
UltraLight VM-UNet: Parallel Vision Mamba Significantly Reduces Parameters for Skin Lesion Segmentation

Renkai Wu¹ Yinghao Liu² Pengchen Liang^{1*} Qing Chang^{1*}
¹Shanghai University ²University of Shanghai for Science and Technology

Abstract

Traditionally for improving the segmentation performance of models, most approaches prefer to use adding more complex modules. And this is not suitable for the medical field, especially for mobile medical devices, where computationally loaded models are not suitable for real clinical environments due to computational resource constraints. Recently, state-space models (SSMs), represented by Mamba, have become a strong competitor to traditional CNNs and Transformers. In this paper, we deeply explore the key elements of parameter influence in Mamba and propose an UltraLight Vision Mamba UNet (UltraLight VM-UNet) based on this. Specifically, we propose a method for processing features in parallel Vision Mamba, named PVM Layer, which achieves excellent performance with the lowest computational complexity while keeping the overall number of processing channels constant. We conducted comparisons and ablation experiments with several state-of-the-art lightweight models on three skin lesion public datasets and demonstrated that the UltraLight VM-UNet exhibits the same strong performance competitiveness with parameters of only 0.049M and GFLOPs of 0.060. In addition, this study deeply explores the key elements of parameter influence in Mamba, which will lay a theoretical foundation for Mamba to possibly become a new mainstream module for lightweighting in the future. The code is available from <https://github.com/wurenkai/UltraLight-VM-UNet>.

1 Introduction

With the continuous development of computer technology and hardware computing power, computer-aided diagnosis has been widely used in the medical field, and medical image segmentation is an important part of it. Medical image segmentation is usually realized using deep learning networks represented by convolution and Transformers. Convolution has excellent localized feature extraction capabilities, but is deficient in establishing the correlation of remote information [20, 21, 32]. In the previous work [14, 5], researchers proposed to utilize large convolutional kernels to alleviate this problem. As for the network architecture based on Transformers, in recent years, researchers [9, 8] have deeply investigated its methods. Although the self-attention mechanism can solve the problem of remote information extraction by means of sequences of consecutive patches, it also introduces more computational complexity. This is because the quadratic complexity of the self-attention mechanism is closely related to the image size.

In addition, in terms of improving the accuracy of computer-aided diagnosis, the number of parameters of the algorithmic model is often enriched to improve the predictive power of the model [1, 34]. However, for clinical and real healthcare environments, realistic computational power and memory constraints often need to be considered. Low parameters and minimal computational memory

*Corresponding authors.

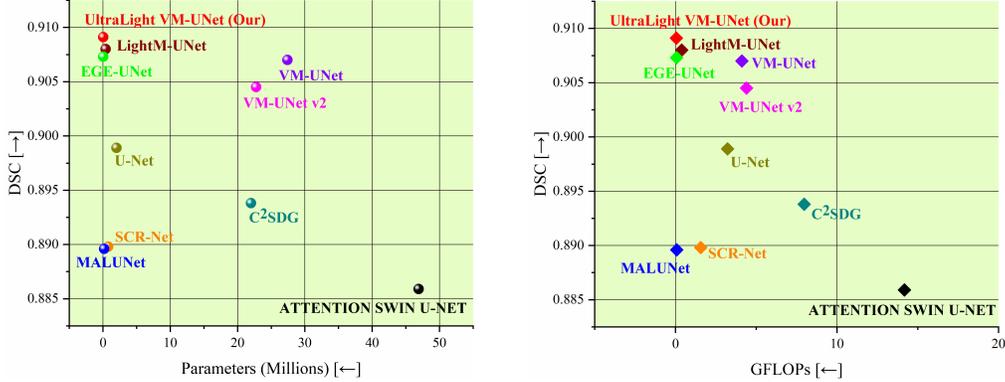


Figure 1: Visualization of the comparison results for the ISIC2017 dataset. X-axis corresponds to parameters and GFLOPs, the fewer the better. Y-axis corresponds to segmentation performance (DSC), the higher the better.

footprint are essential considerations in mHealth tasks [25]. Therefore, an algorithmic model with low computational complexity and good performance is urgently needed for future mobile medical devices.

Recently, state-space models (SSMs) have shown linear complexity in terms of input size and memory occupation [13, 37], which makes them key to lightweight model foundations. In addition, SSMs excel at capturing remote dependencies, which can critically address the problem of convolution for extracting information over long distances. In Gu et al. [7], time-varying parameters were introduced into SSM to obtain Mamba, and it was demonstrated that Mamba is able to process textual information with lower parameters than Transformers. On the vision side, the introduction of Vision Mamba [37] has once again furthered people’s understanding of Mamba, which saves 86.8% of memory when reasoning about images of 1248×1248 size without the need for an attentional mechanism. With the outstanding work of the researchers mentioned above, we are more confident that Mamba will occupy a major position in the future as a basic building block for lightweight models.

In this paper, we are building a lightweight model based on Vision Mamba. We deeply explore the critical memory footprint of Mamba and the performance tradeoffs, and propose an UltraLight Vision Mamba UNet (UltraLight VM-UNet). To the best of our knowledge, the proposed UltraLight VM-UNet is the most lightweight Vision Mamba model available (with parameters of 0.049M and GFLOPs of 0.060) and exhibits very competitive performance in three skin lesion segmentation tasks. Specifically, we delve into the keys affecting the computational complexity in Mamba, and conclude that the number of channels is a key factor in the explosive memory occupation for Mamba computation. We build on this finding of ours by proposing a parallel Vision Mamba approach for processing deep features, named PVM Layer, which simultaneously keeps the overall processing channel count constant. The proposed PVM Layer achieves excellent performance with surprisingly low parameters. In addition, the deep feature extraction of the proposed UltraLight VM-UNet we implement using only the PVM Layer containing Mamba, as shown in Figure 2. In the Methods section, we will present the details of the proposed UltraLight VM-UNet as well as the key factors of the parameter effects in Mamba and the performance balancing approach.

Although similar parallel connection of modules has been mentioned in previous studies [28, 11, 29, 27], the impact of using parallel connection in Mamba is still unknown. Does parallel connection lead to a significant performance degradation of Mamba when utilizing the SSM selection mechanism? This is because parallel connections lead to a reduction in the number of feature channels learned per SSM. In this paper, we will give the answer in detail. Parallel Vision Mamba or Mamba not only remains competitive in terms of performance, but gets a significant reduction in parameters and computational complexity.

Our contributions and findings can be summarized as follows:

- An UltraLight Vision Mamba UNet (UltraLight VM-UNet) is proposed for skin lesion segmentation. To the best of our knowledge, the UltraLight VM-UNet is the lightest Vision Mamba model available (parameters of 0.049M, GFLOPs of 0.060).

- A parallel Vision Mamba method for processing deep features, named PVM Layer, is proposed, which achieves excellent performance with the lowest computational complexity while keeping the overall number of processing channels constant. This can be generalized to the parallel connection of any Mamba variant.
- We provide an in-depth analysis of the key factors influencing the parameters of Mamba, and provide a theoretical basis for Mamba to become a mainstream module for lightweight modeling in the future.
- The proposed UltraLight VM-UNet parameters are 99.82% lower than the traditional pure Vision Mamba UNet model (VM-UNet) and 87.84% lower than the parameters of the current lightest Vision Mamba UNet model (LightM-UNet). In addition, the UltraLight VM-UNet maintains strong performance competitiveness in all three publicly available skin lesion segmentation datasets.

2 Related Work

With the significant improvement of computer computing power, computer vision has become an important field in computer technology nowadays. And deep learning development has led to the full convolutional model (FCN) [15] showing excellent performance in image segmentation methods. Soon after the emergence of FCN, another fully convolutional model (U-Net) [22] emerged to generate renewed excitement. The skip-connection operation in U-Net is able to merge high-level and low-level features well, which is especially important for image segmentation, especially for medical image segmentation that requires fine-grained segmentation.

