# Nd-BiMamba2: A Unified Bidirectional Architecture for Multi-Dimensional Data Processing

Hao Liu, *College of Information Science and Technology, Qingdao University of Science and Technology*

*Abstract*—**Deep learning models often require specially designed architectures to process data of different dimensions, such as 1D time series, 2D images, and 3D volumetric data. Existing bidirectional models mainly focus on sequential data, making it difficult to scale effectively to higher dimensions. To address this issue, we propose a novel multi-dimensional bidirectional neural network architecture, named Nd-BiMamba2, which efficiently handles 1D, 2D, and 3D data. Nd-BiMamba2 is based on the Mamba2 module and introduces innovative bidirectional processing mechanisms and adaptive padding strategies to capture bidirectional information in multi-dimensional data while maintaining computational efficiency. Unlike existing methods that require designing specific architectures for different dimensional data, Nd-BiMamba2 adopts a unified architecture with a modular design, simplifying development and maintenance costs. To verify the portability and flexibility of Nd-BiMamba2, we successfully exported it to ONNX and TorchScript and tested it on different hardware platforms (e.g., CPU, GPU, and mobile devices). Experimental results show that Nd-BiMamba2 runs efficiently on multiple platforms, demonstrating its potential in practical applications. The code is open-source: https://github.com/Human9000/nd-Mamba2-torch.**

*Index Terms*—**mamba2, nd-mamba2, bimamba2, attention, multi-dimensional learning, deep learning, model deployment, ONNX, TorchScript, cross-platform**

## I. INTRODUCTION

Deep learning has made significant progress in many fields, but data of different dimensions (e.g., 1D time series, 2D images, and 3D volumetric data) often require specially designed model architectures. For instance, convolutional neural networks (CNNs) [1] excel at processing image data, recurrent neural networks (RNNs) [2] are suitable for sequential data, while 3D CNNs are used for volumetric data. This domain-specific model design paradigm leads to increased development and maintenance costs and limits the generalization ability of models.

Although bidirectional models, such as bidirectional RNNs (BiRNNs) [3], have been successful in sequential data modeling, they struggle to scale effectively to higher-dimensional data and face challenges in cross-platform deployment. The sequential processing nature of BiRNNs limits their parallelization capabilities, making them inefficient for long sequences and high-dimensional data, and they are prone to gradient vanishing issues. Moreover, the recurrent structure of RNNs makes it difficult to convert them into formats like ONNX or TorchScript for cross-platform deployment. On the other hand, self-attention mechanisms like Transformers [4] can capture long-range dependencies, but their computational complexity becomes prohibitive and memory consumption increases when processing high-dimensional data, complicating deployment.

While the existing Mamba [5] model strikes a balance between efficiency and performance, most are limited to unidirectional processing or data of specific dimensions. To overcome these limitations, this paper proposes a novel multi-dimensional bidirectional neural network architecture, Nd-BiMamba2. The core innovations of Nd-BiMamba2 include: 1) extending the Mamba2 module to support efficient bidirectional processing that can be applied effectively to 1D, 2D, and 3D data; 2) introducing an adaptive padding strategy that adjusts padding size based on input data dimensions, improving computational efficiency and reducing memory consumption.

The main contributions of this paper are as follows:

- We propose Nd-BiMamba2, a unified bidirectional network architecture that can efficiently process multi-dimensional data.
- We design an innovative bidirectional processing mechanism that effectively captures bidirectional information in high-dimensional data.
- We introduce an adaptive padding strategy to improve computational efficiency and reduce memory consumption.
- We validate the portability and deployment capability of Nd-BiMamba2 across different hardware platforms.

The following sections will provide detailed descriptions of the network structure and implementation details of Nd-BiMamba2, experimental results, its performance on multi-dimensional tasks, and discuss the model's advantages and potential applications.

## II. RELATED WORK

Multi-dimensional data modeling is a key research direction in deep learning, encompassing various scenarios such as 1D time series, 2D images, and 3D volumetric data. To efficiently model multi-dimensional data, researchers have proposed various methods, including convolutional neural networks (CNNs), recurrent neural networks (RNNs), self-attention mechanisms, and recently emerging modular architectures such as Mamba. However, these methods have limitations to varying degrees and struggle to balance the efficiency and generalizability required for multi-dimensional feature modeling.

