# P-Mamba: Marrying Perona Malik Diffusion with Mamba for Efficient Pediatric Echocardiographic Left Ventricular Segmentation

1st Zi Ye
*Institute of Intelligent Software*
*Guangzhou, China*
yezi1022@gmail.com

2nd Tianxiang Chen
*University of Science and Technology of China*
*China*
txchen@mail.ustc.edu.cn

3rd Fangyijie Wang
*School of Medicine*
*University College Dublin*
Dublin, Ireland
fangyijie.wang@ucdconnect.ie

4th Hanwei Zhang
*Saarland University*
*Germany*
zhang@depend.uni-saarland.de

5th Lijun Zhang
*SKLCS, Institute of Software*
*University of Chinese Academy of Sciences*
China
zhanglj@ios.ac.cn

*Abstract*—In pediatric cardiology, the accurate and immediate assessment of cardiac function through echocardiography is crucial since it can determine whether urgent intervention is required in many emergencies. However, echocardiography is characterized by ambiguity and heavy background noise interference, causing more difficulty in accurate segmentation. Present methods lack efficiency and are prone to mistakenly segmenting some background noise areas, such as the left ventricular area, due to noise disturbance. To address these issues, we introduce P-Mamba, which integrates the Mixture of Experts (MoE) concept for efficient pediatric echocardiographic left ventricular segmentation. Specifically, we utilize the recently proposed ViM layers from the vision mamba to enhance our model's computational and memory efficiency while modeling global dependencies. In the DWT-based Perona-Malik Diffusion (PMD) Block, we devise a PMD Block for noise suppression while preserving the left ventricle's local shape cues. Consequently, our proposed P-Mamba innovatively combines the PMD's noise suppression and local feature extraction capabilities with Mamba's efficient design for global dependency modeling. We conducted segmentation experiments on two pediatric ultrasound datasets and a general ultrasound dataset, namely Echonet-dynamic, and achieved state-of-the-art (SOTA) results. Leveraging the strengths of the P-Mamba block, our model demonstrates superior accuracy and efficiency compared to established models, including vision transformers with quadratic and linear computational complexity.

*Index Terms*—Left Ventricular Segmentation, Mamba, Mixture of Experts, Pediatric Echocardiography, Perona–Malik Diffusion

Fig. 1. Visualization of the noise interference challenge among P-Mamba and other models on the pediatric A4C dataset.

## I. INTRODUCTION

CONGENITAL heart diseases (CHD) pose significant health risks, necessitating precise diagnostic tools like the Echocardiogram for early detection and treatment in children [1]. Among these indices, the Left Ventricular Ejection Fraction (LVEF) is the most commonly used and vital metric for assessing systolic function [2]. The LVEF is predominantly calcu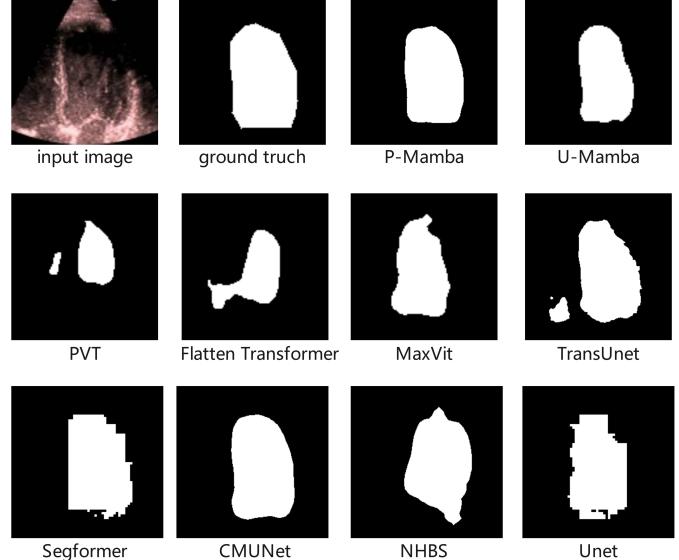lated using the biplane Simpson's standard protocol method in clinical settings [3]. This technique involves the manual delineation of the left ventricular endocardium by physicians in specific frames of the apical two-chamber (A2C) and apical four-chamber (A4C) echocardiographic views, a process essential for the identification of the end-systolic volume (LVESV) and end-diastolic volume (LVEDV).

Machine learning and artificial intelligence have significantly improved the dependability and precision of evaluating left ventricular (LV) function using echocardiography in adults, as shown by multiple research projects. However, machine learning is more difficult in youngsters due to diverse

anatomical anomalies, heart rate, stature, and cooperative capacity. Various factors influence the spatial and temporal resolutions, eventually influencing echocardiographic imaging quality [4]. As a result, there is concern about how well machine learning models built on adult datasets can be applied in pediatric echocardiography owing to the more significant variabilities.

State Space Models (SSMs) have recently extracted broad interest from researchers. As the first proposed basic model built by SSMs, Mamba [5] has achieved superior performance compared to transformers in long-range dependency modeling, even with linear complexity. Vision Mamba [6] was later proposed to apply Mamba to the vision domain and achieves superb efficiency-accuracy balance compared with DeiT. U-Mamba [7] was the first Mamba-based medical image segmentation method and boasted its efficient design. Based on the above works, we try to apply the Mamba structure to our model to guarantee model efficiency.
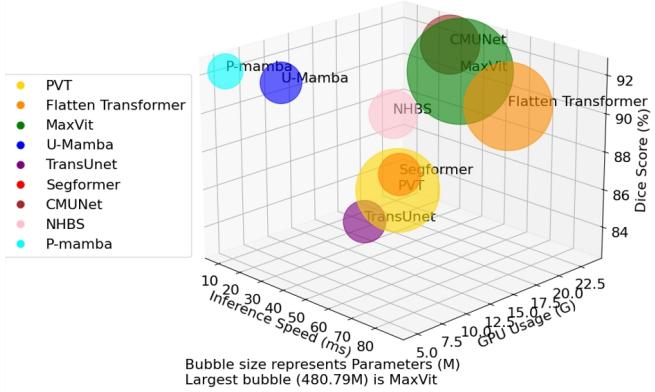


Fig. 2. The efficiency comparisons between P-Mamba and other models on the pediatric PSAX dataset

On the other hand, echocardiography often encounters challenges, including significant speckle noise, limiting its imaging technique. As shown in Fig. 1, current approaches easily mistake some background noise areas for the target area. To suppress background noise while maintaining target structural features, we draw inspiration from Perona–Malik Diffusion (PMD) [8], initially in de-noising tasks to achieve this goal. Therefore, we tailor P-Mamba for more efficient and accurate pediatric echocardiographic LV segmentation. This model can eliminate noise while preserving the local target boundary details for the best performance. At the same time, our model demonstrates superior efficiency.