Medical image segmentation, as one of the important branches in image segmentation, is also one of the research directions to which many researchers have devoted their efforts. Among them, multi-scale variation problem and feature refinement learning are the key problems in medical image segmentation [35]. And skin lesion segmentation has rich and varied feature information as well as high lethality caused by its malignant melanoma [26], which has led many researchers to carry out a series of studies around skin lesion segmentation [1, 32, 34].

Medical image segmentation algorithms represented by skin lesions have been rapidly developed after the advent of U-Net. In Aghdam et al. [1], an inhibition operation of the attention mechanism in cascade operation has been proposed for skin lesion segmentation based on Swin U-Net [2]. MHorUNet [32] model proposes a high-order spatial interaction UNet model for skin lesion segmentation. In Wu et al. [34], an adaptive selection of higher order UNet model for order interaction has been proposed for skin lesion segmentation. In addition, there are very many algorithms based on U-Net improved for skin lesion segmentation. However, it is common for researchers to add richer modules to the model to improve the accuracy of recognition, but this also significantly increases the parameters and computational complexity of the model. After the emergence of Vision Mamba, LightM-UNet [12] was proposed to reduce the number of parameters in the model based on Mamba. LightM-UNet further extracts the deep semantics and tele-relationships by using the residual Vision Mamba, and achieves better performance with a smaller number of parameters. In addition, U-Mamba [17] introduced Vision Mamba into the U-framework for the first time, but its having a large number of parameters (173.53M) limits its use in real clinical settings.

In this paper, in order to solve the current problem of large model parameters and to reveal the key factors affecting Mamba parameters. We propose an UltraLight Vision Mamba UNet (UltraLight VM-UNet) based on Mamba with a parameter of only 0.049M. The UltraLight VM-UNet is confirmed to still maintain a strong competitive performance in three public datasets of skin lesions. In the next section, our method is described in detail.

3 Method

3.1 Architecture Overview

The proposed UltraLight Vision Mamba UNet (UltraLight VM-UNet) is shown in Figure 2. The UltraLight VM-UNet has a total of 6-layer structure consisting of a U-shaped structure (encoder, decoder, and skip-connection path). The number of channels in the 6-layer structure is set to [8, 16, 24, 32, 48, 64]. The extraction of shallow features in the first 3 layers is composed using a convolution

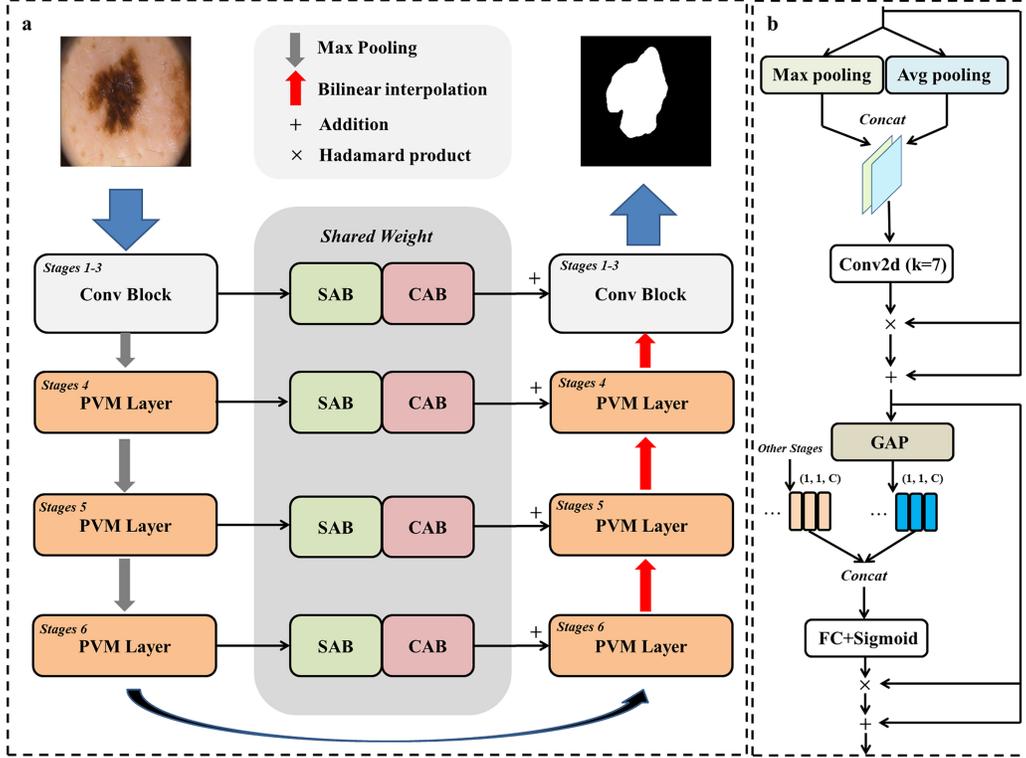


Figure 2: (a) The proposed UltraLight Vision Mamba UNet (UltraLight VM-UNet) model architecture. (b) Multilevel and multiscale information fusion module architecture for skip-connection paths.

module (Conv Block), where each layer includes a standard convolution with a convolution kernel of 3 and a maximum pooling operation. The deeper features from layer 4 to layer 6 are our core part, where each layer consists of our proposed Parallel Vision Mamba Layer (PVM Layer). The decoder part maintains the same setup as the encoder. The skip-connection path utilize the Channel Attention Bridge (CAB) module and the Spatial Attention Bridge (SAB) module for multilevel and multiscale information fusion [24].

3.2 Mamba Parameter Impact Analysis

Vision Mamba for PVM Layer is mainly composed using Mamba combined with residual connections and adjustment factors (Figure 3), which allows traditional Mamba to improve the capture of remote spatial relations without introducing additional parameters and computational complexity [12]. This has a better improvement in the performance of Mamba in visual tasks while keeping the parameter and computational complexity low.

Among SSM-based Mamba, the number of channels, the size of the SSM state dimension, the size of the internal 1D convolutional kernel, the projection dilation multiplier, and the rank of the step size all affect the parameters. And in this, the impact of the channel number is explosive, and its main influence is from the following multiple directions:

First, the d_{inner} of the Mamba internal extended projection channel is determined by the product of the projection expansion multiplier and the number of input channels. This can be specifically expressed by the following equation:

$$d_{inner} = expand * d_{model} \quad (3.1)$$

where d_{inner} is the internal expansion projection channel, $expand$ is the projection expansion multiplier (fixed at 2 by default), and d_{model} is the number of input channels. We can conclude that d_{inner} will double up as the number of channels (d_{model}) per layer in the model increases.

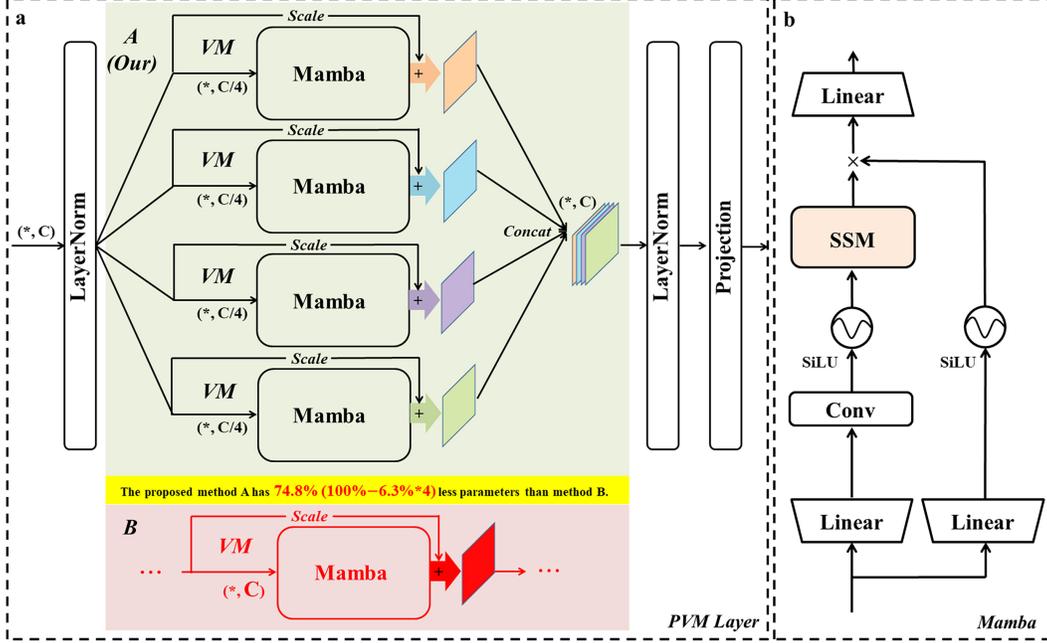


Figure 3: (a) Architecture of the proposed Parallel Vision Mamba Layer (PVM Layer) method. Vision Mamba (VM) is composed by Mamba combined with residual connection and adjustment factor. (b) Mamba composition structure.

