67 | 61.20 | 53.83 | 41.23 | 67.45 | FAIL | 51.39 |
| Synthesizer [29] | 36.99 | 61.68 | 54.67 | 41.61 | 69.45 | FAIL | 52.88 |
| BigBird [30] | 36.05 | 64.02 | 59.29 | 40.83 | 74.87 | FAIL | 55.01 |
| Linear Trans [31] | 16.13 | 65.90 | 53.09 | 42.34 | 75.30 | FAIL | 50.55 |
| Performer [32] | 18.01 | 65.40 | 53.82 | 42.77 | 77.05 | FAIL | 51.41 |
| Our Work | 39.48 | 70.54 | 63.92 | 50.30 | 84.76 | 63.99 | 62.17 |

TABLE IV
COMPARISON OF SPEED AND MEMORY USAGE AMONG VARIOUS
MODELS WITH DIFFERENT INPUT LENGTHS.

| Model | Steps per second | | | | Peak Memory Usage (GB) | | | |
|---|---|---|---|---|---|---|---|---|
| | 1K | 2K | 3K | 4K | 1K | 2K | 3K | 4K |
| Transformer | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| Local Attention | 1.1 | 1.7 | 3.2 | 5.3 | 0.49 | 0.29 | 0.19 | 0.14 |
| Linformer | 1.2 | 1.9 | 3.7 | 5.5 | 0.44 | 0.21 | 0.18 | 0.10 |
| Reformer | 0.5 | 0.4 | 0.7 | 0.8 | 0.56 | 0.37 | 0.28 | 0.24 |
| Sinkhorn Trans | 1.1 | 1.6 | 2.9 | 3.8 | 0.55 | 0.31 | 0.21 | 0.16 |
| Synthesizer | 1.1 | 1.2 | 2.9 | 1.4 | 0.76 | 0.75 | 0.74 | 0.74 |
| BigBird | 0.9 | 0.8 | 1.2 | 1.1 | 0.91 | 0.56 | 0.40 | 0.30 |
| Linear Trans | 1.1 | 1.9 | 3.7 | 5.6 | 0.44 | 0.22 | 0.15 | 0.11 |
| Performer | 1.2 | 1.9 | 3.8 | 5.7 | 0.44 | 0.22 | 0.15 | 0.11 |
| Our Work | 1.0 | 1.6 | 2.2 | 2.7 | 1.24 | 0.87 | 0.63 | 0.50 |

TABLE V
TEST AND TRAIN ACCURACY OF DIFFERENT MODELS ON THE IMAGE
CLASSIFICATION TASK.

| Model | Test Accuracy | Train Accuracy |
|---|---|---|
| Transformer | 42.44 | 69.45 |
| Local Attention | 41.46 | 63.19 |
| Linformer | 38.56 | 97.23 |
| Reformer | 38.07 | 68.45 |
| Sinkhorn Trans | 41.23 | 69.21 |
| Synthesizer | 41.61 | 97.31 |
| BigBird | 40.83 | 71.49 |
| Linear Trans | 42.34 | 65.61 |
| Performer | 42.77 | 73.90 |
| Our Work | 50.30 | 52.88 |

To ensure a fair comparison of different models, we follow the same model configuration as [36], including the structure and number of iterations. Table II summarizes the neural network hyperparameters used in each task. Table III reports the accuracy scores of various tasks in the experiment. Despite the limited number of iterations that prevent the model from converging in these tasks, it still surpasses all other Transformer-based models and baseline methods in terms of average accuracy.

This model outperforms some Transformer-based sequence processing models in terms of computation speed and memory usage. Table IV presents the comparison of speed and memory usage among various models with different input lengths (1K, 2K, 3K, and 4K). Following the experimental setting of [36], we evaluate all models on a byte-level classification task with the same batch size.

Fig. 2 depicts the trade-off among memory usage, computation speed, and performance for different models. Most models that reduce the complexity of the Transformer model sacrifice its performance. However, our model outperforms them with lower time and space complexity.