Our main contributions are as follows:
- We propose P-Mamba, which innovatively combines the PMD's noise suppression and local feature extraction ability with the Vision Mamba's efficient design for global dependency modeling and integrates the Mixture of Experts (MoE) concept to set new performance benchmarks on noisy pediatric echocardiogram datasets.

- Benefiting from PMD, our model excels by suppressing noise while preserving and enhancing target edges in ultrasound images, as shown in Fig. 1.
- Extensive experiments demonstrate that our P-Mamba achieves superior segmentation accuracy and efficiency to other methods, including specialized ultrasound segmentation models [9], [10], vision transformers with quadratic [11] and linear [12], [13] computational complexity. In addition, Fig. 2 visually compares the efficiency among different models.

## II. RELATED WORK

To date, deep learning (DL) development has promoted automatic medical image segmentation, and several well-known deep learning frameworks have provided good ideas for echocardiography segmentation with outstanding performance.

### A. General Deep Learning Segmentation Methods

Deep learning has revolutionized image segmentation in recent years, a crucial task in computer vision. Among the pioneering methods, U-Net [14] stands out for its effective encoder-decoder architecture, which captures context through its contracting path and refines localization through its expansive path. U-Net's simplicity and efficiency have made it a foundational medical image segmentation model.

Pyramid Vision Transformer (PVT) [11] and TransUNet [15] both draw significant inspiration from transformer architecture. PVT captures long-range dependencies through self-attention mechanisms and hierarchical feature representations, making it highly effective for dense prediction tasks. Similarly, TransUNet integrates transformer modules into the traditional U-Net structure, leveraging the transformer's attention mechanisms to enhance feature encoding and significantly improve segmentation accuracy, particularly in complex scenarios. Both models demonstrate the transformative impact of incorporating transformer principles into segmentation tasks.

UniFormer [16], SpectFormer [17], and Segformer [18] represent some of the most contemporary segmentation models. Segformer achieves state-of-the-art performance in multiple segmentation benchmarks with lower computing costs by integrating hierarchical transformers with effective feature fusion approaches. This simplifies the transformer architecture while maintaining its key benefits. UniFormer merges convolutional networks with transformers to balance local feature extraction and global context understanding, efficiently capturing multi-scale features to enhance segmentation performance. SpectFormer delves deeper into the transformer domain by incorporating spectral analysis, emphasizing frequency domain information to complement traditional spatial representations. These models exemplify the cutting-edge advancements in segmentation technology.

### B. Medical Image Segmentation Methods

Recently, innovative deep-learning methods have become increasingly vital in medical analysis and represent the forefront of leveraging complex neural networks for precise medical image segmentation. For instance, CMU-Net [9] and

NHBS-Net [10] are specialized ultrasound segmentation models. CMU-Net uses hybrid convolution and multi-scale attention gates to improve feature extraction and context information. NHBS-Net has advanced attention and fusion modules, making clinical ultrasound applications more accurate and error-free.

However, when applying the present segmentation methods to segment the left ventricle in pediatric echocardiograms, the challenge exists since the irregular shape of the left ventricles still cannot be well segmented since these methods pay insufficient attention to high-frequency boundary details, which can be seen in Fig. 1. Also, present methods lack efficiency, hindering their wider application, as shown in Fig. 2. Another concern of the studies above is the limitations of their reliance on proprietary datasets. The EchoNet-Peds dataset, developed by Stanford University, indicates the first publicly available pediatric echocardiography dataset [19], featuring 4,467 echocardiograms from 1,958 patients, including a 43% female demographic and ages ranging from newborns to 18 years. This comprehensive dataset yielded 7,643 video clips and 17,600 labeled images, primarily from A4C and PSAX view clips. As a result, the video clips were strategically allocated, with 6,114 (80%) for training, 765 (10%) for testing, as well as 764 (10%) for validation.

### C. Mamba

Mamba [5] is a novel deep sequence model architecture addressing the computational inefficiency of traditional Transformers on long sequences. It is based on selective state space models (SSMs), which improve upon previous SSMs by allowing the model parameters to be input functions. Here, some key concepts related to Mamba are explained.

*1) Selective State Space Model:* Consider a structured SSM mapping one-dimensional sequence $x(t) \in \mathbb{R}^L$ to $y(t) \in \mathbb{R}^L$ through a hidden state $h(t) \in \mathbb{R}^N$. With the evolution parameter $\boldsymbol{A} \in \mathbb{R}^{N \times N}$ and the projection parameters $\boldsymbol{B} \in \mathbb{R}^{N \times 1}$, $\boldsymbol{C} \in \mathbb{R}^{1 \times N}$, such a model is formulated using linear ordinary differential equations

$$
\begin{aligned}
h'(t) &= \boldsymbol{A}h(t) + \boldsymbol{B}x(t), \\
y(t) &= \boldsymbol{C}h(t).
\end{aligned} \tag{1}
$$

As a continuous-time model, SSM is discretized with a Zero-Order Hold (ZOH) assumption to adapt to deep learning. Therefore, the continuous-time parameters $\boldsymbol{A}, \boldsymbol{B}$ are transformed to their discretized counterparts $\overline{\boldsymbol{A}}, \overline{\boldsymbol{B}}$ with a timescale parameter $\Delta$ according to

$$
\begin{aligned}
\overline{\boldsymbol{A}} &= \exp(\Delta \boldsymbol{A}), \\
\overline{\boldsymbol{B}} &= (\Delta \boldsymbol{A})^{-1}(\exp(\Delta \boldsymbol{A}) - \boldsymbol{I}) \cdot \Delta \boldsymbol{B}.
\end{aligned} \tag{2}
$$