Second, the parameters of the input projection layer (the same input linear layer is used for both branches) and output projection layer within Mamba will be directly related to the number of input channels. The input projection layer and output projection layer operate as follows:

$$in_proj : nn.Linear(d_model, d_inner * 2, bias = False) \quad (3.2)$$

$$out_proj : nn.Linear(d_inner, d_model, bias = False) \quad (3.3)$$

where then the input projection (in_proj) layer parameter is $d_model * d_inner * 2$, and the output projection layer (out_proj) parameter is $d_inner * d_model$. we can conclude that the number of input channels, d_model , is the key element controlling the parameter, where the internal extended projection channel, d_inner , is also controlled by d_model is controlled.

Further, the intermediate linear projection layers of SSM are also key to the influence of the parameters. The details are as follows:

$$x_proj = nn.Linear(d_inner, (dt_rank + d_state * 2), bias = False) \quad (3.4)$$

$$dt_proj = nn.Linear(dt_rank, d_inner, bias = True) \quad (3.5)$$

where dt_rank is the rank of the step ($dt_rank = d_model/16$), d_state is the size of the state dimension (fixed to 16), where the parameters can be derived as $d_inner * (dt_rank + d_state * 2)$. dt_proj is a linear projection layer for step size, with parameters $(dt_rank * d_inner) + d_inner$, mainly used for linear projection for step size (dt). So, we can conclude that all parameters are still mainly controlled by the number of input channels d_model .

In addition, the internal convolution ($nn.Conv1d(d_inner, d_inner, d_conv, bias = True)$) also provides parametric influence with $d_inner * d_inner * d_conv + d_inner$. In this paper, d_conv is fixed to 4, so the convolution provides a parameter of $4 * d_inner^2 + d_inner$, which is also controlled by the d_model .

Also, the logarithmic form parameters (A_logs) of the transfer matrix A in the SSM module are an important influencing element. A_logs is a parameter matrix of the shape (d_inner, d_state) , so its parameters can be derived as $d_inner * d_state$. In addition, the trainable vector parameter (D) of the process computed within the SSM contains the d_inner parameter, which is used to

selectively integrate the SSM state outputs with the original input signals, thereby enhancing the model’s expressiveness and training stability.

In summary, assuming that the original channel number is 1024, keeping other parameters unchanged, and when the channel number is reduced to one-fourth of the original (channel number 256), the original total parameters can be calculated from the above parameter formula to get the original total parameters reduced from 23,435,264 to 1,484,288. The parameter explosion is reduced by 93.7%, which further confirms that the number of channels has a very critical impact on Mamba parameters.

Based on the above in-depth study of the key elements affecting the Mamba parameters, we propose a method for processing features in parallel Vision Mamba, named PVM Layer. Excellent performance is achieved with the lowest computational complexity, while keeping the overall number of processing channels unchanged. Specifically, it will be detailed in the next section.

3.3 Parallel Vision Mamba Layer

As analyzed in the previous subsection, the number of input channels has an explosive effect on the parameters of Mamba. As shown in Figure 3(a), we propose the Parallel Vision Mamba Layer (PVM Layer) for processing deep features. Specifically, a feature X with channel number C first passes through a LayerNorm layer and then is divided into $Y_1^{C/4}$, $Y_2^{C/4}$, $Y_3^{C/4}$ and $Y_4^{C/4}$ features each with channel number $C/4$. After that, each of the features is then fed into Mamba, and then the output is subjected to residual concatenation and adjustment factor for optimizing the remote spatial information acquisition capability [12]. Finally, the four features are combined into the feature X_{out} with channel number C by concat operation, and then output by LayerNorm and Projection operation respectively. The specific operation can be expressed by the following equation:

$$Y_1^{C/4}, Y_2^{C/4}, Y_3^{C/4}, Y_4^{C/4} = Sp [LN (X_{in}^C)] \quad (3.6)$$

$$VM_Y_i^{C/4} = Mamba (Y_i^{C/4}) + \theta \cdot Y_i^{C/4} \quad i = 1, 2, 3, 4 \quad (3.7)$$

$$X_{out} = Cat (VM_Y_1^{C/4}, VM_Y_2^{C/4}, VM_Y_3^{C/4}, VM_Y_4^{C/4}) \quad (3.8)$$

$$Out = Pro [LN (X_{out})] \quad (3.9)$$

where LN is the LayerNorm, Sp is the Split operation, $Mamba$ is the Mamba operation, θ is the adjustment factor for the residual connection, Cat is the concat operation, and Pro is the Projection operation. From Eq. 3.7, we used parallel Vision Mamba processing features, while ensuring that the total number of channels processed remains constant, maintaining high accuracy while maximizing parameter reduction. As shown in Figure 3(a) for Methods A and B, again assuming a channel count size of 1024, each Vision Mamba (VM) in Method A reduces the parameters by 93.7%. And it contains 4 such operations, so when summed up, the comparison method B parameters are reduced by 74.8% overall. Through our proposed parallel Vision Mamba operation, the parameter reduction is maximized while maintaining strong performance competitiveness.

3.4 Skip-connection Path

The skip-connection path uses the Spatial Attention Bridge (SAB) module and the Channel Attention Bridge (CAB) module proposed by Ruan et al. [24], as shown in Figure 2(b). The combined use of SAB and CAB allows for the fusion of multi-stage features of different scales of the UltraLight VM-UNet. The SAB module contains the maximum pooling, the average pooling, the extended convolution of the shared weights. The CAB module contains global average pooling, concat operation, fully connected layers, and sigmoid activation function. Both SAB and CAB have been shown to be effective in improving the convergence ability of the model and enhancing the sensitivity to lesions in previous work [24, 32, 33].

4 Experiment

4.1 Comparison results

In order to validate the competitive performance of the proposed UltraLight VM-UNet under the 0.049M parameter, we conducted comparison experiments with several state-of-the-art lightweight

Table 1: Experimental comparisons of the proposed UltraLight VM-UNet with the best lightweight and classical models. This includes comparisons on two publicly available large dermatologic lesion datasets (ISIC2017 and ISIC2018), as well as external validation on a small publicly available dataset (PH²). “↑” stands for external validation experiments. In addition, comparisons of the parameters and computational complexity of the individual models are included. Our models are highlighted in gray.