### A. Convolutional Neural Networks (CNN)

CNNs, as classical deep learning methods, have achieved outstanding performance in image processing tasks. Typical models such as LeNet [6], ResNet [7], and U-Net [8] extract local features through multiple layers of convolutions and

progressively expand the receptive field. However, CNNs face the following limitations in multi-dimensional data modeling:

- Inadequate long-range dependency modeling: CNNs struggle to capture global context information when processing long sequences or high-resolution images. - High computational cost for high-dimensional extension: 3D CNNs are effective for spatial feature extraction but significantly increase the parameter scale and computational complexity, limiting their practical applications.

### B. Recurrent Neural Networks (RNN)

RNNs and their variants (such as LSTM [9] and GRU [10]) perform excellently in sequence modeling, especially in capturing long-term and short-term dependencies in time series. For example, bidirectional LSTMs (BiLSTMs) [11] enhance context modeling in natural language processing tasks by fusing bidirectional information.WaveNet [12] introduces a novel deep neural network architecture based on dilated causal convolutions, capable of directly generating high-quality raw audio waveforms and effectively capturing long-range dependencies in audio signals. However, RNNs have the following limitations:

1) Difficulty in parallelization: The sequential processing nature of RNNs makes them inefficient when handling long sequences. 2) Challenges in scaling to high-dimensional data: The recurrent structure does not adapt well to 2D images or 3D volumetric data, leading to increased memory consumption and computational complexity. 3) Training stability issues: RNNs still face gradient vanishing and gradient explosion problems, impacting model performance.

### C. Self-Attention Mechanisms (SA)

Self-attention mechanisms, with their global modeling ability, have been widely applied to natural language processing and computer vision tasks. The Transformer [4] is a representative model, and its extensions such as BERT [13] and ViT [14] have made significant progress in various fields. However, in multi-dimensional data modeling, self-attention mechanisms still face the following challenges:

1) High computational complexity: Although numerous Swin-based attention methods [15], [16] have been proposed to reduce computational complexity in 2D, the quadratic complexity of attention mechanisms leads to a significant increase in memory and computational resource requirements when dealing with high-dimensional data. 2) Poor adaptability to high-dimensional scenarios: While low-rank decomposition methods (such as Linformer [17],Rethinking [18]) reduce complexity, they still do not fully solve the memory bottleneck in high-dimensional data processing.

### D. Mamba Modules

The Mamba module is a lightweight architecture that combines the advantages of convolution and attention mechanisms, which has recently gained prominence in multi-dimensional data modeling. For example, the latest Mamba2 [19] and vssd [20] etc [21], [22]modules significantly improve image classification performance by combining local feature extraction

with global information modeling. However, existing Mamba modules primarily focus on unidirectional feature modeling and have the following limitations:

1) Lack of bidirectional feature modeling: The inability to effectively capture bidirectional information in multi-dimensional data limits its generalization capability.

2) Insufficient adaptation to multi-dimensional data: Current designs mainly target 1D or 2D image data individually, making it challenging to efficiently extend to 3D scenarios.

### E. Summary and Limitations

Existing methods each have their advantages, but still face shortcomings in efficient and multi-dimensional feature modeling:

1) CNNs excel at local feature extraction but struggle to capture global context.

2) RNNs are strong in modeling sequential data but suffer from low computational efficiency and poor scalability.

3) Self-attention mechanisms offer global modeling capabilities but come with high computational complexity.

4) The Mamba module, while excelling in lightweight design, lacks a unified modeling capability for multi-dimensional data.

### F. Innovations of Nd-BiMamba2

To address the above issues, we propose a unified bidirectional modeling architecture, Nd-BiMamba2. By extending the Mamba2 module, it supports efficient modeling of 1D, 2D, and 3D data. The bidirectional processing mechanism fully explores directional information in multi-dimensional data. Dynamic padding adjustment based on input data dimensions improves computational efficiency and reduces memory consumption. We validated the model's efficiency on CPU, GPU, and mobile devices, enhancing its practical application potential.