Thus, (1) can be rewritten as

$$
\begin{aligned}
h_t &= \overline{\boldsymbol{A}}h_{t-1} + \overline{\boldsymbol{B}}x_t, \\
y_t &= \boldsymbol{C}h_t.
\end{aligned} \tag{3}
$$

To enhance computational efficiency and scalability, the iterative process in (3) can be synthesized through a global convolution

$$
\begin{aligned}
\overline{\boldsymbol{K}} &= (\boldsymbol{C}\overline{\boldsymbol{B}}, \boldsymbol{C}\overline{\boldsymbol{A}}\overline{\boldsymbol{B}}, \cdots, \overline{\boldsymbol{A}}^{L-1}\overline{\boldsymbol{B}}), \\
\boldsymbol{y} &= \boldsymbol{x} * \overline{\boldsymbol{K}},
\end{aligned} \tag{4}
$$

where $L$ is the length of the input sequence $\boldsymbol{x}$, $\overline{\boldsymbol{K}} \in \mathbb{R}^L$ serves as the kernel of the SSM and $*$ represents the convolution operation.

Traditional SSM demonstrated linear time complexity, but their representativity of sequence context is inherently limited by time-invariant parameterization. To overcome the existing constraint, Selective SSM introduces a selective scan for interactions among sequential states with

$$
\begin{aligned}
\boldsymbol{B} &= S_{\boldsymbol{B}}(\boldsymbol{x}), \\
\boldsymbol{C} &= S_{\boldsymbol{C}}(\boldsymbol{x}), \\
\Delta &= \tau_{\Delta}(\Delta + S_{\Delta}(\boldsymbol{x})),
\end{aligned} \tag{5}
$$

before (2,3), so that parameters $\boldsymbol{B} \in \mathbb{R}^{B \times L \times N}$, $\boldsymbol{C}^{B \times L \times N}$ and $\Delta^{B \times L \times D}$ are dependent on the input sequence $\boldsymbol{x} \in \mathbb{R}^{B \times L \times D}$, where $B$ represents the batch size, and $D$ represents number of channels. Normally, $S_B$ and $S_C$ are linear parameterized projections to dimension $N$, that is, $Linear_N(\cdot)$, while $S_{\Delta}(\boldsymbol{x}) = Broadcast_D(Linear_1(\boldsymbol{x}))$ and $\tau_{\Delta} = softplus$. The choice of $S_{\Delta}$ and $\tau_{\Delta}$ is caused by a relationship with RNNs gating mechanisms explained later.

*2) Selection Mechanism:* There is a well-established link between discretizing continuous-time systems and RNNs gating [20]. One example of the selection mechanism for SSM is the traditional gating mechanism of RNNs. When $N = 1$, $\boldsymbol{A} = -1$, $\boldsymbol{B} = 1$, $S_{\Delta} = Linear(\boldsymbol{x})$ and $\tau_{\Delta} = softplus$, then the selective SSM recurrence takes the form:

$$
\begin{aligned}
g_t &= \sigma(Linear(x(t))), \\
h_t &= (1 - g_t)h_{t-1} + g_t x_t,
\end{aligned} \tag{6}
$$

*3) Scan:* The selection mechanism is devised to address the constraints of Linear Time Invariance (LTI) models. However, it reintroduces the computation issue in association with SSM. To enhance GPU utilization and efficiently materialize the state $h$ within the memory hierarchy, hardware-aware state expansion is enabled by selective scan. By incorporating kernel fusion and recomputation with parallel scan, the fused selective scan layer can effectively decrease the quantity of memory I/O operations, leading to a significant acceleration compared to conventional implementations.

### III. METHODOLOGY

We illustrate the overall architecture of P-Mamba in Fig. 3. The P-Mamba network, designed for automatically segmenting the left ventricle in echocardiograms, leverages the Mixture of Experts framework. The DWT-based PMD block suppresses background noise interference while preserving target edges to capture more local features. The ViM block [21], sometimes also mentioned as the Bidirectional Mamba block, is designed
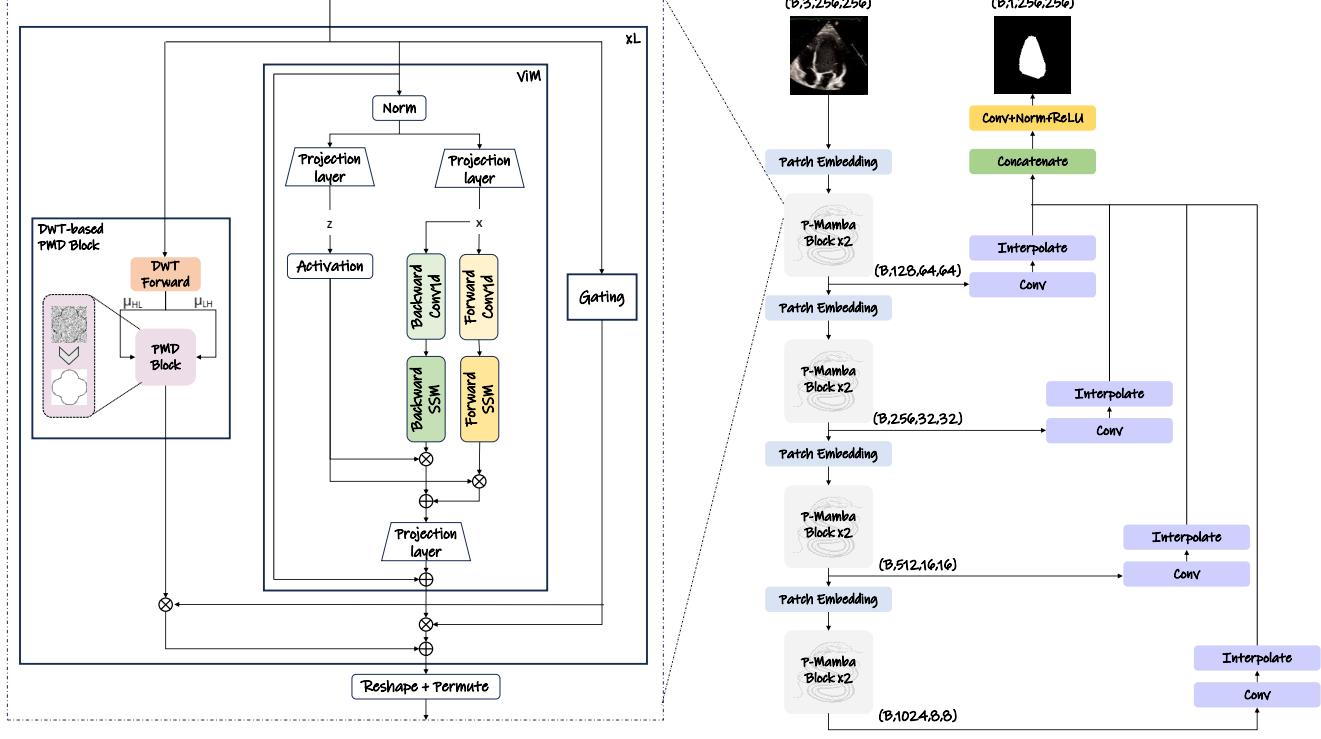
Fig. 3. The structure of our P-Mamba.