Comparison experiments on the ISIC2017 dataset.					↑ Comparison experiments on the PH ² dataset.				
Model	DSC	SE	SP	ACC	Model	DSC	SE	SP	ACC
U-Net	0.8989	0.8793	0.9812	0.9613	U-Net	0.9060	0.9255	0.9440	0.9381
SCR-Net	0.8898	0.8497	0.9853	0.9588	SCR-Net	0.8989	0.9114	0.9446	0.9339
ATTENTION SWIN U-NET	0.8859	0.8492	0.9847	0.9591	ATTENTION SWIN U-NET	0.8850	0.8886	0.9363	0.9213
C ² SDG	0.8938	0.8859	0.9765	0.9588	C ² SDG	0.9030	0.9137	0.9476	0.9367
VM-UNet	0.9070	0.8837	0.9842	0.9645	VM-UNet	0.9033	0.9131	0.9483	0.9369
VM-UNet v2	0.9045	0.8768	0.9849	0.9637	VM-UNet v2	0.9050	0.9160	0.9485	0.9380
MALUNet	0.8896	0.8824	0.9762	0.9583	MALUNet	0.8865	0.8922	0.9425	0.9263
LightM-UNet	0.9080	0.8839	0.9846	0.9649	LightM-UNet	0.9156	0.9129	0.9613	0.9457
EGE-UNet	0.9073	0.8931	0.9816	0.9642	EGE-UNet	0.9086	0.9198	0.9502	0.9404
UltraLight VM-UNet (Our)	0.9091	0.9053	0.9790	0.9646	UltraLight VM-UNet (Our)	0.9265	0.9345	0.9606	0.9521

Comparison experiments on the ISIC2018 dataset.					Comparison of parameters and computational consumption.		
Model	DSC	SE	SP	ACC	Model	Params	GFLOPs
U-Net	0.8851	0.8735	0.9744	0.9539	U-Net	2.009M	3.224
SCR-Net	0.8886	0.8892	0.9714	0.9547	SCR-Net	0.801M	1.567
ATTENTION SWIN U-NET	0.8540	0.8057	0.9826	0.9480	ATTENTION SWIN U-NET	46.910M	14.181
C ² SDG	0.8806	0.8970	0.9643	0.9506	C ² SDG	22.001M	7.972
VM-UNet	0.8891	0.8809	0.9743	0.9554	VM-UNet	27.427M	4.112
VM-UNet v2	0.8902	0.8959	0.9702	0.9551	VM-UNet v2	22.771M	4.400
MALUNet	0.8931	0.8890	0.9725	0.9548	MALUNet	0.175M	0.083
LightM-UNet	0.8898	0.8829	0.9765	0.9555	LightM-UNet	0.403M	0.391
EGE-UNet	0.8819	0.9009	0.9638	0.9510	EGE-UNet	0.053M	0.072
UltraLight VM-UNet (Our)	0.8940	0.8680	0.9781	0.9558	UltraLight VM-UNet (Our)	0.049M	0.060

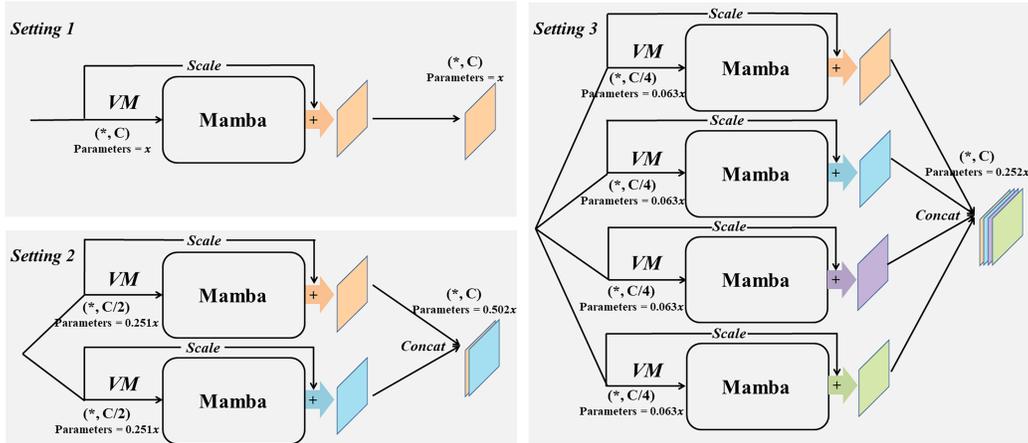


Figure 4: Settings for ablation experiments with Vision Mamba used in different parallel ways (PVM Layer).

and classical medical image segmentation models. Specifically, they include U-Net [22], SCR-Net [31], ATTENTION SWIN U-NET [1], C²SDG [10], VM-UNet [23], VM-UNet v2 [36], MALUNet [24], LightM-UNet [12] and EGE-UNet [25].

Table 1 show the experimental results on ISIC2017, ISIC2018 and PH² datasets, respectively. As shown in the table, the parameters of our model are 99.82% lower than those of the traditional pure Vision Mamba UNet model (VM-UNet) and 87.84% lower than those of the current lightest Vision Mamba UNet model (LightM-UNet). In addition, the GFLOPs of our model are 98.54% lower than VM-UNet and 84.65% lower than LightM-UNet. With such a large reduction in parameters and GFLOPs, the performance of our model still maintains excellent and highly competitive performance. In addition, MALUNet is a lightweight model proposed based on convolution, and although it has lower parameters and GFLOPs than VM-UNet and LightM-UNet, the parameters and GFLOPs of our model are still 72.0% and 27.71% lower than them, respectively. In particular, the performance

of both MALUNet, the proposed lightweight model based on convolution, is much lower than that of the Mamba-based model, which reflects that it is difficult for the convolution-based lightweight model to balance the relationship between performance and computational complexity.

Table 2: Ablation experiments on the effect of Vision Mamba in different parallel connections.

Settings	Params	GFLOPs	DSC	SE	SP	ACC
1	0.136M	0.060	0.9069	0.8861	0.9834	0.9644
2	0.070M	0.060	0.9073	0.8866	0.9835	0.9645
3	0.049M	0.060	0.9091	0.9053	0.9790	0.9646

Table 3: **Impact of adopting parallelism for different SSM variants.** (*P*) shows the adoption of quadruple parallelism to replace the original form. ■ indicates experiments on the ISIC2017 dataset, ♠ indicates the ISIC2018 dataset, and ♣ indicates the PH² dataset.

Methods	Parallelism	SSM-Variants	Params	GFLOPs	DSC	SE	SP	ACC
■ VM-UNet	False	SS2D	27.427M	4.112	0.9070	0.8837	0.9842	0.9645
♠ VM-UNet	False	SS2D	27.427M	4.112	0.8891	0.8809	0.9743	0.9554
♣ VM-UNet	False	SS2D	27.427M	4.112	0.9033	0.9131	0.9483	0.9369
■ VM-UNet	True	(<i>P</i>) SS2D	5.116M	1.564	0.9159	0.8843	0.9886	0.9682
♠ VM-UNet	True	(<i>P</i>) SS2D	5.116M	1.564	0.8977	0.8996	0.9734	0.9584
♣ VM-UNet	True	(<i>P</i>) SS2D	5.116M	1.564	0.9185	0.9340	0.9525	0.9465
■ H-vmunet	False	H-SS2D	8.973M	0.742	0.9172	0.9056	0.9831	0.9680
♠ H-vmunet	False	H-SS2D	8.973M	0.742	0.8966	0.8952	0.9741	0.9581
♣ H-vmunet	False	H-SS2D	8.973M	0.742	0.9146	0.9295	0.9509	0.9440
■ H-vmunet	True	(<i>P</i>) H-SS2D	1.764M	0.318	0.9066	0.8825	0.9843	0.9644
♠ H-vmunet	True	(<i>P</i>) H-SS2D	1.764M	0.318	0.8948	0.9067	0.9695	0.9567
♣ H-vmunet	True	(<i>P</i>) H-SS2D	1.764M	0.318	0.9181	0.9254	0.9569	0.9467

4.2 Ablation experiments

Vision Mamba with different levels of parallelism. In order to verify the validity of the proposed method of Vision Mamba (VM) with different parallelism, we performed a series of ablation experiments. As shown in Figure 4, we performed 3 different Settings. Setting 1 is a conventional connection of VM, Setting 2 is a connection using parallel connection of two VM with half the number of channels, and Setting 3 is a connection using parallel connection of four VM each with $C/4$ the number of channels. By analyzing the parameters of Mamba in section 3.2 of this study, assuming that the parameter of this module is x for Setting 1 of the traditional VM connection method, Setting 2 can be calculated with a parameter of $0.502x$ and Setting 3 with a parameter of $0.252x$. Table 2 shows the results of this ablation experiment, and it should be noted that the parameters here refer to the parameters of the overall model (which contains the Conv Block and the skip-connection part). The parameters of Setting 2 and Setting 3 are 51.47% and 36.03%, respectively, of the parameters of Setting 1 for the traditional VM connection method, while the GFLOPs as a whole do not change much. In terms of performance, the lowest parameter of Setting 3 still maintains better segmentation performance, so in this paper, we adopt Setting 3 as the key structure of the proposed parallel Vision Mamba Layer (PVM Layer).