A common drawback of most Transformer models is their weak generalization ability. In contrast, our model establishes an attention mechanism based on adaptive inductive bias, which enables the model to learn position-dependent inductive relations. This endows it with superior robustness. Table V presents the training and testing accuracy of various models
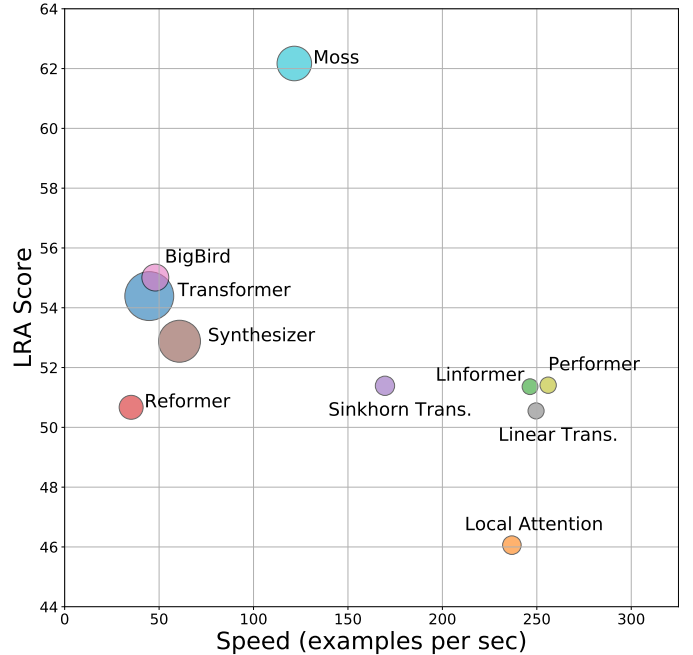


Fig. 2. Performance ($y$ axis), speed ($x$ axis), and memory footprint (size of the circles) of different models.

on the image classification benchmark. It is evident that all Transformer-like models suffer from severe overfitting, with a large gap between the two scores, despite using strong reg-

TABLE VI
STRUCTURAL PARAMETERS AND EXPERIMENTAL RESULTS OF MOSS AND GPT-2 IN WIKITEXT-103.

| Model | $N$ | $d_{model}$ | $d_{ff}$ | Parameters | PPL |
|---|---|---|---|---|---|
| GPT-2 small | 12 | 768 | 3072 | 114M | 37.50 |
| GPT-2 medium | 24 | 1024 | 4096 | 345M | 26.37 |
| GPT-2 large | 36 | 1280 | 5120 | 752M | 22.05 |
| Moss-small | 12 | 768 | 3072 | 117M | 21.27 |
| Moss-medium | 24 | 1024 | 4096 | 328M | 19.42 |
| Moss-large | 36 | 1280 | 5120 | 712M | 17.88 |

ularization mechanisms. Without applying any regularization algorithm, our model achieves lower training accuracy than all Transformer models, but higher testing accuracy than all Transformer models and is very close to the training accuracy level. This demonstrates that our model is much more robust than Transformer models.

### B. Natural Language Modelling

We tested the decoder of our model on various natural language data sets, including WikiText-103, Text8, and enwik8, without any additional training data. The results indicate that our model significantly outperforms the GPT-2 [4] model with the same structural parameters. Moreover, our model outperforms the current SOTA in character-level language modeling task of enwik8.

*1) Language Modelling on WikiText-103:* In this section, we compare Moss with GPT-2's small and medium models. GPT-2 is a Transformer-based language model that has been trained unsupervised on a large corpus of English text and can generate coherent and fluent text. We use WikiText-103— a large-scale English text data set with 1.03 million tokens extracted from Wikipedia—as our evaluation benchmark. We follow the same text encoding method and training settings as in the GPT-2 paper. The input sequence consists of 1024 consecutive tokens randomly cropped from the document, and the batch size is 1. In addition, we use the Adam optimizer and adjust the learning rate according to the GPT-2 paper.

The structural parameters of Moss in this experiment are consistent with GPT-2's small medium and large models. The structural parameters and experiment results are presented in Table VI. Without using any additional training data, Moss-small achieves a perplexity (PPL) of 21.27 on WikiText-103, Moss-medium achieves a PPL of 19.42, and Moss-large achieves a PPL of 17.88, both of which are much better than GPT-2 small model's 37.50, GPT-2 medium model's 26.37 and GPT-2 large model's 22.05. This demonstrates that Moss has a stronger ability and generalization performance on English text-generation tasks.

*2) Character-Level Language Modelling:* Character-level language modeling greatly increases the length and complexity of the token sequence, which challenges the ability of language models to capture long-distance dependencies. Conventional sequence models struggle with this task. We experimented on Text8, and enwik8 data sets without using any extra training data. In the Text8 and enwik8 data sets, we set the last 10M characters as the test set. In the language modeling task of enwik8, we also attempted to add the last 900M characters

TABLE VII
BPC AND PARAMETERS OF MOSS-LARGE AND GPT-2.

| Dataset | GPT-2 | | Moss-large | |
|---|---|---|---|---|
| | BPC | Parameters | BPC | Parameters |
| Text8 | 1.02 (large) [4] | 752M | 0.99 | 712M |
| enwik8 | 0.97 (large) [4] | 752M | 0.915 | 712M |
| enwik8 | 0.93 (SOTA) [4] | 1542M | 0.915 | 712M |

of enwik9 [62] as an additional training set. In addition, we randomly cropped 1024 consecutive characters as the neural network input. Table VII compares Moss-large with GPT-2 large and SOTA (GPT-2 xl) in terms of bits per character (BPC). Moss outperforms GPT-2 in all tasks.

