# Vivim: a Video Vision Mamba for Medical Video Segmentation

Yijun Yang, Zhaohu Xing, Lequan Yu, *Member, IEEE,* Chunwang Huang, Huazhu Fu, *Senior Member, IEEE,* Lei Zhu, *Member, IEEE*

*Abstract*—Medical video segmentation gains increasing attention in clinical practice due to the redundant dynamic references in video frames. However, traditional convolutional neural networks have a limited receptive field and transformer-based networks are mediocre in constructing long-term dependency from the perspective of computational complexity. This bottleneck poses a significant challenge when processing longer sequences in medical video analysis tasks using available devices with limited memory. Recently, state space models (SSMs), famous by Mamba, have exhibited impressive achievements in efficient long sequence modeling, which develops deep neural networks by expanding the receptive field on many vision tasks significantly. Unfortunately, vanilla SSMs failed to simultaneously capture causal temporal cues and preserve non-casual spatial information. To this end, this paper presents a Video Vision Mamba-based framework, dubbed as Vivim, for medical video segmentation tasks. Our Vivim can effectively compress the long-term spatiotemporal representation into sequences at varying scales with our designed Temporal Mamba Block. We also introduce an improved boundary-aware affine constraint across frames to enhance the discriminative ability of Vivim on ambiguous lesions. Extensive experiments on thyroid segmentation, breast lesion segmentation in ultrasound videos, and polyp segmentation in colonoscopy videos demonstrate the effectiveness and efficiency of our Vivim, superior to existing methods. The code is available at: https://github.com/scott-yjyang/Vivim. The dataset will be released once accepted.

*Index Terms*—Thyroid segmentation, Breast lesion segmentation, polyp segmentation, State space model, Ultrasound videos.

## I. INTRODUCTION

Automatic segmentation of lesions and tissues is essential for computer-aided clinical examination and treatment [1], such as ultrasound lesion segmentation, polyp segmentation. However, segmenting medical objects is usually challenging due to inherent factors, including ambiguous lesion boundaries, inhomogeneous distributions, diverse motion patterns,

Yijun Yang and Zhaohu Xing are with Robotics and Autonomous Systems Thrust, The Hong Kong University of Science and Technology, Guangzhou, China (e-mail: yyang018@connect.hkust-gz.edu.cn; zxing565@connect.hkust-gz.edu.cn).

Lequan Yu is with the Department of Statistics and Actuarial Science, The University of Hong Kong, Hong Kong SAR, China (e-mail: lqyu@hku.hk).

Chunwang Huang is with Guangdong Provincial People's Hospital, Guangzhou, China (e-mail: huangchunwang@126.com).

Huazhu Fu is with the Institute of High Performance Computing, A*STAR, Singapore (e-mail: hzfu@ieee.org).

Lei Zhu is with Robotics and Autonomous Systems Thrust, Hong Kong University of Science and Technology (Guangzhou), China, and the Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology, Hong Kong SAR, China (e-mail: leizhu@ust.hk).

Lei Zhu is the corresponding author of this work.

and dynamic changes in complex environments [2]. Medical videos, essentially sequences of medical images, offer a richer and more detailed context for locating ambiguous lesions and tissues. This additional information makes video-based segmentation better handle unexpected complexities by providing a continuous view, allowing for a more accurate and comprehensive analysis. Consequently, to consider more object context, expanding the deep model's receptive field in the spatiotemporal space is highly desired in medical video analysis. Traditional convolutional neural networks [3]–[6] often struggle to capture global information compared to recent transformer-based architectures. The transformer architecture, which utilizes the Multi-Head Self Attention (MSA) [7] to extract global information, has attracted much attention from the community of generic video object segmentation [8]–[10]. Considering that neighboring frames offer beneficial hints to the segmentation, these methods usually introduce some elaborated modules on the self-attention mechanism to exploit the temporal information. However, exploring the additional temporal dimension often leads to increased complexity and greater demands on resources, posing significant challenges for implementation due to the strict environmental conditions and inherently high-dimensional characteristics of medical videos. For example, the incorporation of temporal self-attention modules can unintentionally trigger a quadratic increase in complexity relative to the time dimension, resulting in substantial computational challenges. The marked rise in the number of tokens within lengthy video sequences introduces considerable computational strains when employing Multi-head Self-Attention (MSA) techniques for temporal information modeling [9].