to guarantee the high efficiency of our model while encapsulating global dependencies. A gating mechanism inspired by the MoE structure dynamically selects appropriate pathways for feature processing, enhancing efficiency and segmentation accuracy.

The P-Mamba encoder, incorporating multiple P-Mamba blocks, extracts hierarchical features at descending resolutions of 1/4, 1/8, 1/16, and 1/32 relative to the original input size while progressively increasing the channels at each level. In the following sections, we will introduce the detailed structure of P-Mamba.

### A. DWT-based PMD Block

Initially used in image de-noising tasks, Perona-Malik Diffusion (PMD) can suppress noise disturbance while preserving boundary details. Considering that echocardiograms contain heavy background noise that may interfere with segmentation accuracy, we propose the DWT-based PMD Block to act on feature maps so that the background noise can be filtered. At the same time, some target boundary cues can be preserved.

Given an input feature map $u$, its PMD equation can be expressed as:

$$\frac{\partial u}{\partial t} = div\left(g\left(|\nabla u|\right)\nabla u\right) \quad (7)$$

where $g(|\nabla u|) = \frac{1}{1+\left(\frac{|\nabla u|}{k}\right)^2}$ is the diffusion coefficient; $t$ is the diffusion step and can be regarded as the layer of the

feature map; $k$ is a positive constant to control the degree of diffusion and is set to 1 by default in our experiments. Notably, (7) is an anisotropic diffusion equation. In the flat or smooth regions where the gradient magnitude is small ($|\nabla u| \to 0$), the diffusion coefficient $g$ is large, meaning that the diffusion is strong and (7) acts as Gaussian smoothing to remove the noise interference. Somewhere near the target's boundary, the gradient magnitude is large ($|\nabla u| \to 1$), so the coefficient $g$ is near zero, meaning the diffusion is weak, and the boundary details can be preserved. Equation (7) can also be rewritten to the following form:

$$\frac{\partial u}{\partial t} = \frac{\partial}{\partial x}\left\{g\left(\sqrt{\left(\frac{\partial u_k}{\partial x}\right)^2 + \left(\frac{\partial u_k}{\partial y}\right)^2}\right)\frac{\partial u_k}{\partial x}\right\} + \frac{\partial}{\partial y}\left\{g\left(\sqrt{\left(\frac{\partial u_k}{\partial x}\right)^2 + \left(\frac{\partial u_k}{\partial y}\right)^2}\right)\frac{\partial u_k}{\partial y}\right\} \quad (8)$$

where $\frac{\partial u}{\partial x}$ and $\frac{\partial u}{\partial y}$ represent the gradients of the feature map in horizontal and vertical directions. On the other hand, the Discrete Wavelet Transform (DWT) of an input feature map can be expressed as:

$$u_i = DWT(u), i \in \{u_{LL}, u_{LH}, u_{HL}, u_{HH}\} \quad (9)$$

where $u_{LL}$ is the low-frequency part of the feature map, while $u_{LH}$, $u_{HL}$ and $u_{HH}$ are the high-frequency parts in horizontal, vertical and diagonal directions of the feature map which

mainly contain the boundary details. By approximating the derivative terms $\frac{\partial u}{\partial x}$ with $u_{LH}$ and $\frac{\partial u}{\partial y}$ with $u_{HL}$ and setting the diffusion step size $\delta t$ to one, we can transform (8) to the discrete format:

$$
\begin{aligned}
u_k = & u_{k-1} + \left[ g\left( \sqrt{u_{LH}^2 + u_{HL}^2} \right) \cdot u_{LH} \right]_{LH} \\
& + \left[ g\left( \sqrt{u_{LH}^2 + u_{HL}^2} \right) \cdot u_{HL} \right]_{HL}
\end{aligned} \tag{10}
$$

After enhancing the feature map by PMD, we feed the diffusion output into a basic ResNet block. Piling multiple DWT-based PMD blocks in all layers of the encoder branch, our P-Mamba can suppress background noise disturbance while preserving the target boundary features.

### B. Vision Mamba Block

We mainly adopt the ViM layer to improve our model's computing and memory efficiency. Also, it helps capture global dependencies complementary to the local shape cues extracted by our DWT-based PMD Block.

Before processing into the P-Mamba block, a 2-D input $\in \mathbb{R}^{H \times W \times C}$ is transformed into flattened patches $P_N$ with dimensions $M \times (N^2 \cdot C)$. Here $(H, W)$ is the size of the original input, $C$ stands for the number of channels, and $N$ and $M$ denote the size and total count of segmented patches, respectively. Then, the $P_N$ is linearly projected to vectors of dimension D and add position embedding $E_{\text{pos}} \in \mathbb{R}^{M \times D}$. This process can be described as follows:

$$
X_0 = [x^1 W; x^2 W; \ldots; x^M W] + E_{\text{pos}} \tag{11}
$$

where $x^M$ is the $M^{th}$ patch of $P_N$, $W \in \mathbb{R}^{(N^2 \cdot C) \times D}$ is the learnable projection matrix. The token sequence from the patch embedding layer, $X_{pe}$, is processed by the layer ViM to obtain the output $X_{vim}$, which is expressed by (12).

$$
X_{vim} = Vim(X_{pe}) + X_{pe} \tag{12}
$$

Subsequently, the outputs from the DWT-based PMD module and the ViM module are iteratively combined using the gating mechanism across several layers within the p-Mamba block. The final output is then reshaped and permuted to serve as the input for the next stage, where the process is repeated to obtain the next stage's output.