## V. CONCLUSION AND PROSPECTS

In this paper, we introduced Moss, a novel neural network that exhibits superior fitting ability and robustness. Moss models long-distance dependencies through an adaptive inductive bias attention mechanism that has near-linear complexity and enables parallel computation. In both encoder and decoder experiments, the model surpasses all Transformer-based models in terms of accuracy with the same number of structural parameters. The model also excels in language modeling tasks with long sequences and achieves SOTA results, even with a small parameter size in certain tasks.

As future work, we plan to apply this model to image recognition, translation, and ultra-long character-level text generation, among other domains. We also aim to explore more applications of Moss models in multimodal learning and cross-lingual transfer learning. Due to its remarkable fitting ability and robustness, we believe this model has great potential in various research fields. However, we also acknowledge that our model has a few limitations and challenges, such as the scalability of very large data sets and vocabularies, the interpretability of the attention mechanism, and the ethical issues of generating realistic texts. Moreover, we note that our current experiments are conducted on medium-scale networks due to the limited computational resources. We intend to investigate the performance of larger-scale networks in the future. We hope our work can inspire more research on these aspects.

by the AI "Moss" in the movie The Wandering Earth, a science fiction film that depicts the human struggle to save the Earth from a dying sun. We hope that our work can contribute to the advancement of AI and its applications for the benefit of humanity.

## REFERENCES

[1] S. Bai, J. Z. Kolter, and V. Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," *arXiv preprint arXiv:1803.01271*, 2018.

[2] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[3] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever *et al.*, "Improving language understanding by generative pre-training," 2018.

[4] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.

[5] N. Zhang, "Learning adversarial transformer for symbolic music generation," *IEEE transactions on neural networks and learning systems*, 2020.

[6] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.

[7] Y. Luan and S. Lin, "Research on text classification based on cnn and lstm," in *2019 IEEE international conference on artificial intelligence and computer applications (ICAICA)*. IEEE, 2019, pp. 352–355.

[8] S. N. Tran, A. d. Garcez, T. Weyde, J. Yin, Q. Zhang, and M. Karunanithi, "Sequence classification restricted boltzmann machines with gated units," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 11, pp. 4806–4815, 2020.

[9] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," *Advances in neural information processing systems*, vol. 27, 2014.

[10] M.-T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," *arXiv preprint arXiv:1508.04025*, 2015.

[11] B. Zhang, D. Xiong, J. Xie, and J. Su, "Neural machine translation with gru-gated attention model," *IEEE transactions on neural networks and learning systems*, vol. 31, no. 11, pp. 4688–4698, 2020.

[12] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[13] H. Purwins, B. Li, T. Virtanen, J. Schlüter, S.-Y. Chang, and T. Sainath, "Deep learning for audio signal processing," *IEEE Journal of Selected Topics in Signal Processing*, vol. 13, no. 2, pp. 206–219, 2019.

[14] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.

[15] Y. Liu, T. Han, S. Ma, J. Zhang, Y. Yang, J. Tian, H. He, A. Li, M. He, Z. Liu *et al.*, "Summary of chatgpt/gpt-4 research and perspective towards the future of large language models," *arXiv preprint arXiv:2304.01852*, 2023.

[16] R. M. Schmidt, "Recurrent neural networks (rnns): A gentle introduction and overview," *arXiv preprint arXiv:1912.05911*, 2019.

[17] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[18] Y. Tay, M. Dehghani, D. Bahri, and D. Metzler, "Efficient transformers: A survey," *ACM Computing Surveys*, vol. 55, no. 6, pp. 1–28, 2022.

[19] M. Hahn, "Theoretical limitations of self-attention in neural sequence models," *Transactions of the Association for Computational Linguistics*, vol. 8, pp. 156–171, 2020.

[20] T. Lin, Y. Wang, X. Liu, and X. Qiu, "A survey of transformers," *AI Open*, 2022.

[21] T. Dao, D. Fu, S. Ermon, A. Rudra, and C. Ré, "Flashattention: Fast and memory-efficient exact attention with io-awareness," *Advances in Neural Information Processing Systems*, vol. 35, pp. 16344–16359, 2022.