Very recently, to address the ill-posed issue concerning long sequence modeling, Mamba [11], inspired by state space models (SSMs) [12], has been developed. Its main idea is to efficiently capture long-range dependencies by implementing a selective scan mechanism for 1-D sequence interaction. Based on this, U-Mamba [13] designed a hybrid CNN-SSM block, which is mainly composed of Mamba modules, to handle the long sequences in biomedical image segmentation tasks. Vision Mamba [14] provided a new generic vision backbone with bidirectional Mamba blocks on image classification and semantic segmentation tasks. As suggested by them, relying on the self-attention module is not necessary to achieve efficient visual representation learning. It can be replaced by Mamba when exploring long-term temporal dependency in video scenarios. The crucial aspect of adapting the Vision Mamba model for video applications lies in the ability to
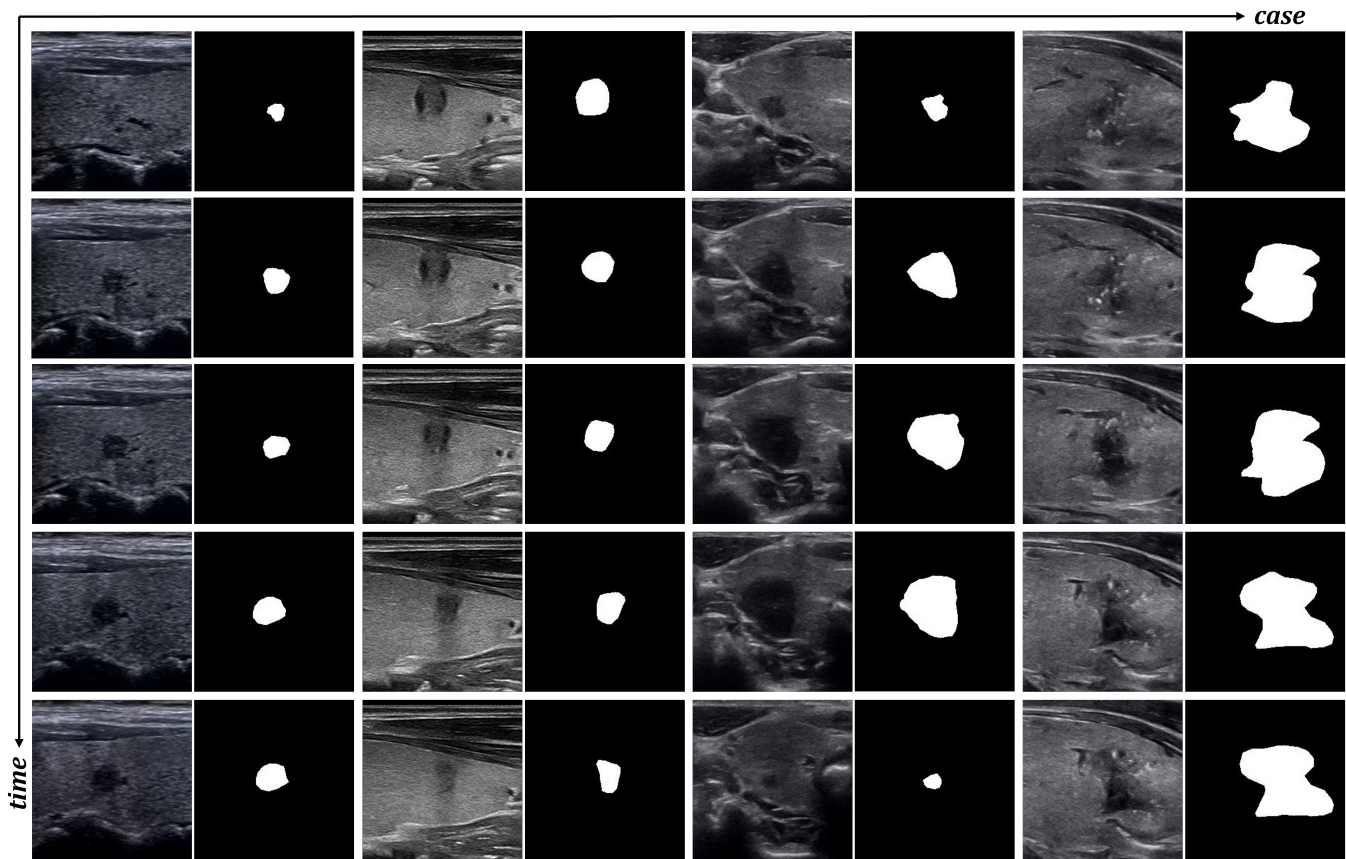
Fig. 1. **Several cases of our collected VTUS dataset.** All videos are taken from patients with thyroid nodules. They are taken by ultrasound doctors with more than 10 years of clinical experience to ensure the image quality. These videos are cross-annotated by three experts with over three years of experience in thyroid diagnosis.

concurrently capture causal temporal cues while maintaining the integrity of non-causal spatial information.

Motivated by this, we present an SSMs-based framework Vivim that integrates Mamba into the multi-level transformer architecture to exploit spatiotemporal information in videos with linear complexity. *To the best of our knowledge, this is the first work to incorporate SSMs into the task of medical video segmentation, facilitating faster and greater performance.* In our Vivim, drawing inspiration from the architecture of modern transformer blocks, we design a novel Temporal Mamba Block. A hierarchical encoder consisting of multiple Temporal Mamba Blocks is introduced to investigate the correlation between spatial and temporal dependency at various scales. As the structured state space sequence models with selective scan (S6) [11] causally process input data, they can only capture information within the scanned portion of the data. This aligns S6 with NLP and video tasks involving temporal data but poses challenges when addressing non-causal data like 2D images within medical videos. To this end, the structured state space sequence model with spatiotemporal selective scan, ST-Mamba, is designed and incorporated into each scale of the model's encoder, replacing the self-attention or window-attention module to achieve efficient video visual representation learning. Finally, we employ an improved boundary-aware affine constraint to improve the discrimination of Vivim on

ambiguous tissues in medical videos at the training stage.

It is worth noting that *there is no public dataset with pixel-level annotated ultrasound videos for thyroid segmentation*, as it is expensive to delineate the boundaries of ambiguous lesions in low-contrast ultrasound videos in a frame-by-frame spirit. In this work, we collect a thyroid segmentation dataset VTUS with 100 annotated transverse viewed and longitudinal viewed ultrasound videos and a total of 9342 frames with pixel-level ground truth to facilitate the benchmarking evaluation. Several examples are displayed in Fig. 1. We conduct extensive experiments on three popular medical video segmentation tasks, *i.e.*, thyroid segmentation in ultrasound videos, breast lesion segmentation in ultrasound videos, and polyp segmentation in colonoscopy videos. The superior results validate the effectiveness, efficiency and versatility of our framework Vivim.

Our contributions can be summarized as follows:

- We develop a medical video segmentation framework consisting of a Mamba-based encoder and a CNN-based decoder to obtain holistic understanding of medical videos and preserve local details, respectively. This is the first work to introduce state space models into medical video scenarios.
- Instead of simply adapting Mamba to medical tasks, we design spatio-temporal selective scan to enhance the

global perception capability in videos of our Temporal Mamba Block.

- We employ an improved boundary-aware constraint based on the optimization of the affine transformation to mitigate ambiguous boundary prediction of our model.
- We collect the first video ultrasound thyroid segmentation dataset with pixel-level annotation, which facilitates the benchmarking evaluation of medical video segmentation methods. Our model achieves promising segmentation results on diverse modalities but maintains decent efficiency superior to Transformer-based methods.

## II. RELATED WORKS

### A. Medical Video Segmentation

Recent approaches have introduced innovative hybrid transformer-based algorithms that fuse transformative and convolutional layer techniques for medical image segmentation (*e.g.*, breast lesion, polyp) [3], [6], [15]–[20]. For thyroid segmentation in ultrasound images, Jeremy *et al.* [21] developed a novel spatio-temporal recurrent deep learning network to automatically segment the thyroid gland in ultrasound cineclips by leveraging time sequence information. Ma *et al.* [22] utilized the region proposal network (RPN) for initial deep feature extraction and incorporated the spatial pyramid RoIAlign as a segmentation head to capture global and local information in ultrasound images. Chi *et al.* [23] developed a 2D Transformer-UNet for thyroid gland segmentation, combining high-level features from decoding layers with lower-level features from encoding layers using a multiscale cross-attention transformer module. These algorithms skillfully manage the representations derived from high-definition medical images, however, they grapple with computational difficulties owing to complexity issues. Additionally, the direct application of such image segmentation methods may inadvertently overlook critical temporal context, thereby inducing temporal inconsistencies. In order to address temporal modeling in video-level segmentation, the innovative method of Space-Time Memory Networks (STM) [24] and its variants [10], [25], [26] are introduced, employing a memory network to extract vital information from a time-based buffer composed of all previous video sequences. Building upon this methodology, DPSTT [27] integrates a memory bank with decoupled transformers to track temporal lesion movement in medical ultrasound videos. However, DPSTT calls for substantial data augmentation to avoid overfitting and is marked by a sluggish processing speed, stressing some potential limitations. FLA-Net [28] presents a frequency and location feature aggregation network with a large amount of memory occupancy for ultrasound video breast lesion segmentation. Thus, the challenge in medical video segmentation revolves around efficiently harnessing the wealth of temporal data available.