In conclusion, Nd-BiMamba2 provides a general and efficient solution, opening new directions for multi-dimensional data modeling and cross-platform deployment.

## III. ALGORITHM DESIGN

## IV. ALGORITHM DESIGN OPTIMIZATION

### A. Design Objectives and Challenges

To handle data of different dimensions (1D, 2D, 3D) and optimize model performance, the design objectives of Nd-BiMamba2 include the following:

- **Generality:** The algorithm needs to provide a unified processing framework to accommodate multi-dimensional data.
- **Efficiency:** To reduce computational redundancy on high-dimensional data, convolution operations need to be designed for adaptation.
- **Boundary Handling:** To avoid boundary effects in multi-dimensional scenarios, tailored padding strategies must be designed.

## B. Core Algorithm Design

*1) Input Representation:* To simplify the processing of data with different dimensions, the input tensor is uniformly represented as:

$$X \in \mathbb{R}^{B \times C \times D_1 \times D_2 \times D_3} \quad (1)$$

where $B$ is the batch size, $C$ is the number of channels, and $D_1, D_2, D_3$ represent the sizes of the three dimensions. For 1D and 2D data, this is maintained consistently by setting $D_3 = 1$ or $D_2 = D_3 = 1$.

This unified representation reduces the complexity of handling logical branches between different dimensional data, allowing subsequent convolution and activation operations to reuse the same logic.

*2) Core Convolution Calculation Formula:*

*a) Design Philosophy::* To capture local features, dimension-adaptive convolution operations are employed. The core calculation formulas are as follows:

$$\mathrm{F}(X) = \sigma(W_f * X + b_f) \quad (2)$$
$$\mathrm{B}(X) = \sigma(W_b * X + b_b) \quad (3)$$

where $W_f, W_b$ are the convolution kernels for the forward and backward paths, $b_f, b_b$ are the biases, $*$ denotes dimension-adaptive convolution, and $\sigma$ is the activation function.

To enhance the model's ability to handle directional information, separate forward and backward paths are designed. Additionally, an activation function $\sigma$ is included to improve the model's nonlinear modeling capabilities.

*b) Dimensional Differences and Optimizations::*

1) **1D Data:** For processing sequential data, the convolution kernel is designed with shape $(k, 1, 1)$, sliding only along the $D_1$ direction:

$$Y[i] = \sum_{j=0}^{k-1} W[i,j] \cdot X[i \cdot s + j] + b[i] \quad (4)$$

where $k$ is the kernel size, and $s$ is the stride. To reduce computational redundancy in other dimensions, the convolution operation slides only along the $D_1$ direction, improving computational efficiency.

2) **2D Data:** For processing image data, the convolution kernel is designed with shape $(k_1, k_2, 1)$, sliding along both $D_1$ and $D_2$:

$$Y[i,j] = \sum_{m=0}^{k_1-1} \sum_{n=0}^{k_2-1} W[i,j,m,n] \cdot$$
$$X[i \cdot s_1 + m, j \cdot s_2 + n] + b[i,j]. \quad (5)$$

To effectively capture local pattern information, the convolution operation slides simultaneously along $D_1$ and $D_2$, which is suitable for extracting image features.

3) **3D Data:** For processing high-dimensional spatial data, the convolution kernel is designed with shape $(k_1, k_2, k_3)$, sliding along $D_1, D_2, D_3$ simultaneously:

$$Y[i,j,k] = \sum_{m=0}^{k_1-1} \sum_{n=0}^{k_2-1} \sum_{l=0}^{k_3-1} W[i,j,k,m,n,l] \cdot$$
$$X[i \cdot s_1 + m, j \cdot s_2 + n, k \cdot s_3 + l] + b[i,j,k]. \quad (6)$$

To capture the complex features in 3D space, convolution operations are evenly distributed across the three dimensions, improving feature extraction capabilities.