Parallelization of different SSM variants. In the Methods section, we detail the key to the influence of the parameters of Mamba, represented by SSM. However, many current studies propose improvements based on SSM to adapt it to 2D image processing. In Liu et al. [13], researchers proposed 2D Selective Scan (SS2D) for visual image processing. In Runa et al. [23], researchers proposed VM-UNet for medical image segmentation by combining SS2D with UNet framework. In addition, based on SS2D, Wu et al. [33] proposed High-order SS2D (H-SS2D) for medical image segmentation. The parameter and performance effects of SS2D and H-SS2D using parallel approach are shown in Table 3. From the table, it is concluded that the parameters and GFLOPs of (*P*) SS2D are reduced by 81.35% and 61.97%, respectively, while those of (*P*) H-SS2D are reduced by 80.34% and 57.14%, respectively. The above results surface that the parallel approach is effective in reducing parameters and GFLOPs not only for Mamba represented by SSM, but also for different variants of SSM at the same time.

Table 4: **Comparative experiments where PVM Layer directly replaces modules from different models**, with results including parameters, computational complexity, and performance. *Italics* in the method column (e.g., *_Conv*) indicate that the PVM Layer replaces the convolution module. In particular, we replace modules containing Convolution, Vision Transformer, Mamba, and Vision Mamba, covering the four most commonly used base modules. D^P and D^G denote the percentage decrease in parameters and GFLOPs, respectively, after replacing with PVM Layer.

Methods	Parameters and Computational Complexity				Performance Evaluation			
	Params	D^P	GFLOPs	D^G	DSC	SE	SP	ACC
UNet [22]	2.009M		3.224		0.8989	0.8793	0.9812	0.9613
UNet_Conv	0.394M	80.38%	0.686	78.74%	0.8974	0.8667	0.9842	0.9612
Att UNet [19]	3.581M		8.575		0.8821	0.8423	0.9836	0.9560
Att UNet_Conv_block	1.966M	45.10%	6.036	29.61%	0.8857	0.8414	0.9857	0.9575
SCR-Net [31]	0.812M		1.567		0.8898	0.8497	0.9853	0.9588
SCR-Net_Conv	0.147M	81.90%	0.427	72.75%	0.9058	0.8847	0.9833	0.9640
Swin-UNet [2]	27.176M		7.724		0.8670	0.8427	0.9754	0.9494
Swin-UNet_Vision Transformer	6.491M	76.11%	1.894	75.48%	0.8734	0.8445	0.9783	0.9521
META-UNet [30]	22.209M		5.14		0.9068	0.8801	0.9836	0.9639
META-UNet_ResNet Layer	1.161M	94.77%	0.416	91.91%	0.9047	0.8915	0.9807	0.9633
MHorUNet [32]	9.585M		0.864		0.9132	0.8974	0.9834	0.9666
MHorUNet_Horblock	0.925M	90.35%	0.156	81.94%	0.9107	0.8806	0.9870	0.9662
HSH-UNet [34]	18.803M		9.362		0.9108	0.8907	0.9864	0.9654
HSH-UNet_HSHB	2.888M	84.64%	0.882	90.58%	0.9001	0.8938	0.9776	0.9612
MALUNet [24]	0.175M		0.083		0.8896	0.8824	0.9762	0.9583
MALUNet_DGA	0.128M	26.86%	0.074	10.84%	0.9025	0.8623	0.9882	0.9635
EGE-UNet [25]	0.053M		0.072		0.9073	0.8931	0.9816	0.9642
EGE-UNet_GHPA	0.052M	1.89%	0.071	1.39%	0.9092	0.8941	0.9823	0.9650
VM-UNet [23]	27.427M		4.112		0.9070	0.8837	0.9842	0.9645
VM-UNet_Vision Mamba	10.713M	60.94%	1.805	56.1%	0.8885	0.8517	0.9841	0.9582
VM-UNet v2 [36]	22.771M		4.400		0.9045	0.8768	0.9849	0.9637
VM-UNet v2_Vision Mamba	5.991M	73.69%	1.521	65.43%	0.8916	0.8692	0.9804	0.9586
LightM-UNet [12]	0.403M		0.391		0.9080	0.8839	0.9846	0.9649
LightM-UNet_Mamba Layer	0.109M	72.95%	0.390	0.26%	0.9045	0.8668	0.9878	0.9642

Plug-and-play PVM Layer. The proposed PVM Layer can simply replace the base building blocks of any model, which include but are not limited to Convolution, Vision Transformers, Mamba, Vision Mamba and so on. Table 4 shows the powerful plug-and-play capabilities of the PVM Layer. From the table, it can be concluded that after replacing the base building blocks or key functional modules of any model with PVM Layer, there is a significant decrease in parameters and GFLOPs, and the performance is still clearly competitive. In addition, MALUNet and EGE-UNet, as the most advanced lightweight models, after replacing the key modules in PVM Layer, the parameters and GFLOPs can still be effectively reduced, and the performance is improved. With the above results, it is shown that the powerful plug-and-play feature of PVM Layer is used to significantly reduce the parameters and GFLOPs of arbitrary models.

5 Conclusion

In this study, we deeply analyze the key factors of parameter influence in Mamba, and based on this, we propose a parallel Vision Mamba Layer (PVM Layer) for processing the deep features. The PVM Layer uses four Vision Mamba (VM) in parallel for processing the features, and the number of channels processed by each VM is one-fourth of the initial number of channels. This is due to the fact that the number of input channels to the SSM in Mamba has an explosive effect on the number of parameters, and the VM parameters for processing a quarter of the number of channels are 6.3% of the original VM parameters, which is an explosive reduction of 93.7%. In addition, the PVM Layer can be generalized to any Mamba variant (PxM Layer). Based on the PVM Layer, we propose the UltraLight Vision Mamba UNet (UltraLight VM-UNet) with a parameter of only 0.049M and GFLOPs of only 0.060. The UltraLight VM-UNet parameters are 99.82% lower than those of the traditionally pure Vision Mamba UNet model (VM-UNet) and 87.84% lower than those of the lightest Vision Mamba UNet model available (LightM-UNet). In addition, we experimentally demonstrated on three publicly available skin lesion datasets that the UltraLight VM-UNet has equally strong performance competitiveness with such low parameters. Based on this paper, in the future, Mamba may become a new mainstream module for lightweight modeling.