### B. State Space Models

Recently, State Space Models (SSMs) [12] have demonstrated notable efficiency in utilizing state space transformations [29] to manage long-term dependencies within language sequences. S4 [30] introduced a structured state-space sequence model to exploit long-range dependencies with the benefit of linear complexity. Based on this, Mamba [11] integrates efficient hardware design and a selection mechanism employing parallel scan (S6), thereby surpassing Transformers in processing extensive natural language sequences. Subsequently, S4ND [31] explores SSMs' continuous-signal modeling of multi-dimensional data like images and videos. More recently, Vision Mamba [14] and Vmamba [32] pioneered generic vision tasks and outperformed transformer-based methods in effectiveness and efficiency, introducing bidirectional scan and cross-scan mechanisms to tackle the directional sensitivity challenge in SSMs. U-Mamba [13] designed a hybrid CNN-SSM block, which is mainly composed of Mamba modules, to handle the long sequences in biomedical image segmentation tasks. To the best of our knowledge, SSMs have not yet been explored in medical video segmentation tasks.

## III. METHOD

### A. Overview

In this part, we elaborate on a Mamba-based solution Vivim for medical video segmentation tasks. Our Vivim mainly consists of two modules: A hierarchical encoder with the stacked Temporal Mamba Blocks to extract coarse and fine feature sequences at different scales, and a lightweight CNN-based segmentation head to fuse multi-level feature sequences and predict segmentation masks. Fig. 2 illustrates the flowchart of our proposed Vivim. Specifically, given a video clip with $T$ frames, *i.e.*, $\mathbf{V} = \{I^1, ..., I^T\}$, we first divide these frames into patches of size $4 \times 4$ by overlapped patch embedding. We then feed the sequence of patches into our hierarchical Temporal Mamba Encoder to obtain multi-level spatiotemporal features with resolution $\{1/4, 1/8, 1/16, 1/32\}$ of the original frame. Finally, we pass multi-level features to the CNN-based segmentation head to predict the segmentation results. The Boundary-aware Affine Constraint is deployed on the results only during training as shown in Fig. 4. Please refer to the following sections for details of our proposed module.

### B. Preliminaries: State Space Models

State Space Models (SSMs) are commonly considered as linear time-invariant systems, which map a 1-D function or sequence $x(t) \in \mathbb{R} \mapsto y(t) \in \mathbb{R}$ through a hidden state $h(t) \in \mathbb{R}^{\mathbb{N}}$. This system is typically formulated as linear ordinary differential equations (ODEs), which uses $\mathbf{A} \in \mathbb{R}^{\mathbb{N} \times \mathbb{N}}$ as the evolution parameter and $\mathbf{B} \in \mathbb{R}^{\mathbb{N} \times 1}$, $\mathbf{C} \in \mathbb{R}^{1 \times \mathbb{N}}$ as the projection parameters.

$$h'(t) = \mathbf{A}h(t) + \mathbf{B}x(t), \ y(t) = \mathbf{C}h(t). \tag{1}$$

The discretization is introduced to primarily transform the ODE into a discrete function. This transformation is crucial to align the model with the sample rate of the underlying signal embodied in the input data, enabling computationally efficient operations. The structured state space sequence models (S4) and Mamba are the classical discrete versions of the continuous system, which include a timescale parameter $\mathbf{\Delta}$ to
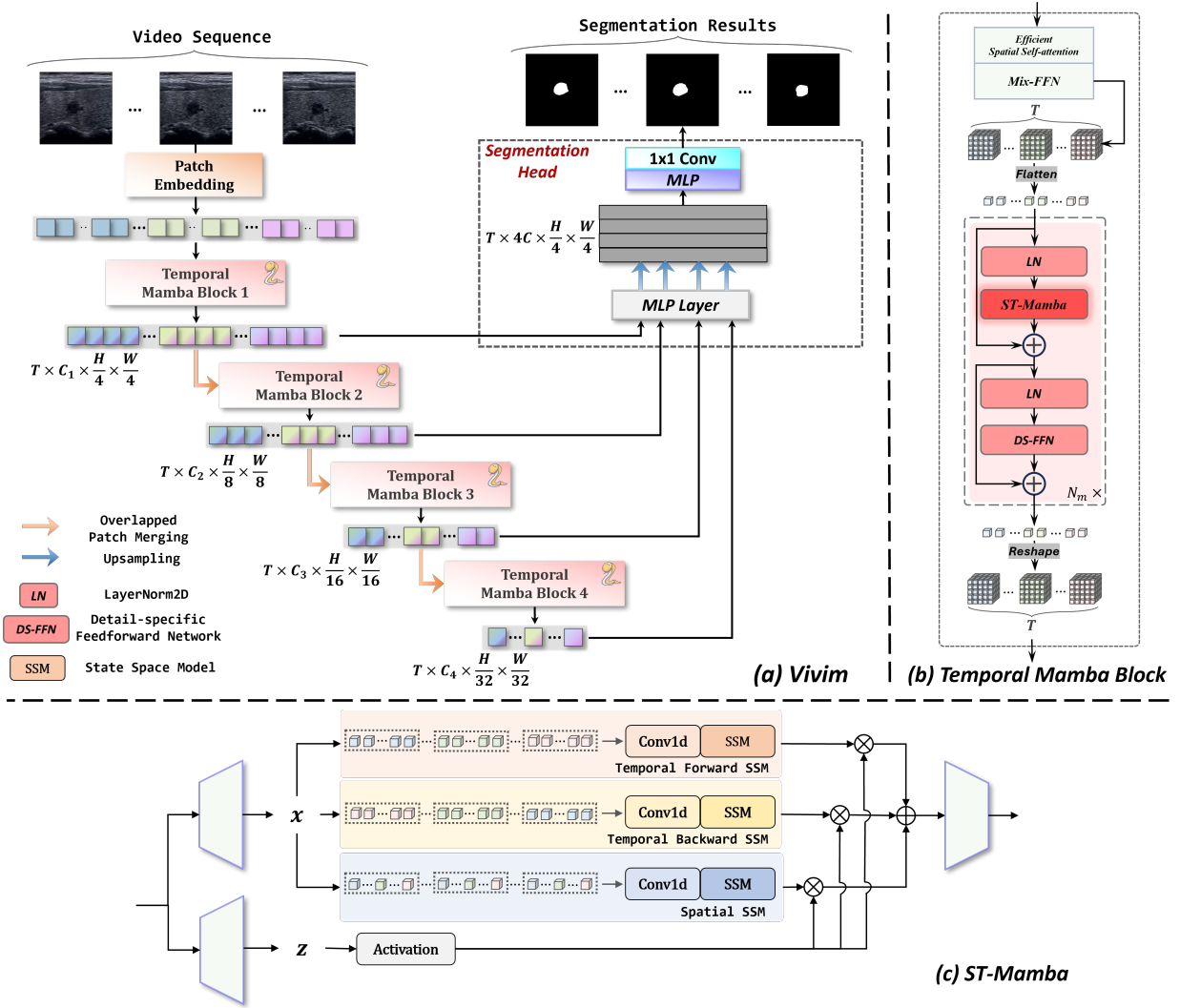
Fig. 2. (a ) **The overview of the proposed Vivim for medical video segmentation.** The video sequence is first fed into patch embedding and multi-scale Temporal Mamba Blocks for encoding. Then, the feature sequences are aggregated to predict the segmentation results by a CNN-based segmentation head. (b) The fundamental building block of Vivim, namely Temporal Mamba Block. While Efficient Spatial Self-attention conducts initial spatial modeling, ST-Mamba explores spatiotemporal dependency in a linear complexity. (c) ST-Mamba incorporates spatiotemporal selective scan for long sequence modeling of video vision tasks in a multi-way spirit.