### C. Decoder

The decoder consists of multiple interpolation layers interspersed with Conv2d operations, progressively upsampling the encoded features to generate high-resolution segmentation maps. Specifically, the decoder processes the output feature map from the different encoder stages through several Conv2d layers, followed by interpolation operations that align the feature maps to a common resolution. These feature maps are then concatenated and passed through a series of Conv2d, batch normalization, and ReLU layers to produce the final segmentation map.

## IV. EXPERIMENT

### A. Datasets

*1) Pediatric Dataset:* The dataset comprises echocardiographic evaluations from patients at Lucile Packard Children's Hospital Stanford from 2014 to 2021, authorized by the Stanford University Institutional Review Board. The dataset contains 4467 echocardiograms collected from 1958 patients, 43% female, aged between 0 and 18 years (mean ± SD: 10 ± 5.4 years). The patients were classified into two groups based on their echocardiographic results: those with structurally normal hearts and average ejection fraction (EF) and those with structurally normal hearts but systolic dysfunction (including dilated cardiomyopathy, chemotherapy-induced systolic dysfunction) without congenital heart disease [19].

After additional processing, the dataset was employed to obtain apical four-chamber (A4C) and parasternal short-axis (PSAX) video clips, totaling 7643 video clips and 17600 annotated pictures. The video clips were partitioned into training (80%, n=6114), testing (10%, n=765), and validation (10%, n=764) sets for machine learning purposes. In addition, 86% of the trials had an ejection fraction (EF) equal to or greater than 55%.

*2) EchoNet-Dynamic Dataset:* To ensure the model's effectiveness on general echocardiogram datasets, not just pediatric ones, we also conducted experiments on the EchoNet-Dynamic dataset [22]. EchoNet-Dynamic, obtained from Stanford University Hospital, is the largest publicly available dataset of apical four-chamber (A4C) cardiac echocardiograms. It includes 10030 echocardiogram clips and 9989 video samples that remained after data cleansing. Each video was examined solely for the end-diastolic and end-systolic frames. As a result, 96 images were excluded due to poor-quality ground truth, leaving 14846 out of 19882 images for the training set. The remaining images were divided into validation (2563 images) and testing (2473 images) sets. End-diastolic and end-systolic frames from the same individuals were grouped for the experiments.

### B. Implementation Details

The computational setup consists of a single Tesla V100-32GB GPU, a 12-core CPU, and 61GB of RAM. The system operates on an Ubuntu 18 environment with CUDA 11.0 and Pytorch 1.13 software.

The network was trained for 50 epochs, beginning with an initial learning rate 1e-4. Batch sizes of 24 were selected for training to obtain a compromise between computational efficiency and model accuracy. The model's performance was assessed every five epochs, and early halting with a patience parameter of 10 was implemented to prevent overfitting. The network architecture was organized with layers configured in depth [2, 2, 2, 2].

## V. RESULTS AND DISCUSSION

Table I offers a quantitative comparison of our P-Mamba with various state-of-the-art methods on the pediatric LV 2D

segmentation task and includes experiments on the general ultrasound dataset EchoNet-Dynamic. The comparison includes CNN-based methods like U-Net [14] with FCN and PSPNet backbones, as well as ViT-based approaches such as PVT [11], Flatten Transformer [12], and MaxVit [13]. Moreover, we also compared Uniformer [16] and SpectFormer [17], both newly introduced segmentation models from the last year, showcasing advanced capabilities in image segmentation. U-Mamba [7] represents a straightforward plan Mamba model characterized by a U-shaped structure designed specifically for effective segmentation tasks. TransUnet [15] and Segformer [18] are recognized as classic segmentation models, and they are widely used in the field due to their robust performance. Additionally, CMUNet [9] and NHBS [10] are specialized models tailored for ultrasound image segmentation, addressing the unique challenges of medical imaging.

As a result, Table I demonstrates the superiority of our P-Mamba model. It achieves the highest average Dice Similarity Coefficient (DSC) on the PSAX and A4C pediatric datasets, with values of 0.9221 and 0.9056, respectively. Furthermore, our model also excels on the EchoNet-Dynamic dataset, with an impressive DSC of 0.9314. These results highlight the effectiveness and robustness of our approach across different datasets, outperforming the listed state-of-the-art methods.

*A. Ablation Studies*

We first ablate the DWT-based PMD Block in Table II across various configurations: 'Ours w/o PMD' means removing the DWT-based PMD part; 'Ours w/ Sobel' means replacing the DWT-based PMD part with a Sobel operator, which is only for edge preservation. Our DWT-based PMD block achieves the best performance on those three datasets. The DWT-based PMD part encompasses the Sobel operator due to its noise suppression and edge preservation capability, while the Sobel operator cannot finely preserve edges during noise removal.

In addition, Table III studies the effect of the Vision Mamba block. We replaced the ViM block with ViT structures of quadratic (PVT) and linear complexity (Flatten ViT, MaxViT). The results, summarized across the PSAX, A4C, and EchoNet-Dynamic datasets, indicate that Vision Mamba consistently achieves a higher Dice coefficient. Our findings indicate that Vision Mamba can be more accurate than other models with linear or quadratic complexity.

To further validate the effectiveness of the mixture of expert's gating network convergence, we compared it with the results of simply adding the outputs of the DWT-based PMD Block and the ViM module. Table IV shows that the gating structure improves the Dice coefficient across three datasets, demonstrating that the gating structure implemented with a linear layer significantly improves precision.