*3) Padding Strategy:*

*a) Formula Definition::* To handle boundary effects, the padding size $p_i$ in the $i$-th dimension is calculated as:

$$p_i = \max(0, \lceil \frac{D_i - 1 \cdot s_i + k_i - 1}{2} \rceil) \quad (7)$$

where $k_i$ and $s_i$ are the kernel size and stride for the $i$-th dimension, respectively.

*b) Dimensional Differences and Advantages::*

1) **1D Data:** To preserve the original data characteristics, padding is minimized only along the $D_1$ direction.
2) **2D Data:** To enhance the effectiveness of the boundary regions, a mirroring padding strategy is applied along both $D_1$ and $D_2$.
3) **3D Data:** To balance boundary handling with computational complexity in high-dimensional scenarios, padding is uniformly distributed across $D_1, D_2, D_3$.

*4) Activation Function Selection:* To improve the model's nonlinear expression capabilities, Nd-BiMamba2 uses the GELU (Gaussian Error Linear Unit) activation function, defined as:

$$\sigma(x) = x \cdot \Phi(x) \quad (8)$$

where $\Phi(x)$ is the cumulative distribution function of the standard normal distribution.

The GELU activation function was selected to more smoothly handle the distribution of input values, especially exhibiting stronger feature extraction abilities in high-dimensional data.

Overall, Nd-BiMamba2 retains the advantages of Bi-Mamba2 when processing sequential and image data, and by incorporating support for three-dimensional data, along with more refined partitioning and feature fusion techniques, it extends the application scope. This improvement enables Nd-BiMamba2 to provide more efficient and accurate modeling capabilities when dealing with more complex input data.

Nd-BiMamba2's modules and functional layers are shown in Table I:

## C. Comparative Analysis

To highlight the advantages of Nd-BiMamba2, the following Table II summarizes its comparison with other models:

TABLE I
ALGORITHM MODULES, LAYERS, AND FUNCTIONAL DESCRIPTIONS

| Module | Contained Layers | Functional Description |
|---|---|---|
| Data Preprocessing | Input Padding | Pads the input data to meet processing requirements: 1D padding to a multiple of 4, 2D padding to a multiple of 8, 3D padding to a multiple of 4. |
| | Dimension Adjustment | Rearranges the data according to its dimensions, flattening 2D or 3D data into a 2D matrix to conform to the network structure. |
| | Channel Mapping | Uses a linear layer $\text{FC}_{\text{in}}$ to map the input channel size $c$ to the target model's dimension $d_{\text{model}}$. |
| Bi-Directional Modeling | Forward Feature Extraction | Extracts features through the forward Mamba2 network to obtain the forward feature representation $H_{\text{forward}}$. |
| | Backward Feature Extraction | Reverses the input data and inputs it into the backward Mamba2 network to extract the backward feature representation $H_{\text{backward}}$, then restores the original order. |
| | Feature Fusion | Fuses the forward and backward feature representations using an addition operation to obtain the final feature representation: $H_{\text{fused}} = H_{\text{forward}} + H_{\text{backward}}$. |
| Output Generation | Linear Transformation | Uses a linear layer $\text{FC}_{\text{out}}$ to map the fused feature representation back to the target channel size $c'$. |
| | Padding Removal | Removes the additional data added during padding to restore the original shape of the input data. |

TABLE II
MODEL COMPARISON ANALYSIS

| Model | Applicable Data Dimensions | Cross-Platform Computational Efficiency | Modular Design | Deployment |
|---|---|---|---|---|
| BiLSTM | 1D | Medium | No | Difficult |
| Transformer | 1D/2D/3D | Low | No | Difficult |
| Mamba2 | 1D/2D | High | Yes | Fairly Easy |
| **Nd-BiMamba2** | 1D/2D/3D | **High** | **Yes** | **Easy** |

### D. Summary

By optimizing strategies for padding, dimension rearrangement, channel adjustment, and feature fusion across different dimensions, the model can efficiently extract features from 1D, 2D, and 3D data while maintaining consistency in the output dimension with the input data. These steps are clearly described through mathematical symbols to ensure the correctness and efficiency of multi-dimensional optimization.