transform the continuous parameters $\mathbf{A}$, $\mathbf{B}$ to discrete parameters $\overline{\mathbf{A}}$, $\overline{\mathbf{B}}$. The commonly used method for transformation is zero-order hold (ZOH), which is defined as follows:

$$\overline{\mathbf{A}} = \exp\left(\mathbf{\Delta A}\right), \ \overline{\mathbf{B}} = (\mathbf{\Delta A})^{-1}(\exp\left(\mathbf{\Delta A}\right) - \mathbf{I}) \cdot \mathbf{\Delta B}. \quad (2)$$

After the discretization of $\overline{\mathbf{A}}$, $\overline{\mathbf{B}}$, the discretized version of Eq. (1) can be rewritten as:

$$h_t = \overline{\mathbf{A}}h_{t-1} + \overline{\mathbf{B}}x_t, \ y_t = \mathbf{C}h_t. \quad (3)$$

At last, the models compute output through a global convolution.

$$\overline{\mathbf{K}} = (\mathbf{C}\overline{\mathbf{B}}, \mathbf{C}\overline{\mathbf{A}}\overline{\mathbf{B}}, \dots, \mathbf{C}\overline{\mathbf{A}}^{\mathtt{M}-1}\overline{\mathbf{B}}), \ \mathbf{y} = \mathbf{x} * \overline{\mathbf{K}}, \quad (4)$$

where $\mathtt{M}$ is the length of the input sequence $\mathbf{x}$, and $\overline{\mathbf{K}} \in \mathbb{R}^{\mathtt{M}}$ is a structured convolutional kernel.

## C. Overall Architecture

*1) Hierarchical feature representation:* Multi-level features provide both high-resolution coarse features and low-resolution fine-grained features that significantly improve the segmentation results, especially for medical images. To this end, unlike Vivit [9], our encoder extracts multi-level multi-scale features given input video frames. Specifically, we perform patch merging frame-by-frame at the end of each Temporal Mamba Block, resulting in the $i$-th feature embedding $\mathcal{F}_i$ with a resolution of $\frac{H}{2^{i+1}} \times \frac{W}{2^{i+1}}$.

*2) Temporal Mamba Block:* Exploring temporal information is critically important for medical video segmentation by providing dynamic appearance and motion cues. However, MSA in vanilla Transformer architectures has quadratic complexity concerning the number of tokens [7]. This complexity is pertinent for long feature sequences from videos, as the number of tokens increases linearly with the number of input frames. Motivated by this, we develop a more efficient block,
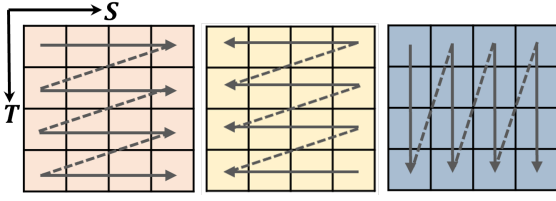
Fig. 3. The illustration of the proposed spatiotemporal selective scan, including temporal forward scan, temporal backward scan and spatial scan.

Temporal Mamba Block, to simultaneously exploit spatial and temporal information by structured state space sequence models.

As illustrated in Fig. 2 (b), in the Temporal Mamba Block, an efficient spatial-only self-attention module is first introduced to provide the initial aggregation of spatial information followed by a Mix-FeedForwoard layer. We leverage the sequence reduction process introduced in [33], [34] to improve its efficiency. For the $i$-level feature embedding $\mathcal{F}_i \in \mathbb{R}^{T \times C_i \times H \times W}$ of the given video clip, we transpose the channel and temporal dimension, and flatten the spatiotemporal feature embedding into 1D long sequence $h_i \in \mathbb{R}^{C_i \times THW}$. Then, the flattened sequence $h_i$ is fed into layers of a Spatio-Temporal Mamba module (ST-Mamba) and a Detail-Specific Feedforward (DSF). The ST-Mamba module establishes the intra- and inter-frame long-range dependencies while the DSF preserves fine-grained details by a depth-wise convolution with a kernel size of $3 \times 3 \times 3$. The procedure in the stacked Mamba Layer can be defined as, where $l \in [1, N_m]$:

$$
\begin{aligned}
h^l &= \text{ST-Mamba}\left(\text{LN}\left(h^{l-1}\right)\right) + h^{l-1}, \\
h^l &= \text{DSF}\left(\text{LN}\left(h^l\right)\right) + h^l.
\end{aligned}
\tag{5}
$$

Finally, we return the output feature sequence to the original shape and employ overlapped patch merging to down-sampling the feature embedding.

*3) Decoder:* To predict the segmentation mask from the multi-level feature embeddings, we introduce a CNN-based segmentation head. While our hierarchical Temporal Mamba encoder has a large effective receptive field across spatial and temporal axes, the CNN-based segmentation head further refines the details of local regions. To be specific, the multi-level features $\{\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3, \mathcal{F}_4\}$ from the temporal mamba blocks are passed into an MLP layer to unify the channel dimension. These unified features are up-sampled to the same resolution and concatenated together. Third, a MLP layer is adopted to fuse the concatenated features $\mathcal{F}$. Finally, The fused feature goes through a $1 \times 1$ convolutional layer to predict the segmentation mask $\mathcal{M}$. The segmentation loss $\mathcal{L}_{seg}$ consisting of pixel-wise cross-entropy loss and IoU loss is applied during training.