References

- [1] Ehsan Khodapanah Aghdam, Reza Azad, Maral Zarvani, and Dorit Merhof. Attention swin u-net: Cross-contextual attention mechanism for skin lesion segmentation. In *2023 IEEE 20th International Symposium on Biomedical Imaging (ISBI)*, pages 1–5. IEEE, 2023. 1, 3, 7
- [2] Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning Wang. Swin-unet: Unet-like pure transformer for medical image segmentation. In *European conference on computer vision*, pages 205–218. Springer, 2022. 3, 9
- [3] Noel Codella, Veronica Rotemberg, Philipp Tschandl, M Emre Celebi, Stephen Dusza, David Gutman, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, et al. Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1902.03368*, 2019. 14
- [4] Noel CF Codella, David Gutman, M Emre Celebi, Brian Helba, Michael A Marchetti, Stephen W Dusza, Aadi Kalloo, Konstantinos Liopyris, Nabin Mishra, Harald Kittler, et al. Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi), hosted by the international skin imaging collaboration (isic). In *2018 IEEE 15th international symposium on biomedical imaging (ISBI 2018)*, pages 168–172. IEEE, 2018. 14
- [5] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11963–11975, 2022. 1
- [6] Stefan Elfving, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural networks*, 107:3–11, 2018. 12
- [7] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023. 2, 12
- [8] Ali Hatamizadeh, Vishwesh Nath, Yucheng Tang, Dong Yang, Holger R Roth, and Daguang Xu. Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images. In *International MICCAI Brainlesion Workshop*, pages 272–284. Springer, 2021. 1
- [9] Ali Hatamizadeh, Yucheng Tang, Vishwesh Nath, Dong Yang, Andriy Myronenko, Bennett Landman, Holger R Roth, and Daguang Xu. Unetr: Transformers for 3d medical image segmentation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 574–584, 2022. 1
- [10] Shishuai Hu, Zehui Liao, and Yong Xia. Devil is in channels: Contrastive single domain generalization for medical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 14–23. Springer, 2023. 7
- [11] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr, 2015. 2
- [12] Weibin Liao, Yinghao Zhu, Xinyuan Wang, Cehngwei Pan, Yasha Wang, and Liantao Ma. Lightm-unet: Mamba assists in lightweight unet for medical image segmentation. *arXiv preprint arXiv:2403.05246*, 2024. 3, 4, 6, 7, 9
- [13] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and Yunfan Liu. Vmamba: Visual state space model. *arXiv preprint arXiv:2401.10166*, 2024. 2, 8, 12
- [14] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986, 2022. 1
- [15] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 3

- [16] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 14
- [17] Jun Ma, Feifei Li, and Bo Wang. U-mamba: Enhancing long-range dependency for biomedical image segmentation. *arXiv preprint arXiv:2401.04722*, 2024. 3
- [18] Teresa Mendonça, Pedro M Ferreira, Jorge S Marques, André RS Marcal, and Jorge Rozeira. Ph 2-a dermoscopic image database for research and benchmarking. In *2013 35th annual international conference of the IEEE engineering in medicine and biology society (EMBC)*, pages 5437–5440. IEEE, 2013. 14
- [19] Ozan Oktay, Jo Schlemper, Loic Le Folgoc, Matthew Lee, Mattias Heinrich, Kazunari Misawa, Kensaku Mori, Steven McDonagh, Nils Y Hammerla, Bernhard Kainz, et al. Attention u-net: Learning where to look for the pancreas. *arXiv preprint arXiv:1804.03999*, 2018. 9
- [20] Yongming Rao, Wenliang Zhao, Yansong Tang, Jie Zhou, Ser Nam Lim, and Jiwen Lu. Hornet: Efficient high-order spatial interactions with recursive gated convolutions. *Advances in Neural Information Processing Systems*, 35:10353–10366, 2022. 1
- [21] Yongming Rao, Wenliang Zhao, Zheng Zhu, Jiwen Lu, and Jie Zhou. Global filter networks for image classification. *Advances in neural information processing systems*, 34:980–993, 2021. 1
- [22] Olaf Ronneberger, Philipp Fischer, and Thomas 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*, pages 234–241. Springer, 2015. 3, 7, 9
- [23] Jiacheng Ruan and Suncheng Xiang. Vm-unet: Vision mamba unet for medical image segmentation. *arXiv preprint arXiv:2402.02491*, 2024. 7, 8, 9
- [24] Jiacheng Ruan, Suncheng Xiang, Mingye Xie, Ting Liu, and Yuzhuo Fu. Malunet: A multi-attention and light-weight unet for skin lesion segmentation. In *2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 1150–1156. IEEE, 2022. 4, 6, 7, 9, 15
- [25] Jiacheng Ruan, Mingye Xie, Jingsheng Gao, Ting Liu, and Yuzhuo Fu. Ege-unet: an efficient group enhanced unet for skin lesion segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 481–490. Springer, 2023. 2, 7, 9
- [26] RL Siegel, KD Miller, HE Fuchs, and A Jemal. Cancer statistics, 2022. *CA: a Cancer Journal for Clinicians*, 72(1):7–33, 2022. 3
- [27] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017. 2
- [28] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015. 2
- [29] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. 2
- [30] Huisi Wu, Zebin Zhao, and Zhaoze Wang. Meta-unet: Multi-scale efficient transformer attention unet for fast and high-accuracy polyp segmentation. *IEEE Transactions on Automation Science and Engineering*, 2023. 9
- [31] Huisi Wu, Jiafu Zhong, Wei Wang, Zhenkun Wen, and Jing Qin. Precise yet efficient semantic calibration and refinement in convnets for real-time polyp segmentation from colonoscopy videos. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 2916–2924, 2021. 7, 9

- [32] Renkai Wu, Pengchen Liang, Xuan Huang, Liu Shi, Yuandong Gu, Haiqin Zhu, and Qing Chang. Mhorunet: High-order spatial interaction unet for skin lesion segmentation. *Biomedical Signal Processing and Control*, 88:105517, 2024. 1, 3, 6, 9
- [33] Renkai Wu, Yinghao Liu, Pengchen Liang, and Qing Chang. H-vmunet: High-order vision mamba unet for medical image segmentation. *arXiv preprint arXiv:2403.13642*, 2024. 6, 8
- [34] Renkai Wu, Hongli Lv, Pengchen Liang, Xiaoxu Cui, Qing Chang, and Xuan Huang. Hsh-unet: Hybrid selective high order interactive u-shaped model for automated skin lesion segmentation. *Computers in Biology and Medicine*, 168:107798, 2024. 1, 3, 9
- [35] Hanguang Xiao, Li Li, Qiyuan Liu, Xiuhong Zhu, and Qihang Zhang. Transformers in medical image segmentation: A review. *Biomedical Signal Processing and Control*, 84:104791, 2023. 3
- [36] Mingya Zhang, Yue Yu, Limei Gu, Tingsheng Lin, and Xianping Tao. Vm-unet-v2 rethinking vision mamba unet for medical image segmentation. *arXiv preprint arXiv:2403.09157*, 2024. 7, 9
- [37] Lianghai Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision mamba: Efficient visual representation learning with bidirectional state space model. *arXiv preprint arXiv:2401.09417*, 2024. 2

Appendix

A Mamba variant (SS2D) parameter impact analysis

2D Selective Scan (SS2D) have been developed based on SSM which are more suitable for visual tasks, and SS2D are usually embedded in a Visual State Space (VSS) Block [13] as shown in Figure 5. The VSS Block consists of two main branches, the first one mainly consists of a linear layer and SiLU activation function [6]. The second branch is mainly composed of linear layers, convolution, SiLU activation function, 2D Selective Scan (SS2D) and LayerNorm. Finally, the two branches merge the outputs by element-by-element multiplication.

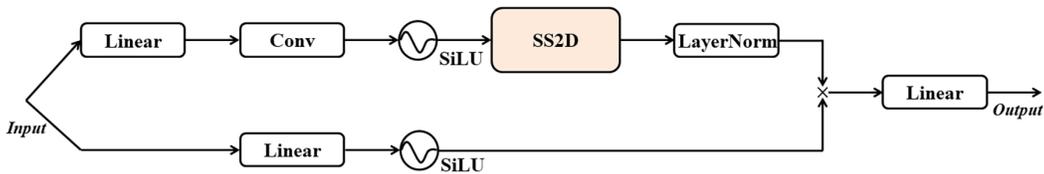


Figure 5: Visual State Space (VSS) Block composition structure.