### E. Model Export and Deployment

To enhance model portability and deployment capabilities, Nd-BiMamba2 supports multiple export formats:

- **ONNX Export**: Supports converting the model to ONNX format for running on various hardware platforms.
- **TorchScript Export**: Supports converting the model to TorchScript format to ensure efficient inference in production environments.

Through this modular design and multi-dimensional optimization, Nd-BiMamba2 achieves efficient and unified modeling for 1D, 2D, and 3D data, providing powerful support for multi-modal data processing.

## V. EXPERIMENTS

### A. Experimental Setup

All experiments were conducted on the following hardware platform:

- **Processor (CPU)**: Intel Core i9-11900K, 8 cores, 16 threads, 3.5 GHz base frequency. The high clock speed and multi-core design of this processor allow it to efficiently handle parallel computing tasks, particularly for processing large amounts of data and task scheduling, significantly enhancing overall computational performance.
- **Graphics Processing Unit (GPU)**: NVIDIA RTX 4090D, 24GB GDDR6X VRAM. As one of the latest high-performance GPUs, the RTX 4090D provides powerful parallel computing capabilities for deep learning model training and inference, especially for large-scale datasets and complex models. The 24GB of VRAM ensures the processing of large models and high-resolution data, effectively mitigating memory bottlenecks.
- **Memory (RAM)**: 64GB DDR4 3200 MHz. The ample memory capacity ensures efficient data reading and caching during model training, preventing computational bottlenecks due to memory limitations. This is especially important when handling large-scale data, maintaining high data throughput.
- **Storage**: 1TB NVMe SSD (used for data storage and intermediate result caching). The high-speed SSD improves data read/write speed, significantly reducing I/O latency, especially when training involves large amounts of data input and output, ensuring efficient operation during the training process.

### B. Feature Representation Ability of Nd-BiMamba2

The bidirectional modeling module of the Nd-BiMamba2 model enhances its feature perception ability by incorporating both forward and backward information flows. In traditional unidirectional modeling, the model can only rely on information from one direction of the input sequence for inference. In contrast, bidirectional modeling considers both forward and backward information flows, allowing for the capture of more

comprehensive features. The advantages of bidirectional modeling are particularly evident in various data dimensions (1D, 2D, and 3D), especially in capturing long-range dependencies and local features.

Through comparative experiments across different data dimensions (1D, 2D, and 3D), we have validated the improvement in feature representation by bidirectional modeling, demonstrating that this approach is more efficient than traditional unidirectional modeling when dealing with complex data. The experimental model configuration parameters were set as follows: $c_{in} = 64, c_{out} = 64, d_{\text{model}} = 128$, ensuring the ability to handle high-dimensional data and perform sufficient feature extraction.

TABLE III
PERFORMANCE COMPARISON BETWEEN ND-BIMAMBA2 AND TRADITIONAL UNIDIRECTIONAL MODELING (DIMENSIONS: 1D/2D/3D, FLOPS IN GMAC, TIME IN MILLISECONDS, PARAMETERS IN THOUSANDS)

| Bi. | Size | FLOPs (GMac) | Time (ms) | Params (k) |
|-----|------|--------------|-----------|------------|
| No | 1024 | 0.15 | 1.69 | |
| | $128 \times 128$ | 2.47 | 1.53 | 150.8 |
| | $32 \times 32 \times 32$ | 4.93 | 4.36 | |
| Yes | 1024 | 0.29 | 2.43 | |
| | $128 \times 128$ | 4.66 | 3.15 | 285.21 |
| | $32 \times 32 \times 32$ | 9.33 | 8.11 | |

Note: The number of parameters is independent of input size and is only affected by the use of Bi.

As shown in Table III, enabling bidirectional modeling leads to a significant increase in FLOPs (floating point operations) and computation time. Particularly with 3D data, the increase in FLOPs and computation time is more pronounced, though the growth in parameter count remains relatively small. This result indicates that while bidirectional modeling increases computational overhead, it captures more feature information and improves the model's expressive power.

### C. Flexibility and Adaptability from Modular Design

The modular design of Nd-BiMamba2 provides strong support for model flexibility and adaptability. Through experiments on 1D, 2D, and 3D data, the model can adaptively adjust padding strategies according to different data dimensions, ensuring computational efficiency and flexibility. With this design, Nd-BiMamba2 can dynamically adjust input size and padding strategy, achieving good computational efficiency and performance across different data dimensions.