[22] C. Wu, F. Wu, T. Qi, Y. Huang, and X. Xie, "Fastformer: Additive attention can be all you need," *arXiv preprint arXiv:2108.09084*, 2021.

[23] M. Arar, A. Shamir, and A. H. Bermano, "Learned queries for efficient local attention," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10841–10852.

[24] R. Child, S. Gray, A. Radford, and I. Sutskever, "Generating long sequences with sparse transformers," *arXiv preprint arXiv:1904.10509*, 2019.

[25] I. Beltagy, M. E. Peters, and A. Cohan, "Longformer: The long-document transformer," *arXiv preprint arXiv:2004.05150*, 2020.

[26] S. Wang, B. Z. Li, M. Khabsa, H. Fang, and H. Ma, "Linformer: Self-attention with linear complexity," *arXiv preprint arXiv:2006.04768*, 2020.

[27] N. Kitaev, Ł. Kaiser, and A. Levskaya, "Reformer: The efficient transformer," *arXiv preprint arXiv:2001.04451*, 2020.

[28] Y. Tay, D. Bahri, L. Yang, D. Metzler, and D.-C. Juan, "Sparse sinkhorn attention," in *International Conference on Machine Learning*. PMLR, 2020, pp. 9438–9447.

[29] Y. Tay, D. Bahri, D. Metzler, D. Juan, Z. Zhao, and C. Zheng, "Synthesizer: Rethinking self-attention in transformer models. arxiv 2020," *arXiv preprint arXiv:2005.00743*, vol. 2, 2020.

[30] M. Zaheer, G. Guruganesh, K. A. Dubey, J. Ainslie, C. Alberti, S. Ontanon, P. Pham, A. Ravula, Q. Wang, L. Yang *et al.*, "Big bird: Transformers for longer sequences," *Advances in neural information processing systems*, vol. 33, pp. 17283–17297, 2020.

[31] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret, "Transformers are rnns: Fast autoregressive transformers with linear attention," in *International Conference on Machine Learning*. PMLR, 2020, pp. 5156–5165.

[32] K. Choromanski, V. Likhosherstov, D. Dohan, X. Song, A. Gane, T. Sarlos, P. Hawkins, J. Davis, D. Belanger, L. Colwell *et al.*, "Masked language modeling for proteins via linearly scalable long-context transformers," *arXiv preprint arXiv:2006.03555*, 2020.

[33] S. Hao, D.-H. Lee, and D. Zhao, "Sequence to sequence learning with attention mechanism for short-term passenger flow prediction in large-scale metro system," *Transportation Research Part C: Emerging Technologies*, vol. 107, pp. 287–300, 2019.

[34] D. Ju, S. Roller, S. Sukhbaatar, and J. E. Weston, "Staircase attention for recurrent processing of sequences," *Advances in Neural Information Processing Systems*, vol. 35, pp. 13203–13213, 2022.

[35] A. Hernández and J. M. Amigó, "Attention mechanisms and their applications to complex systems," *Entropy*, vol. 23, no. 3, p. 283, 2021.

[36] Y. Tay, M. Dehghani, S. Abnar, Y. Shen, D. Bahri, P. Pham, J. Rao, L. Yang, S. Ruder, and D. Metzler, "Long range arena: A benchmark for efficient transformers," *arXiv preprint arXiv:2011.04006*, 2020.

[37] B. Li, Y. Hu, X. Nie, C. Han, X. Jiang, T. Guo, and L. Liu, "Dropkey," *arXiv preprint arXiv:2208.02646*, 2022.

[38] T. Mikolov, M. Karafiát, L. Burget, J. Cernockỳ, and S. Khudanpur, "Recurrent neural network based language model." in *Interspeech*, vol. 2, no. 3. Makuhari, 2010, pp. 1045–1048.

[39] M. Mahoney, "Large text compression benchmark," 2011.

[40] "enwik8 ·datasets at hugging face," https://huggingface.co/datasets/enwik8, accessed on May 26, 2023.

[41] Y. Kim, C. Denton, L. Hoang, and A. M. Rush, "Structured attention networks," *arXiv preprint arXiv:1702.00887*, 2017.

[42] Z. Niu, G. Zhong, and H. Yu, "A review on the attention mechanism of deep learning," *Neurocomputing*, vol. 452, pp. 48–62, 2021.

[43] G. W. Lindsay, "Attention in psychology, neuroscience, and machine learning," *Frontiers in computational neuroscience*, vol. 14, p. 29, 2020.

[44] Q. Lai, S. Khan, Y. Nie, H. Sun, J. Shen, and L. Shao, "Understanding more about human and machine attention in deep neural networks," *IEEE Transactions on Multimedia*, vol. 23, pp. 2086–2099, 2020.