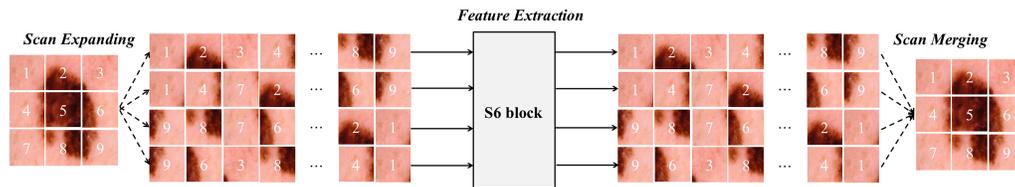


Figure 6: Image interpretation of 2D Selective Scan.

The components of SS2D are shown in Figure 6. They include scan expansion operation, S6 block feature extraction and scan merge operation. The sequence is first expanded in four directions from top-left to bottom-right, bottom-right to top-left, top-right to bottom-left and bottom-left to top-right by scan expansion operation. Then it is input to the S6 block [7] for feature extraction. Finally, it is merged back to the size of the original initial image by scan merge operation.

Among SS2D-based Mamba (VSS Block), the number of input channels, the size of the state dimension of the S6 block, the size of the internal convolution kernel, the projection dilation multiplier, and the rank of the projection matrix all affect the parameters. And among them, the influence of the number of input channels is explosive, and its influence mainly comes from the following aspects:

First, the d_{inner} of the VSS Block internal extended projection channel is determined by the product of the projection expansion multiplier and the number of input channels. This can be specifically expressed by the following equation:

$$d_{inner} = expand * d_{model} \quad (A.1)$$

where d_{inner} is the internal expansion projection channel, $expand$ is the projection expansion multiplier (fixed at 2 by default), and d_{model} is the number of input channels. We can see that d_{inner} will rise exponentially as the number of channels per layer in the model increases dramatically.

Second, the parameters of the input projection layer (the same input linear layer is used for both branches) and output projection layer within VSS Block will be directly related to the number of input channels. The input projection layer and output projection layer operate as follows:

$$in_proj : nn.Linear(d_{model}, d_{inner} * 2, bias = False) \quad (A.2)$$

$$out_proj : nn.Linear(d_{inner}, d_{model}, bias = False) \quad (A.3)$$

where then the input projection (in_proj) layer parameter is ($d_{model} * d_{inner} * 2$), and the output projection layer (out_proj) parameter is ($d_{inner} * d_{model}$). In addition, the output section has a layer normalization operation (LayerNorm) with parameter ($d_{inner} + d_{inner}$). We can see that the number of input channels, d_{model} , is the key element controlling the parameter, where the internal extended projection channel, d_{inner} , is also controlled by d_{model} .

Further, the linear projection layers in the S6 block of SS2D are also key to the parameter effects. Each linear projection layer is specified as follows:

$$x_proj = nn.Linear(d_{inner}, (dt_rank + d_{state} * 2), bias = False) \quad (A.4)$$

$$dt_proj = nn.Linear(dt_rank, d_{inner}, bias = True) \quad (A.5)$$

where dt_rank is the rank of the projection matrix ($dt_rank = d_{model}/16$), d_{state} is the size of the S6 block state dimension (fixed to 16), and the parameters for each linear projection layer are ($d_{inner} * (dt_rank + d_{state} * 2)$). However, there are 4 linear projection layers in total, so the total parameters are $4 * (d_{inner} * (dt_rank + d_{state} * 2))$. In addition, dt_proj is a linear projection layer for step size with parameters $dt_rank * d_{inner} + d_{inner}$, which is mainly used for linear projection for step size (dt). dt_proj also has 4 layers, with a total parameter of $4 * (dt_rank * d_{inner} + d_{inner})$. So, from the above, we can know that all parameters are still mainly controlled by the number of input channels d_{model} .

In addition, the internal convolution ($nn.Conv2d(d_{inner}, d_{inner}, d_{conv}, bias = True)$) also provides parametric influence with $d_{inner} * d_{inner} * d_{conv} * d_{conv} + d_{inner}$. In this paper, d_{conv} is fixed to 3, so the convolution provides a parameter of $3^2 * d_{inner}^2 + d_{inner}$, which is also controlled by the d_{model} .

Also, the A_logs of the parameter matrix controlling the attention weights of the different states of the S6 block in the SS2D module is an important influencing element. A_logs is a parameter matrix of the shape ($K * d_{inner}, d_{state}$), and K is a hyperparameter which is usually fixed to 4. Therefore, the parameter A_logs can be derived as $d_{inner} * d_{state} * 4$. In addition, the trainable vector parameter (Ds) of the process computed within the SS2D contains the $4 * d_{inner}$ parameter, which is used to selectively integrate the SS2D state outputs with the original input signals, thereby enhancing the model's expressiveness and training stability.

In summary, assuming that the original number of input channels is 1024, keeping the other parameters unchanged, and reducing the number of channels to a quarter of the original (the number of input channels becomes 256), the original total parameters can be calculated by the above parameter formulae from 45,504,512 to 2,921,984. The parameter explosion reduces the number of channels by 93.6%, which further confirms that the number of input channels has a very critical impact on the VSS Block parameters.

B Experimental settings

B.1 Datasets

To validate that the proposed UltraLight VM-UNet also achieves competitive performance at the parameter of 0.049M, we conducted experiments on three publicly available dermatologic lesion datasets. ISIC2017 [4] and ISIC2018 [3] datasets are two large datasets published by the International Skin Imaging Collaboration (ISIC), respectively. The PH² [18] dataset is a small public dataset of skin lesions, so we used PH² as an external validation to train the weights using the ISIC2017 dataset.

For the ISIC2017 dataset we acquired 2000 images as well as dermatoscope images with segmentation mask labels. Among them, the dataset was randomly divided, 1250 were used for model training, 150 images were used for model validation, and 600 images were used for model testing. The initial size of the images is 576×767 pixels, and we standardize the size to 256×256 pixels when inputting the model.

For the ISIC2018 dataset we acquired 2594 images as well as dermatoscope images with segmentation mask labels. Among them, the dataset was randomly divided, 1815 were used for model training, 259 images were used for model validation, and 520 images were used for model testing. The initial size of the images is 2016×3024 pixels, and we standardize the size to 256×256 pixels when inputting the model.

For the PH² dataset we acquired 200 images as well as dermatoscope images with segmentation mask labels. All 200 images were used for external validation. The initial size of the images was 768×560 pixels and we standardized the size to 256×256 pixels for inputting into the model.

B.2 Implementation details

The experiments were all implemented based on Python 3.8 and Pytorch 1.13.0. A single NVIDIA V100 GPU with 32GB of memory was used for all experiments. All experiments used the same data augmentation operations to more fairly determine the performance of the model, including horizontal and vertical flips, and random rotation operations. BceDice loss function was used, with AdamW [16] as the optimizer, a training epoch of 250, a batch size of 8, and a cosine annealing learning rate scheduler with an initial learning rate of 0.001 and a minimum learning rate set to 0.00001.

B.3 Evaluation metrics

Dice similarity coefficient (DSC), sensitivity (SE), specificity (SP) and accuracy (ACC) are the most commonly used evaluation metrics for medical image segmentation. DSC is used to measure the degree of similarity between the ground truth and the predicted segmentation map. SE is mainly used to measure the percentage of true positives in true positives and false negatives. SP is mainly used to measure the percentage of true negatives in true negatives and false positives. ACC is mainly used to measure the percentage of correct classification.

$$DSC = \frac{2TP}{2TP + FP + FN} \tag{B.1}$$

$$ACC = \frac{TP + TN}{TP + TN + FP + FN} \tag{B.2}$$

$$Sensitivity = \frac{TP}{TP + FN} \tag{B.3}$$

$$Specificity = \frac{TN}{TN + FP} \tag{B.4}$$

where TP denotes true positive, TN denotes true negative, FP denotes false positive and FN denotes false negative.