### D. Spatiotemporal Selective Scan

Despite the causal nature of S6 for temporal data, videos differ from texts in that they not only contain temporal redundant information but also accumulate non-causal 2D spatial information. To address this problem of adapting to
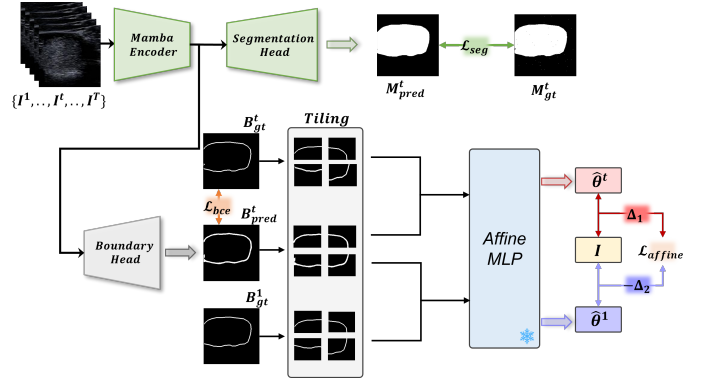


Fig. 4. **The overview of the training strategy.** Specifically, our proposed patch-level boundary-aware affine constraint $\mathcal{L}_{affine}$ is introduced to optimize Vivim jointly with the segmentation loss $\mathcal{L}_{seg}$ and the boundary cross-entropy loss $\mathcal{L}_{bce}$. The pre-trained MLP for computing the affine transformation is frozen during training.

non-causal data and fully exploring temporal information, we introduce ST-Mamba as shown in Fig. 2 (c), which incorporates spatiotemporal sequence modeling for video vision tasks.

Specifically, to explicitly explore the relationship among frames, we first unfold patches of each frame along rows and columns into sequences, and then concatenate the frame sequences to constitute the temporal-first sequence $h_i^t \in \mathbb{R}^{C_i \times T(HW)}$. We parallelly proceed with scanning along the forward and backward directions to explore bidirectional temporal dependency. This approach allows the models to compensate for each other's receptive fields without significantly increasing computational complexity. Simultaneously, we stack patches along the temporal axis and construct the spatial-first sequence $h_i^s \in \mathbb{R}^{C_i \times (HW)T}$. We proceed with scanning to integrate information of each pixel from all frames. The spatiotemporal selective scan mechanism with three directions is also vividly demonstrated in Fig. 3. Our mechanism explicitly considers both single-frame spatial coherence and cross-frame coherence, and leverages parallel SSMs to establish the intra- and inter-frame long-range dependencies. The structured state space sequence models with spatiotemporal selective scan (ST-Mamba), serve as the core element to construct the Temporal Mamba block, which constitutes the fundamental building block of Vivim.

**Computational-Efficiency.** SSMs in ST-Mamba and self-attention in Transformer both provide a crucial solution to model spatiotemporal context adaptively. Given a video visual sequence $\mathbf{K} \in \mathbb{R}^{1 \times T \times M \times D}$, the computation complexity of a global self-attention and SSM are:

$$
\Omega(\text{self-attention}) = 4(\text{TM})\text{D}^2 + 2(\text{TM})^2\text{D}, \tag{6}
$$

$$
\Omega(\text{SSM}) = 4(\text{TM})(2\text{D})\text{N} + (\text{TM})(2\text{D})\text{N}^2, \tag{7}
$$

where the default expansion ratio is 2, N is a fixed parameter and set to 16. As observed, self-attention is quadratic to the whole video sequence length (TM), and SSM is linear to that. Such computational efficiency makes ST-Mamba a better solution for long-term video applications. This is also validated by the experimental analysis on the efficiency of ST-Mamba in Sec. IV-E.

## E. Boundary-aware Affine Constraint

The network optimized only by the segmentation supervision tends to generate ambiguous and unstructured predictions, and overfit on training data. To mitigate these issues, we introduce a patch-level boundary-aware affine constraint inspired by InverseForm [35] to enforce the predicted boundary structure. Specifically, as illustrated in Fig. 4, we address this constraint task by optimizing the affine transformation between ground-truth boundaries and edges in feature maps towards identity transformation matrix. The ground truth edges within the patches are derived from applying the Sobel operator [36] on ground truth masks, while an auxiliary boundary head consisting of three convolutional layers processes the feature patches from the Mamba encoder to obtain the predicted edge. We calculate the affine transform matrix $\hat{\theta}_i^t$ for the $i$-th patch between ground-truth edge $B_{gt}^t$ and predicted edge $B_{pred}^t$ of the target frame $I^t$ in a video clip, by a pre-trained MLP. Simultaneously, we calculate another affine transform matrix $\hat{\theta}_i^1$ for the $i$-th patch between ground-truth edge $B_{gt}^1$ of $I^1$ and predicted edge $B_{pred}^t$ of $I^t$ in a video clip. This MLP is trained in advance with edge masks and not optimized during our method's training. We optimize the matrix $\hat{\theta}_i^t$, and adversarially optimize $\hat{\theta}_i^1$ towards identity matrix $\mathbb{I}$ by:

$$\mathcal{L}_{affine} = \frac{1}{N_p} \sum_{i=1}^{N_p} \left( \Delta_1 \cdot \left| \hat{\theta}_i^t - \mathbb{I} \right|_F - \Delta_2 \cdot \left| \hat{\theta}_i^1 - \mathbb{I} \right|_F \right), \quad (8)$$

where $N_p$ denotes the number of patches and $|\cdot|_F$ is Frobenius norm. $\Delta_1$ and $\Delta_2$ is two balancing hyper-parameters to control the effects of $\hat{\theta}_i^t$ and $\hat{\theta}_i^1$, empirically set as 1.00 and 0.01. In this objective, $B_{pred}^t$ is pushed toward $B_{gt}^t$ and pulled away from $B_{gt}^1$ to improve the target boundary and maintain the subtle inter-frame discrepancy in lesion structure.

We also employ the binary cross entropy loss $\mathcal{L}_{bce}$ between the whole predicted boundary and corresponding ground truths of the target frame to optimize the boundary detection further. Finally, the overall loss to optimize during training is as follows, where the scaling parameters $\lambda_1, \lambda_2$ are both empirically set as 0.3:

$$\mathcal{L}_{total} = \mathcal{L}_{seg} + \lambda_1 \mathcal{L}_{affine} + \lambda_2 \mathcal{L}_{bce}. \quad (9)$$

## IV. EXPERIMENTS

### A. Dataset

We evaluate our Vivim on three medical video segmentation tasks, *i.e.*, video thyroid ultrasound segmentation, video breast lesion ultrasound segmentation and video polyp segmentation.