[45] S. Ghaffarian, J. Valente, M. Van Der Voort, and B. Tekinerdogan, "Effect of attention mechanism in deep learning-based remote sensing image processing: A systematic literature review," *Remote Sensing*, vol. 13, no. 15, p. 2965, 2021.

[46] A. Galassi, M. Lippi, and P. Torroni, "Attention in natural language processing," *IEEE transactions on neural networks and learning systems*, vol. 32, no. 10, pp. 4291–4308, 2020.

[47] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *International conference on machine learning*. PMLR, 2015, pp. 2048–2057.

[48] X.-Y. Zhang, C. Li, H. Shi, X. Zhu, P. Li, and J. Dong, "Adapnet: Adaptability decomposing encoder-decoder network for weakly supervised action recognition and localization," *IEEE transactions on neural networks and learning systems*, 2020.

[49] Z. Ji, Y. Zhao, Y. Pang, X. Li, and J. Han, "Deep attentive video summarization with distribution consistency learning," *IEEE transactions on neural networks and learning systems*, vol. 32, no. 4, pp. 1765–1775, 2020.

[50] S. Takase and S. Kiyono, "Lessons on parameter sharing across layers in transformers," *arXiv preprint arXiv:2104.06022*, 2021.

[51] L. Wu, J. Li, Y. Wang, Q. Meng, T. Qin, W. Chen, M. Zhang, T.-Y. Liu *et al.*, "R-drop: Regularized dropout for neural networks," *Advances in Neural Information Processing Systems*, vol. 34, pp. 10 890–10 905, 2021.

[52] Z. Zhang, J. Yang, and H. Zhao, "Retrospective reader for machine reading comprehension," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 16, 2021, pp. 14 506–14 514.

[53] M. Lippi, M. A. Montemurro, M. Degli Esposti, and G. Cristadoro, "Natural language statistical features of lstm-generated texts," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 11, pp. 3326–3337, 2019.

[54] J. Baxter, "A model of inductive bias learning," *Journal of artificial intelligence research*, vol. 12, pp. 149–198, 2000.

[55] P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner *et al.*, "Relational inductive biases, deep learning, and graph networks," *arXiv preprint arXiv:1806.01261*, 2018.

[56] A. Goyal and Y. Bengio, "Inductive biases for deep learning of higher-level cognition," *Proceedings of the Royal Society A*, vol. 478, no. 2266, p. 20210068, 2022.

[57] N. Rahaman, A. Baratin, D. Arpit, F. Draxler, M. Lin, F. Hamprecht, Y. Bengio, and A. Courville, "On the spectral bias of neural networks," in *International Conference on Machine Learning*. PMLR, 2019, pp. 5301–5310.

[58] A. L. Maas, A. Y. Hannun, A. Y. Ng *et al.*, "Rectifier nonlinearities improve neural network acoustic models," in *Proc. icml*, vol. 30, no. 1. Atlanta, Georgia, USA, 2013, p. 3.

[59] S. Merity, C. Xiong, J. Bradbury, and R. Socher, "Pointer sentinel mixture models," *arXiv preprint arXiv:1609.07843*, 2016.

[60] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[61] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[62] Marcus Hutter, "Hutter prize for lossless compression of human knowledge," 2006, http://prize.hutter1.net/, Last accessed on 2023-05-26.

[63] G. Melis, T. Kočiský, and P. Blunsom, "Mogrifier lstm," *arXiv preprint arXiv:1909.01792*, 2019.

**Lingkai Hu** received a bachelor's degree from the School of Mechanical and Electrical Engineering, Guangzhou University, China, in 2021. He is currently studying for a master's degree in mechanical engineering at Guangzhou University, Guangzhou, China. His current research interests include deep learning (DL), sequence analysis, object detection, attention mechanism, and multilevel feature fusion.



**Wenkai Huang** received his B.S. degree and M.S. from Guangdong University of Technology. He received his Ph.D. from Guangzhou University. In 2007, he joined the School of Mechanical and Electrical Engineering at Guangzhou University, where he is currently an Associate Professor. His research interests are AI, robot vision, medical image processing, and soft robotics.



**Feng zhan** is currently pursuing an undergraduate degree in Electrical Engineering at Guangzhou University. His research focuses on fault diagnosis using deep learning techniques, with a particular emphasis on aspects such as few-shot learning, diffusion models, and self-attention.



**Weiming Gan** is currently studying for a master's degree in machinery engineering at Guangzhou University, Guangzhou, China. His current research interests include industrial data mining recognition and analysis, deep learning (DL), and image processing technology.