C Visualization results

To more directly represent the competitive nature of UltraLight VM-UNet in terms of segmentation performance, we visualized the segmentation results (Figure 7). In addition, we have also visualized the segmentation results of several state-of-the-art lightweight and classical medical image

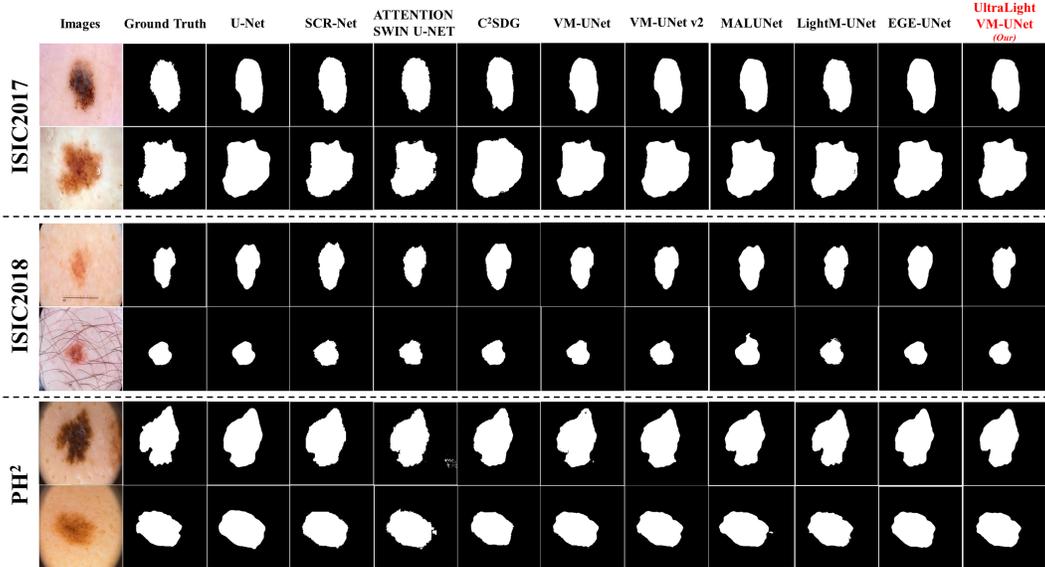


Figure 7: Visualization of segmentation graphs for comparison experiments of three publicly available skin lesion segmentation datasets.

segmentation models. From the visualizations, it can be concluded that the segmentation results of UltraLight VM-UNet have smooth, clear and more accurate boundaries. This shows that UltraLight VM-UNet not only outperforms the rest of the current lightweight models in terms of parameters and computational complexity, but also remains competitive in terms of performance.

Table 5: Ablation experiments on the effect of each module in the UltraLight VM-UNet.

Methods	Params	GFLOPs	DSC	SE	SP	ACC
UltraLight VM-UNet (Baseline)	0.049M	0.060	0.9091	0.9053	0.9790	0.9646
Baseline_Encoder_Conv	0.080M	0.071	0.9033	0.8643	0.9880	0.9638
Baseline_Decoder_Conv	0.080M	0.064	0.8958	0.8512	0.9881	0.9612
Baseline_(En+De)_Conv	0.123M	0.075	0.9065	0.8784	0.9855	0.9645
Baseline_SCAB not applicable	0.033M	0.058	0.9029	0.8767	0.9841	0.9631

D More analysis

D.1 Impact of components in the UltraLight VM-UNet

A series of ablation experiments were performed to verify the impact of each module in the UltraLight VM-UNet. As shown in Table 5, we replace the PVM Layer of the encoder, decoder, respectively, with a standard convolution with convolution kernel 3. In addition, we also replace the PVM Layer of the encoder and decoder at the same time with a standard convolution. Also, *Baseline_SCAB not applicable* indicates that the skip-connection part of the UltraLight VM-UNet does not use the SAB and CAB modules. From the table, we can conclude that in replacing the PVM Layer of the encoder and decoder separately, the parameters increased by 63.26% and the GFLOPs increased in both, while the performance decreased in both. In particular, after replacing the PVM Layer of the encoder and decoder simultaneously, the parameters increase by 151% and the GFLOPs are increased by 25%. In summary, it is shown that after replacing the PVM Layer there is a decrease in all performance aspects and an increase in both parameters and GFLOPs. This again proves the crucial role of PVM Layer. What’s more, although the parameters and GFLOPs were further reduced with the removal of the SAB and CAB modules. However, the performance also exhibits some degradation, which is due to the ability of the SAB and CAB modules to combine multi-scale feature information for learning, which improves the sensitivity to lesions and accelerates the model convergence [24]. Therefore, in order to balance the performance and computational com-

plexity relationship, UltraLight VM-UNet employs the SAB and CAB modules as a bridge for skip connections to further improve segmentation performance. Even so, the performance, parameters and GFLOPs of UltraLight VM-UNet are still in the forefront compared to the rest of the state-of-the-art lightweight models available currently.

Table 6: Ablation experiments on the effect of combining different channel numbers in the UltraLight VM-UNet. D^P denotes the percentage of parameter reduction with the parallel approach. 74.80% denotes the theoretical percentage reduction of localized module parameters derived from the Mamba analysis in the methods section. ✖ indicates no parallel connection method and ★ indicates a quadruple parallel connection method.

Methods	Params	D^P (74.80%)	GFLOPs	DSC	SE	SP	ACC
✖ [8, 16, 24, 32, 48, 64]	0.136M	63.98% (-10.82%)	0.060	0.9069	0.8861	0.9834	0.9644
★ [8, 16, 24, 32, 48, 64]	0.049M		0.060	0.9091	0.9053	0.9790	0.9646
✖ [8, 16, 32, 64, 128, 256]	0.909M	75.47% (+0.67%)	0.074	0.9018	0.8749	0.9840	0.9627
★ [8, 16, 32, 64, 128, 256]	0.223M		0.074	0.9079	0.8965	0.9809	0.9644
✖ [16, 32, 64, 128, 256, 512]	3.479M	75.63% (+0.83%)	0.259	0.9049	0.8981	0.9789	0.9630
★ [16, 32, 64, 128, 256, 512]	0.848M		0.259	0.9056	0.8906	0.9815	0.9637
✖ [32, 64, 128, 256, 512, 1024]	13.607M	75.71% (+0.91%)	0.970	0.9053	0.8878	0.9821	0.9637
★ [32, 64, 128, 256, 512, 1024]	3.305M		0.970	0.9012	0.8812	0.9819	0.9622

D.2 Impact of different channel numbers in UltraLight VM-UNet

In addition, in order to verify the impact of different channel numbers in the UltraLight VM-UNet, we conducted a series of ablation experiments. As shown in Table 6, we set up four different combinations of channel numbers, including [8, 16, 24, 32, 48, 64], [8, 16, 32, 64, 128, 256], [16, 32, 64, 128, 256, 512] and [32, 64, 128, 256, 512, 1024], respectively. In addition, we conducted experiments both in the parallel-free manner (labeled by ✖) and in the parallel manner (labeled by ★). From the table, it can be concluded that the parameters and GFLOPs at [8, 16, 24, 32, 48, 64] are the lowest, while the overall difference in performance is not very large. However, as the number of channels increases, both parameters and GFLOPs show a significant rise. The parameter increases from 0.136M to 13.607M for the group without parallelism and from 0.049M to 3.305M for the group with parallelism. In particular, the average decrease in parameters for the same number of channel combinations using the parallel approach over the no-parallel approach is 72.70%, which is extremely close to the percentage of decrease in the four-parallel approach over the no-parallel approach analyzed in the Methods section (74.80%). In addition, at [8, 16, 24, 32, 48, 64], the parameters are more susceptible to the rest of the modules, so the decrease is not as large as for the rest of the channel number combinations. However, its overall parameters and GFLOPs reach an impressive 0.049M and 0.060, and show strong competitive segmentation performance.