To observe the model's performance with adaptive padding strategies across different data dimensions (1D, 2D, 3D), we conducted comparative experiments on multi-dimensional adaptive padding strategies. This verified that the strategy automatically adjusts padding methods for various input sizes to ensure dimensional consistency and efficient computation.

As seen in Table IV, the model demonstrates excellent flexibility under the adaptive padding strategy. Especially in 2D and 3D data processing, the adaptive padding proves particularly important. It effectively improves computational efficiency while maintaining high accuracy across different input sizes. This shows that Nd-BiMamba2 has strong adaptability in processing multi-dimensional data, adjusting itself according to the different characteristics of the data.

TABLE IV
PERFORMANCE OF MULTI-DIMENSIONAL ADAPTIVE PADDING STRATEGY ACROSS DIFFERENT FEATURE SIZES

| Dim. | Input | Auto-Padding | Mamba2 | Equal |
|------|-------|--------------|--------|-------|
| 1D | 1024 | 1024 | 1024 | TRUE |
| | 1029 | 1088 | 1088 | FALSE |
| | 1001 | 1024 | 1024 | FALSE |
| 2D | $128 \times 128$ | $128 \times 128$ | 16384 | TRUE |
| | $129 \times 127$ | $136 \times 128$ | 17408 | FALSE |
| | $113 \times 128$ | $120 \times 128$ | 15360 | FALSE |
| 3D | $32 \times 32 \times 32$ | $32 \times 32 \times 32$ | 32768 | TRUE |
| | $27 \times 33 \times 32$ | $28 \times 32 \times 36$ | 32256 | FALSE |
| | $37 \times 29 \times 31$ | $40 \times 32 \times 32$ | 40960 | FALSE |

### D. Conclusion

Through the analysis and experiments on the nd-BiMamba2 model, several significant advantages have been identified:

- **Bidirectional Modeling:** Bidirectional modeling significantly enhances the model's ability to perceive features, especially in capturing long-range dependencies and local characteristics.
- **Modular Design:** The modular design provides flexibility and adaptability, allowing the model to automatically adjust input sizes and padding strategies based on different data dimensions, ensuring computational efficiency and model flexibility.
- **Efficient Performance:** Despite the increased computational overhead from bidirectional modeling and adaptive padding, the model still performs excellently across multiple data dimensions, demonstrating its advantage in processing complex data.

Overall, nd-BiMamba2 exhibits strong performance in high-dimensional data processing, feature extraction accuracy, and computational efficiency, proving its effectiveness in complex data analysis, long-range dependency modeling, and large-scale data handling.

### APPENDIX

**Algorithm 1** Nd-BiMamba2 Algorithm

---

**Input:** $X \in \mathbb{R}^{c \times d_1 \times d_2 \times \cdots \times d_n}$
**Output:** $H_{\text{output}} \in \mathbb{R}^{c' \times d'_1 \times d'_2 \times \cdots \times d'_n}$
**Step 1: Data Preprocessing**

- $X_{\text{padded}} \leftarrow \text{Pad}(X)$       *Padding input data*
- $X_{\text{reshaped}} \leftarrow \text{Reshape}(X_{\text{padded}})$   *Adjusting dimensions*
- $X_{\text{mapped}} \leftarrow \text{FC}_{\text{in}}(X_{\text{reshaped}})$ *Mapping the channel count*

**Step 2: Bidirectional Modeling**

- $H_{\text{forward}} \leftarrow \text{Mamba2}_{\text{for}}(X_{\text{mapped}})$    *Forward feature extraction*
- $H_{\text{backward}} \leftarrow \text{Flip}(\text{Mamba2}_{\text{back}}(\text{Flip}(X_{\text{mapped}})))$   *Backward feature extraction*
- $H_{\text{fused}} \leftarrow H_{\text{forward}} + H_{\text{backward}}$     *Fusing features*

**Step 3: Output Generation**

- $H_{\text{fc\_out}} \leftarrow \text{FC}_{\text{out}}(H_{\text{fused}})$   *Restoring the channel count*
- $H_{\text{output}} \leftarrow \text{Trim}(H_{\text{fc\_out}})$   *Removing padded parts*

**Return:** $H_{\text{output}}$

---

## References

[1] W. Yao, J. Bai, W. Liao, Y. Chen, M. Liu, and Y. Xie, "From cnn to transformer: A review of medical image segmentation models," *Journal of Imaging Informatics in Medicine*, pp. 1–19, 2024.