*B. Model Efficiency Comparison*

Table V presents the model efficiency comparison results of P-Mamba with various state-of-the-art methods, considering parameters, inference speed (ms), GPU memory usage (GB),
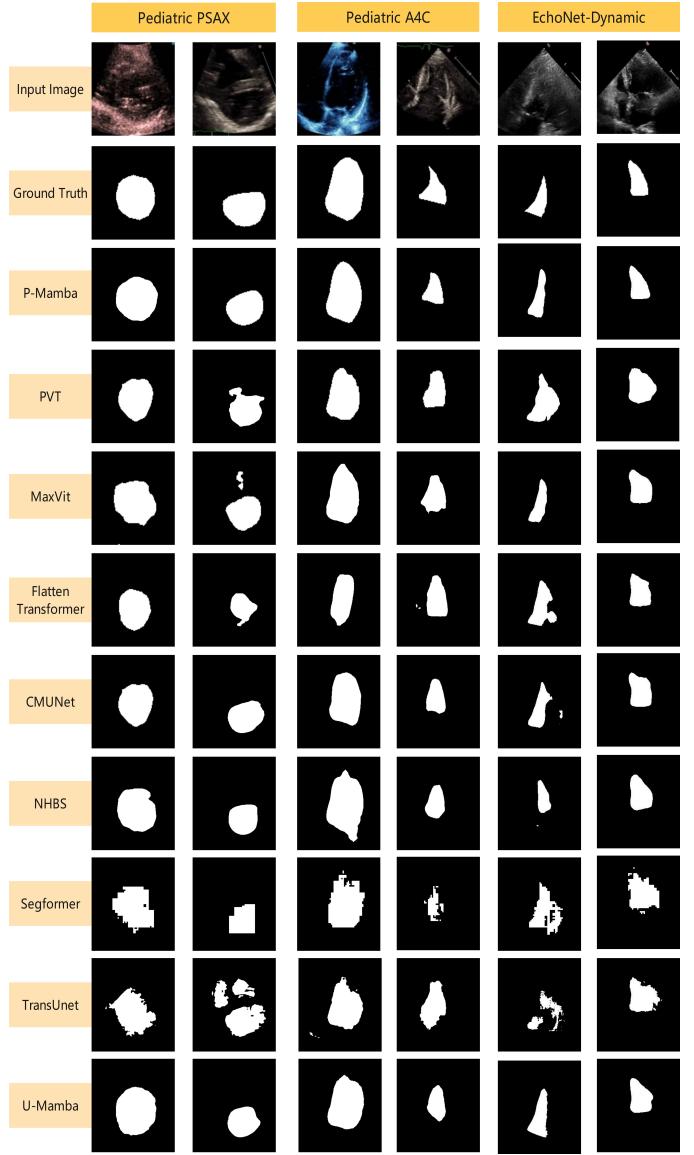


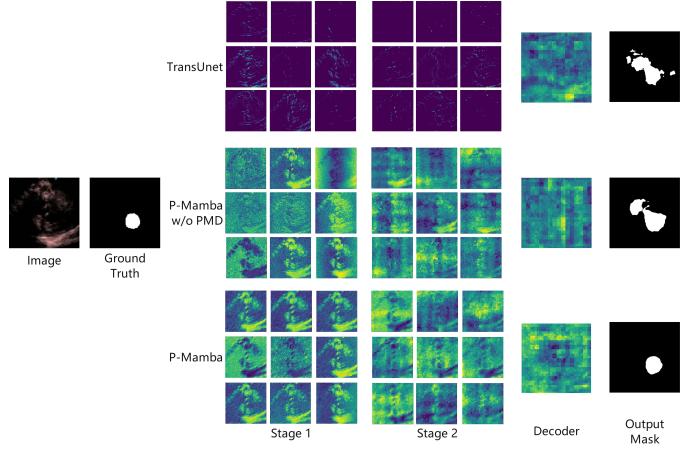Fig. 4. The visual comparison of different methods.



Fig. 5. The Visualization of the features of TransUnet, P-Mamba w/o PMD, and P-Mamba.

TABLE I
COMPARISON WITH THE STATE-OF-THE-ART METHODS.

| Methods | Pediatric PSAX | | | Pediatric A4C | | | EchoNet-Dynamic | | |
|---|---|---|---|---|---|---|---|---|---|
| | Precision | Recall | Dice | Precision | Recall | Dice | Precision | Recall | Dice |
| UNet (FCN) [14] | 0.8492 | 0.8761 | 0.8624 | 0.8279 | 0.8345 | 0.8312 | 0.8837 | 0.8773 | 0.8805 |
| UNet (PSPNet) [14] | 0.8651 | 0.8700 | 0.8675 | 0.8204 | 0.8607 | 0.8401 | 0.8810 | 0.8899 | 0.8854 |
| PVT [11] | 0.8731 | 0.8371 | 0.8547 | 0.8670 | 0.8366 | 0.8515 | 0.9165 | 0.9037 | 0.9101 |
| UniFormer [16] | 0.9100 | 0.9134 | 0.9073 | 0.8969 | 0.8915 | 0.8918 | 0.896 | 0.9157 | 0.9058 |
| SpectFormer [17] | 0.9127 | 0.9063 | 0.9161 | 0.9076 | 0.9035 | 0.9010 | 0.9195 | 0.9234 | 0.9214 |
| Flatten Transformer [12] | 0.9215 | 0.9122 | 0.9168 | 0.9130 | 0.8832 | 0.8978 | 0.9218 | 0.9277 | 0.9247 |
| MaxVit [13] | 0.9024 | 0.9249 | 0.9135 | 0.8912 | 0.9066 | 0.8988 | 0.9246 | 0.9280 | 0.9263 |
| U-Mamba [7] | 0.9275 | 0.9017 | 0.9144 | 0.9022 | 0.8857 | 0.8939 | 0.9331 | 0.9156 | 0.9243 |
| TransUnet [15] | 0.8555 | 0.8059 | 0.8299 | 0.8273 | 0.8390 | 0.8331 | 0.9264 | 0.8991 | 0.9084 |
| Segformer [18] | 0.8646 | 0.8345 | 0.8493 | 0.8481 | 0.8089 | 0.8280 | 0.8917 | 0.8886 | 0.8902 |
| CMUNet [9] | 0.9197 | 0.9126 | 0.9161 | 0.9073 | 0.8959 | 0.9016 | 0.9303 | 0.9266 | 0.9285 |
| NHBS [10] | 0.8734 | 0.8901 | 0.8817 | 0.8864 | 0.8723 | 0.8793 | 0.9224 | 0.9159 | 0.9192 |
| Ours | **0.9316** | 0.9128 | **0.9221** | 0.9045 | **0.9067** | **0.9056** | 0.9186 | **0.9446** | **0.9314** |