*1) Video thyroid ultrasound segmentation:* We collect a video thyroid ultrasound segmentation dataset *VTUS*. VTUS comprises 100 video sequences, one video sequence per patient, and a total of 9342 frames with pixel-level ground truth. VTUS contains the transverse viewed and the longitudinal viewed B-mode ultrasound videos captured by Mindray resona8/TOSHIBA Aplio500 vendors. These videos are cross-annotated by three experts with over three years of experience in thyroid diagnosis. The number of frames in these videos vary from 31 to 196 for better diversity. The entire dataset is partitioned into training and test sets by 7:3, yielding a total of 70 training videos, 30 test videos.

*2) Video breast lesion ultrasound segmentation:* We conduct experiments on the BUV2022 dataset [27] consisting of 63 video sequences, with one video sequence per person, containing 4619 frames that have been annotated with pixel-level ground truth by experts. Following the approach outlined in [27], the video sequences with spatial resolutions ranging from 580×600 to 600×800 were further cropped to a spatial resolution of 300×200. We follow the official splits for training and testing.

*3) Video polyp segmentation:* We adopt four widely used polyp datasets, including image-based Kvasir [40] and video-based CVC-300 [41], CVC-612 [42] and ASU-Mayo [43]. Following the same protocol as [44], we train our model on Kvasir, ASU-Mayo and the training sets of CVC-300 and CVC-612, and conduct three experiments on test datasets CVC-300-TV, CVC-612-V and CVC-612-T.

### B. Implementation details

The proposed framework was trained on one NVIDIA RTX 4090 GPU and implemented on the Pytorch platform. Our framework is empirically trained for 100 epochs in an end-to-end way and the Adam optimizer is applied. The initial learning rate is set to $1 \times 10^{-4}$ and decayed to $1 \times 10^{-6}$. During training, we resize the video frames to $256 \times 256$, and feed a batch of 4 video clips, each of which has 5 frames, into the network for each iteration.

### C. Comparsion with Other Methods

*1) Results on thyroid and breast lesion US video segmentation:* We employed four segmentation evaluation metrics, including Dice, Jaccard, Precision and Recall; for their precise definitions, please refer to [19]. We also report the inference speed performance by computing the number of frames per second (FPS).

As shown in Tab. I, we quantitatively compare our method with many state-of-the-art methods on VTUS dataset and BUV2022 dataset. These methods including popular medical image segmentation methods (UNet [37], UNet++ [3], TransUNet [18], SETR [8], DAF [19]), and video object segmentation methods (OSVOS [38], ViViT [9], STM [24], AFB-URR [10], DPSTT [27], FLA-Net [28], RMem [39]). For the fairness of comparisons, we reproduce these methods following their publicly available codes. Note that we adopted Vision Transformer as the backbone of FLA-Net. We can observe that video-based methods tend to outperform image-based methods as evidenced by their better performance. This suggests that the exploration of temporal information offers significant advantages for segmenting thyroid nodules and breast lesions in ultrasound videos. More importantly, among all image-based and video-based segmentation methods, our Vivim has achieved the highest performance across all scores by a considerable margin (*e.g.*, 2.61%, 2.74% in Dice, Jaccard on VTUS, 1.01%, 0.86% in Dice, Jaccard on BUV2022 than the second-best method DPSTT). Our Vivim also has the best run-time among all video-based methods observed from FPS. This demonstrates that our solution can simultaneously learn spatial and temporal cues in an efficient way, and achieves

TABLE I

QUANTITATIVE COMPARISON WITH STATE-OF-THE-ART METHODS ON OUR VTUS DATASET (THYROID NODULE) AND THE BUV2022 DATASET (BREAST LESION). DICE, JACCARD, PRECISION AND RECALL ARE ADOPTED AS OUR EVALUATION METRICS. THE BEST SCORES ARE HIGHLIGHTED IN BOLD.

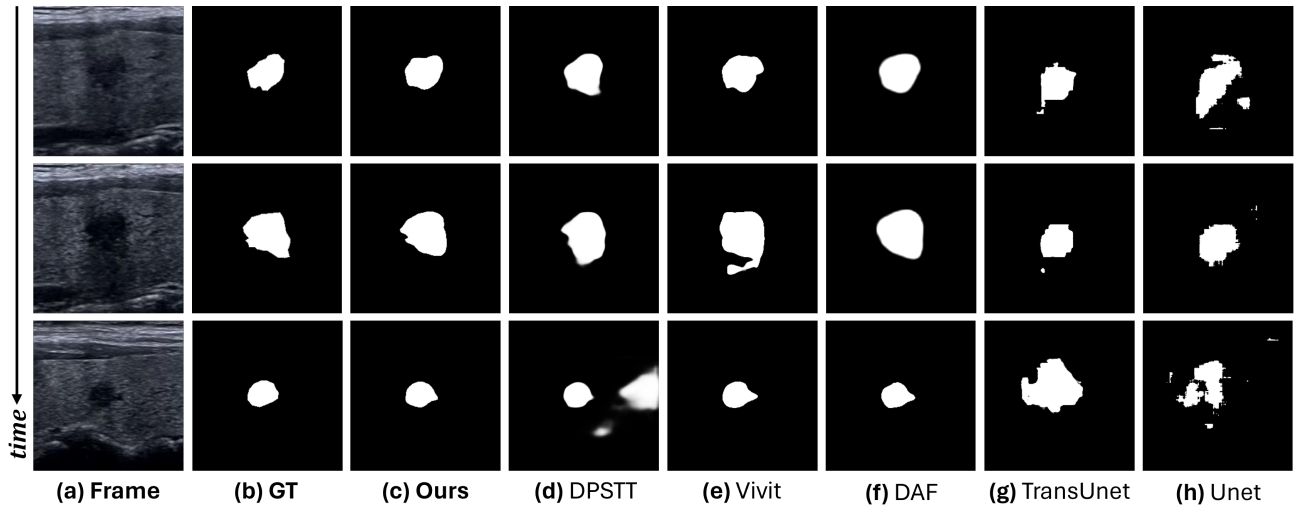| Methods | Venue | Type | VTUS | | | | BUV2022 | | | | FPS |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | Dice | Jaccard | Precision | Recall | Dice | Jaccard | Precision | Recall | |
| UNet [37] | MICCAI15 | image | 0.6662 | 0.5328 | 0.6703 | 0.7471 | 0.7303 | 0.6247 | 0.7946 | 0.7272 | **88.18** |
| UNet++ [3] | DLMIA18 | image | 0.7656 | 0.6486 | 0.7441 | 0.8496 | 0.7179 | 0.6124 | 0.8280 | 0.6884 | 40.90 |
| TransUNet [18] | arXiv21 | image | 0.7461 | 0.6250 | 0.7468 | 0.8321 | 0.6547 | 0.5358 | 0.7167 | 0.6682 | 65.10 |
| SETR [8] | CVPR21 | image | 0.7288 | 0.6010 | 0.7399 | 0.8089 | 0.6649 | 0.5480 | 0.7533 | 0.6643 | 21.61 |
| DAF [19] | MICCAI18 | image | 0.7716 | 0.6583 | 0.7046 | 0.8599 | 0.7890 | 0.6954 | 0.7992 | 0.7979 | 47.62 |
| OSVOS [38] | CVPR17 | video | 0.7769 | 0.6754 | 0.7895 | 0.8241 | 0.7098 | 0.5674 | 0.7778 | 0.6404 | 27.25 |
| ViViT [9] | ICCV21 | video | 0.7610 | 0.6459 | 0.7789 | 0.8252 | 0.6739 | 0.5446 | 0.7554 | 0.6683 | 24.33 |
| STM [24] | ICCV19 | video | 0.7898 | 0.6897 | 0.8112 | 0.8251 | 0.7862 | 0.6858 | 0.8201 | 0.7910 | 23.17 |
| AFB-URR [10] | NIPS20 | video | 0.7930 | 0.6957 | 0.7764 | 0.8429 | 0.8018 | 0.7034 | 0.8008 | 0.8591 | 11.84 |
| DPSTT [27] | MICCAI22 | video | 0.8063 | 0.7117 | 0.8238 | 0.8352 | 0.8255 | 0.7364 | **0.8389** | 0.8455 | 30.50 |
| FLA-Net [28] | MICCAI23 | video | 0.8042 | 0.7075 | 0.8121 | 0.8276 | 0.8232 | 0.7315 | 0.8334 | 0.8422 | 31.22 |
| RMem [39] | CVPR24 | video | 0.7804 | 0.6775 | 0.7821 | 0.8298 | 0.7912 | 0.6901 | 0.8024 | 0.8221 | 29.54 |
| Our method | – – | video | **0.8324** | **0.7391** | **0.8363** | **0.8711** | **0.8356** | **0.7450** | 0.8357 | **0.8869** | 35.33 |