[2] S. M. Al-Selwi, M. F. Hassan, S. J. Abdulkadir, A. Muneer, E. H. Sumiea, A. Alqushaibi, and M. G. Ragab, "Rnn-lstm: From applications to modeling techniques and beyond—systematic review," *Journal of King Saud University-Computer and Information Sciences*, p. 102068, 2024.

[3] Z. Cui, R. Ke, Z. Pu, and Y. Wang, "Stacked bidirectional and unidirectional lstm recurrent neural network for forecasting network-wide traffic state with missing values," *Transportation Research Part C: Emerging Technologies*, vol. 118, p. 102674, 2020.

[4] A. Vaswani, "Attention is all you need," *Advances in Neural Information Processing Systems*, 2017.

[5] A. Gu and T. Dao, "Mamba: Linear-time sequence modeling with selective state spaces," *arXiv preprint arXiv:2312.00752*, 2023.

[6] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[7] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[8] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*. Springer, 2015, pp. 234–241.

[9] A. Sherstinsky, "Fundamentals of recurrent neural network (rnn) and long short-term memory (lstm) network," *Physica D: Nonlinear Phenomena*, vol. 404, p. 132306, 2020.

[10] R. Dey and F. M. Salem, "Gate-variants of gated recurrent unit (gru) neural networks," in *2017 IEEE 60th international midwest symposium on circuits and systems (MWSCAS)*. IEEE, 2017, pp. 1597–1600.

[11] S. Zhang, D. Zheng, X. Hu, and M. Yang, "Bidirectional long short-term memory networks for relation classification," in *Proceedings of the 29th Pacific Asia conference on language, information and computation*, 2015, pp. 73–78.

[12] A. Van Den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, K. Kavukcuoglu *et al.*, "Wavenet: A generative model for raw audio," *arXiv preprint arXiv:1609.03499*, vol. 12, 2016.

[13] J. Devlin, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[14] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," 2021. [Online]. Available: https://arxiv.org/abs/2010.11929

[15] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," 2021. [Online]. Available: https://arxiv.org/abs/2103.14030

[16] Z. Liu, H. Hu, Y. Lin, Z. Yao, Z. Xie, Y. Wei, J. Ning, Y. Cao, Z. Zhang, L. Dong *et al.*, "Swin transformer v2: Scaling up capacity and resolution," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 12 009–12 019.

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

[18] K. Choromanski, V. Likhosherstov, D. Dohan, X. Song, A. Gane, T. Sarlos, P. Hawkins, J. Davis, A. Mohiuddin, L. Kaiser *et al.*, "Rethinking attention with performers," *arXiv preprint arXiv:2009.14794*, 2020.

[19] T. Dao and A. Gu, "Transformers are ssms: Generalized models and efficient algorithms through structured state space duality," *arXiv preprint arXiv:2405.21060*, 2024.

[20] Y. Shi, M. Dong, M. Li, and C. Xu, "Vssd: Vision mamba with non-causal state space duality," *arXiv preprint arXiv:2407.18559*, 2024.

[21] W. Zhou, S.-i. Kamata, H. Wang, M. S. Wong, and H. C. Hou, "Mamba-in-mamba: Centralized mamba-cross-scan in tokenized mamba model for hyperspectral image classification," *Neurocomputing*, p. 128751, 2024.

[22] L. Zhu, B. Liao, Q. Zhang, X. Wang, W. Liu, and X. Wang, "Vision mamba: Efficient visual representation learning with bidirectional state space model," *arXiv preprint arXiv:2401.09417*, 2024.