TABLE II
ABLATION STUDY ON THE DWT-BASED PMD BLOCK DESIGN.

| Methods | PSAX Dataset | | | A4C Dataset | | | EchoNet-Dynamic | | |
|---|---|---|---|---|---|---|---|---|---|
| | Precision | Recall | Dice | Precision | Recall | Dice | Precision | Recall | Dice |
| Ours w/o PMD | 0.8902 | 0.9098 | 0.8999 | 0.8795 | 0.8811 | 0.8803 | 0.9215 | 0.8931 | 0.9071 |
| Ours w/ Soble | 0.9312 | 0.9013 | 0.9160 | 0.9073 | 0.8983 | 0.9028 | 0.9237 | 0.9166 | 0.9201 |
| Ours | **0.9316** | 0.9128 | **0.9221** | 0.9045 | **0.9067** | **0.9056** | 0.9186 | **0.9446** | **0.9314** |

TABLE III
ABLATION STUDY ON THE ViM BLOCK DESIGN.

| Methods | PSAX Dataset | | | A4C Dataset | | | EchoNet-Dynamic | | |
|---|---|---|---|---|---|---|---|---|---|
| | Precision | Recall | Dice | Precision | Recall | Dice | Precision | Recall | Dice |
| Mamba → PVT | 0.9029 | 0.9183 | 0.9105 | 0.9045 | 0.8714 | 0.8876 | 0.9147 | 0.8907 | 0.9025 |
| Mamba → Flatten Transformer | 0.9137 | 0.9157 | 0.9147 | 0.9228 | 0.8752 | 0.8984 | 0.9198 | 0.8974 | 0.9085 |
| Mamba → MaxVit | 0.9253 | 0.9087 | 0.9169 | 0.9062 | 0.9013 | 0.9038 | 0.9211 | 0.9004 | 0.9106 |
| Ours | **0.9316** | 0.9128 | **0.9221** | 0.9045 | **0.9067** | **0.9056** | 0.9186 | **0.9446** | **0.9314** |

TABLE IV
ABLATION STUDY ON THE CONVERGE DESIGN.

| Methods | PSAX Dataset | | | A4C Dataset | | | EchoNet-Dynamic | | |
|---|---|---|---|---|---|---|---|---|---|
| | Precision | Recall | Dice | Precision | Recall | Dice | Precision | Recall | Dice |
| Adding | 0.9203 | 0.9200 | 0.9201 | 0.9090 | 0.8826 | 0.8956 | 0.9320 | 0.9270 | 0.9295 |
| Ours (Gating) | **0.9316** | 0.9128 | **0.9221** | 0.9045 | **0.9067** | **0.9056** | 0.9186 | **0.9446** | **0.9314** |

and GFLOPs. Our P-Mamba significantly outperforms other methods across all metrics. Specifically, P-Mamba has the lowest parameter count (52.77M), fastest inference speed (9.18 ms), least GPU memory usage (4.95 GB), and lowest GFLOPs (24.66). The attention-free design of our Mamba block greatly enhances model efficiency, even compared to models with linear complexity. Additionally, our DWT-based PMD block does not add excessive parameters compared to the pure mamba net like U-Mamba, ensuring high efficiency.

TABLE V
MODEL EFFICIENCY COMPARISON REGARDING PARAMETER NUMBER (M), INFERENCE SPEED (MS), GPU MEMORY (GB), AND GFLOPs.

| Methods | #Params | ms/Inf | GPU memory | GFLOPs |
|---|---|---|---|---|
| PVT [11] | 299.81 | 40.39 | 15.39 | 181.88 |
| Flatten Transformer [12] | 333.43 | 86.12 | 15.55 | 184.08 |
| MaxVit [13] | 480.79 | 48.49 | 20.08 | 259.48 |
| U-Mamba [7] | 74.00 | 18.64 | 8.38 | 37.14 |
| TransUnet [15] | 77.08 | 23.82 | 15.67 | 36.69 |
| Segformer [18] | 77.20 | 24.07 | 19.48 | 38.90 |
| CMUNet [9] | 149.93 | 30.22 | 23.54 | 182.68 |
| NHBS [10] | 102.20 | 21.64 | 19.35 | 119.49 |
| Ours | **52.77** | **9.18** | **4.95** | **24.66** |

## C. Qualitative Comparison

The qualitative results in Fig. 4 illustrate the performance of various models on pediatric PSAX, pediatric A4C, and EchoNet-Dynamic image segmentation tasks. Classic segmentation models like Segformer and TransUnet display poor segmentation results, indicating their unsuitability for echocardiogram data modalities. Due to noise interference, PVT, Flatten Transformer, MaxVit, and U-Mamba struggle to delineate heart structures correctly. Specialized ultrasound segmentation models such as NHBS and CMUNet also fall short in accurately segmenting the A4C views. In contrast, our P-Mamba is the least affected by noise interference and enjoys the best segmentation effect, thanks to our PMD design.

We visualize the learned features of TransUnet, P-Mamba without the DWT-based PMD Block, and P-Mamba to better understand the impact of the proposed method. The feature map outputs of encoder stages 1 and 2 are shown in Fig. 5, with nine feature map channels randomly sampled for each method. With more distinct feature contours, local information is better preserved in P-Mamba than in other methods. Additionally, we visualize the feature map output of the final decoder stage in Fig. 5 to support our analysis. We observe that P-Mamba features are more concentrated and exhibit stronger discriminative power.

## VI. CONCLUSION

We present P-Mamba, a model tailored for efficient left ventricular segmentation in pediatric echocardiography, overcoming the background noise interference challenge. Specifically, the DWT-based Perona-Malik Diffusion Block is a mathematically explainable module focusing on local feature extraction. It can gradually suppress background noise disturbance while preserving the boundary details of the left ventricle. The vision Mamba block, on the other hand, can explore global dependencies. We integrate the two parts by borrowing the Mixture of Experts (MOE) ideas, and our P-Mamba successfully achieves an exceptional balance between segmentation accuracy and efficiency.