Fig. 5. Visual comparison on video ultrasound thyroid segmentation with several competitive image- and video-based methods. Consecutive results of one case are displayed.

significant improvements over those Transformer methods, such as SETR, ViViT and DPSTT. As displayed in Fig. 5, we visualize the thyroid segmentation results on the selected frames. Our model can better locate and segment the target lesions with more accurate boundaries.

*2) Polyp video segmentation:* We adopt six metrics following [44], *i.e.*, maximum Dice (maxDice), maximum specificity (maxSpe), maximum IoU (maxIoU), S-measure [49] ($S_\alpha$), E-measure [50] ($E_\phi$), and mean absolute error (MAE).

We compare our method with existing methods as summarized in Tab. II, including UNet [37], UNet++ [3], ResUNet [45], ACSNet [46], PraNet [47], PNS-Net [44], LD-Net [48] and FLA-Net [28]. We conduct three experiments on CVC-300-TV, CVC-612-V and CVC-612-T to validate the model's performance. CVC-300-TV consists of both validation set and test set including six videos in total, while CVC-612-V and CVC-612-T each contain five videos. On CVC-300-TV,

our Vivim achieves remarkable performance and outperforms all methods by a large margin (*e.g.*, 2.7% in maxDice, 2.2% in maxIoU). On CVC-612-V and CVC-612-T, our Vivim consistently outperforms other SOTAs, especially 1.2% and 1.1% in maxDice, respectively. We also visualize the polyp segmentation results on the consecutive frames of CVC-612-T in Fig. 6. Our model demonstrates improved capability in locating and segmenting polyps with more precise boundaries.

### D. Ablation Study

Extensive experiments are conducted on VTUS dataset to evaluate the effectiveness of our major components. To do so, we construct four baseline networks from our method. The first baseline (denoted as "basic") is to remove all Mamba layers and boundary-aware affine constraint from our network. It means that "basic" equals the vanilla SegFormer [33]. Then, we introduce ST-Mamba layers with temporal forward SSM

TABLE II
QUANTITATIVE RESULTS ON THREE VIDEO POLYP DATASETS. THE BEST SCORES ARE HIGHLIGHTED IN **BOLD**. ↑ INDICATES THE HIGHER THE SCORE THE BETTER, AND VICE VERSA.