## REFERENCES

[1] D. V. D. Linde, E. E. Konings, M. A. Slager, M. Witsenburg, W. A. Helbing, J. J. Takkenberg, and J. W. Roos-Hesselink, "Birth prevalence of congenital heart disease worldwide: a systematic review and meta-analysis," *J. Am. Coll. Cardiol.*, vol. 58, no. 21, pp. 2241–2247, Nov. 2011.

[2] J. F. Silva, J. M. Silva, A. Guerra, S. Matos, and C. Costa, "Ejection fraction classification in transthoracic echocardiography using a deep learning approach," in *Proc. 2018 IEEE 31st Int. Symp. Comput.-Based Med. Syst. (CBMS)*, Karlstad, Sweden, 2018, pp. 123–128.

[3] R. M. Lang, L. P. Badano, V. Mor-Avi, J. Afilalo, A. Armstrong, L. Ernande, F. A. Flachskampf, E. Foster, S. A. Goldstein, and T. Kuznetsova, "Recommendations for cardiac chamber quantification by echocardiography in adults: an update from the american society of echocardiography and the european association of cardiovascular imaging," *Eur. Heart J. Cardiovasc. Imaging*, vol. 16, no. 3, pp. 233–271, Mar. 2015.

[4] A. Power, S. Poonja, D. Disler, K. Myers, D. J. Patton, J. K. Mah, N. M. Fine, and S. C. Greenway, "Echocardiographic image quality deteriorates with age in children and young adults with duchenne muscular dystrophy," *Front. Cardiovasc. Med.*, vol. 4, Dec. 2017.

[5] A. Gu and T. Dao, "Mamba: Linear-time sequence modeling with selective state spaces," 2024, accessed on: May 22, 2024. [Online]. Available: https://arxiv.org/abs/2312.00752

[6] L. Zhu, B. Liao, Q. Zhang, X. Wang, W. Liu, and X. Wang, "Vision mamba: Efficient visual representation learning with bidirectional state space model," 2024, accessed on: May 22, 2024. [Online]. Available: https://arxiv.org/abs/2401.09417

[7] J. Ma, F. Li, and B. Wang, "U-mamba: Enhancing long-range dependency for biomedical image segmentation," 2024, accessed on: May 22, 2024. [Online]. Available: https://arxiv.org/abs/2401.04722

[8] P. Perona and J. Malik, "Scale-space and edge detection using anisotropic diffusion," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 12, no. 7, pp. 629–639, Jul. 1990.

[9] F. Tang, L. Wang, C. Ning, M. Xian, and J. Ding, "Cmu-net: A strong convmixer-based medical ultrasound image segmentation network," in *Proc. 2023 IEEE 20th Int. Symp. Biomed. Imaging (ISBI)*, Cartagena, Colombia, 2023, pp. 1–5.

[10] R. Liu, M. Liu, B. Sheng, H. Li, P. Li, H. Song, P. Zhang, L. Jiang, and D. Shen, "Nhbs-net: A feature fusion attention network for ultrasound neonatal hip bone segmentation," *IEEE Trans. Medical Imaging*, vol. 40, no. 12, pp. 3446–3458, Dec. 2021.

[11] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, and et al., "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," in *Proc. 2021 IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Montreal, QC, Canada, 2021, pp. 568–578.

[12] D. Han, X. Pan, Y. Han, S. Song, and G. Huang, "Flatten transformer: Vision transformer using focused linear attention," in *Proc. 2023 IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Paris, France, 2023, pp. 5961–5971.

[13] Z. Tu, H. Talebi, H. Zhang, F. Yang, P. Milanfar, A. Bovik, and et al., "Maxvit: Multi-axis vision transformer," in *European Conf. Comput. Vis. (ECCV)*. Tel Aviv, Israel: Springer, 2022, pp. 459–479.

[14] D. Gupta, "Image segmentation keras: Implementation of segnet, fcn, unet, pspnet and other models in keras," 2024, accessed on: May 22, 2024. [Online]. Available: https://arxiv.org/abs/2307.13215

[15] J. Chen, Y. Lu, Q. Yu, X. Luo, E. Adeli, Y. Wang, L. Lu, A. L. Yuille, and Y. Zhou, "Transunet: Transformers make strong encoders for medical image segmentation," 2024, accessed on: May 22, 2024. [Online]. Available: https://arxiv.org/abs/2102.04306

[16] K. Li, Y. Wang, J. Zhang, P. Gao, G. Song, Y. Liu, H. Li, and Y. Qiao, "Uniformer: Unifying convolution and self-attention for visual recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 10, pp. 12 581–12 600, Oct. 2023.

[17] B. N. Patro, V. P. Namboodiri, and V. S. Agneeswaran, "Spectformer: Frequency and attention is what you need in a vision transformer," 2024, accessed on: May 22, 2024. [Online]. Available: https://arxiv.org/abs/2304.06446

[18] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, "Segformer: Simple and efficient design for semantic segmentation with transformers," 2024, accessed on: May 22, 2024. [Online]. Available: https://arxiv.org/abs/2105.15203

[19] C. D. Reddy, L. Lopez, D. Ouyang, J. Y. Zou, and B. He, "Video-based deep learning for automated assessment of left ventricular ejection fraction in pediatric patients," *J. Am. Soc. Echocardiogr.*, vol. 36, no. 5, pp. 482–489, May 2023.

[20] C. Tallec and Y. Ollivier, "Can recurrent neural networks warp time?" 2024, accessed on: May 22, 2024. [Online]. Available: https://arxiv.org/abs/1804.11188

[21] L. Zhu, B. Liao, Q. Zhang, X. Wang, W. Liu, and X. Wang, "Vision mamba: Efficient visual representation learning with bidirectional state space model," 2024, accessed on: May 22, 2024. [Online]. Available: https://arxiv.org/abs/2401.09417

[22] D. Ouyang, B. He, A. Ghorbani, M. P. Lungren, E. A. Ashley, D. H. Liang, and J. Y. Zou, "Echonet-dynamic: a large new cardiac motion video data resource for medical machine learning," in *NeurIPS ML4H Workshop*, Vancouver, BC, Canada, 2019.