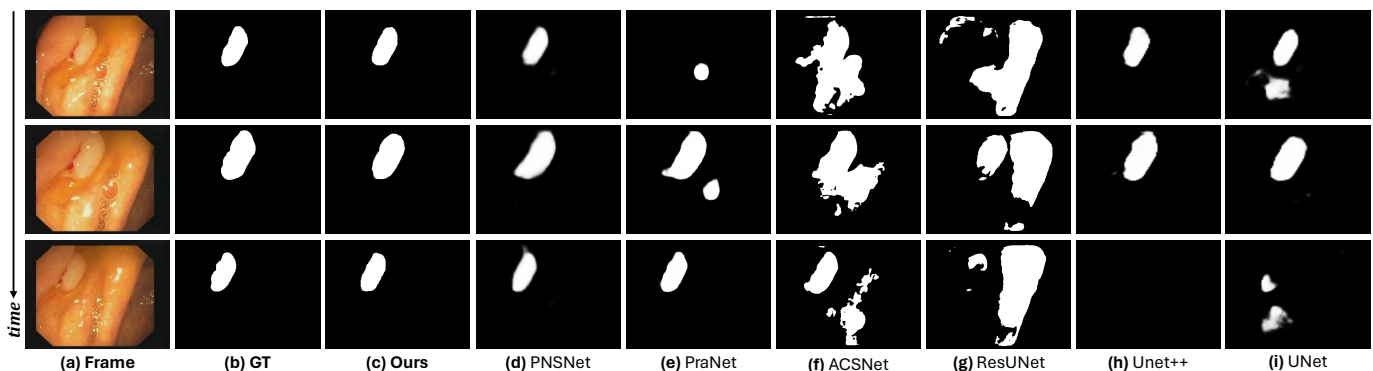| | Metrics | UNet MICCAI [37] | UNet++ TMI [3] | ResUNet ISM [45] | ACSNet MICCAI [46] | PraNet MICCAI [47] | PNS-Net MICCAI [44] | LDNet MICCAI [48] | FLA-Net MICCAI [28] | **Vivim (Ours)** |
|---|---|---|---|---|---|---|---|---|---|---|
| CVC-300-TV | maxDice↑ | 0.639 | 0.649 | 0.535 | 0.738 | 0.739 | 0.840 | 0.835 | 0.874 | **0.901** |
| | maxSpe↑ | 0.963 | 0.944 | 0.852 | 0.987 | 0.993 | 0.996 | 0.994 | 0.996 | **0.997** |
| | maxIoU↑ | 0.525 | 0.539 | 0.412 | 0.632 | 0.645 | 0.745 | 0.741 | 0.789 | **0.831** |
| | $S_\alpha$ ↑ | 0.793 | 0.796 | 0.703 | 0.837 | 0.833 | 0.909 | 0.898 | 0.907 | **0.928** |
| | $E_\phi$ ↑ | 0.826 | 0.831 | 0.718 | 0.871 | 0.852 | 0.921 | 0.910 | **0.969** | 0.958 |
| | $MAE$ ↓ | 0.027 | 0.024 | 0.052 | 0.016 | 0.016 | 0.013 | 0.015 | 0.010 | **0.008** |
| CVC-612-V | maxDice↑ | 0.725 | 0.684 | 0.752 | 0.804 | 0.869 | 0.873 | 0.870 | 0.885 | **0.897** |
| | maxSpe↑ | 0.971 | 0.952 | 0.939 | 0.929 | 0.983 | 0.991 | 0.987 | 0.992 | **0.996** |
| | maxIoU↑ | 0.610 | 0.570 | 0.648 | 0.712 | 0.799 | 0.800 | 0.799 | 0.814 | **0.829** |
| | $S_\alpha$ ↑ | 0.826 | 0.805 | 0.829 | 0.847 | 0.915 | 0.923 | 0.918 | 0.920 | **0.940** |
| | $E_\phi$ ↑ | 0.855 | 0.830 | 0.877 | 0.887 | 0.936 | 0.944 | 0.941 | 0.963 | **0.971** |
| | $MAE$ ↓ | 0.023 | 0.025 | 0.023 | 0.054 | 0.013 | 0.012 | 0.013 | 0.012 | **0.010** |
| CVC-612-T | maxDice↑ | 0.729 | 0.740 | 0.617 | 0.782 | 0.852 | 0.860 | 0.857 | 0.861 | **0.872** |
| | maxSpe↑ | 0.971 | 0.975 | 0.950 | 0.975 | 0.986 | 0.992 | 0.988 | 0.993 | **0.995** |
| | maxIoU↑ | 0.635 | 0.635 | 0.514 | 0.700 | 0.786 | 0.795 | 0.791 | 0.795 | **0.810** |
| | $S_\alpha$ ↑ | 0.810 | 0.800 | 0.727 | 0.838 | 0.886 | 0.903 | 0.892 | 0.904 | **0.915** |
| | $E_\phi$ ↑ | 0.836 | 0.817 | 0.758 | 0.864 | 0.904 | 0.903 | 0.903 | 0.904 | **0.921** |
| | $MAE$ ↓ | 0.058 | 0.059 | 0.084 | 0.053 | 0.038 | 0.038 | 0.037 | 0.036 | **0.033** |



Fig. 6. Qualitative results on the selected frames of CVC-612-T. Our Vivim can better locate and segment polyps with more accurate boundaries than several competitive image- and video-based methods.

TABLE III
ABLATION STUDY OF OUR VIVIM DESIGN ON VTUS DATASET. IN ST-MAMBA, $T^f$ DENOTES TEMPORAL FORWARD SSM, $T^b$ DENOTES TEMPORAL BACKWARD SSM, $S$ DENOTES SPATIAL SSM, WHILE BAC DENOTES BOUNDARY-AWARE AFFINE CONSTRAINT.

| | ST-Mamba | | | BAC | VTUS | | | |
|---|---|---|---|---|---|---|---|---|
| | $T^f$ | $T^b$ | $S$ | | Dice↑ | Jaccard↑ | Precision↑ | Recall↑ |
| basic | – | – | – | – | 0.8144 | 0.7188 | 0.8040 | 0.8572 |
| C1 | ✓ | – | – | – | 0.8159 | 0.7216 | 0.8170 | 0.8704 |
| C2 | ✓ | ✓ | – | – | 0.8213 | 0.7264 | 0.8239 | 0.8670 |
| C3 | ✓ | ✓ | ✓ | – | 0.8259 | 0.7310 | 0.8255 | **0.8753** |
| **Ours** | ✓ | ✓ | ✓ | ✓ | **0.8324** | **0.7391** | **0.8363** | 0.8711 |

that the vanilla SSM helps explore temporal dependency, thereby improving the segmentation performance in videos. Then, the better Dice and Jaccard results of "C2" over "C1" demonstrate that introducing our bidirectional temporal SSMs can critically benefit the cross-frame coherence. Furthermore, by adapting SSMs to non-causal information, "C3" advances "C2" with a significant margin of 0.46% in Dice and 0.83% in Recall. Finally, our method outperforms "C3" in terms of Dice, Jaccard and Precision, which indicates that the boundary-aware affine constraint can further help to enhance the thyroid segmentation results.

### E. Analysis on Efficiency of ST-Mamba

We validate the high efficiency of the proposed ST-Mamba by two ablation studies presented in Tab. IV and Fig. 7. In Tab. IV, we compare against several core modules for the modeling of spatio-temporal dependency, *i.e.*, vanilla self-attention [7], window self-attention [51] and factorized self-attention [9]. We replace ST-Mamba in our full Vivim with the three core modules to construct three variants M1, M2, and M3, respectively. We assessed the efficiency of these models using a single 48G A6000, considering Training Memory (TM), Inference Memory (IM), and Run-time as key metrics. M1, incorporating 3D global self-attention to capture spatial

($T^f$) into "basic" to construct another baseline network "C1", and further equip ST-Mamba with temporal backward SSM ($T^b$) to build a baseline network "C2". Based on "C2", spatial SSM ($S$) is incorporated into the ST-Mamba to construct "C3". Hence, "C3" is equal to removing the boundary-aware affine constraint from the training of our network. Table III reports the results of our method and four baseline networks. While "basic" performs competitively due to the pre-trained SegFormer weights on ADE20K, our proposed modules significantly advance its effectiveness. Compared to "basic", "C1" has a great improvement across all metrics, which indicates

TABLE IV
ABLATION STUDY FOR DIFFERENT ATTENTION MODULES. WE FEED A VIDEO CLIP OF 32 FRAMES WITH 256P INTO THE THREE MODEL VARIANTS AND
OUR METHOD. "TM" DENOTES TRAINING MEMORY, "IM" DENOTES INFERENCE MEMORY, AND "OOM" REPRESENTS OUT-OF-MEMORY. "IS GLOBAL"
DESCRIBES WHETHER THE CORE MODULES ARE GLOBAL MODELING ONES.

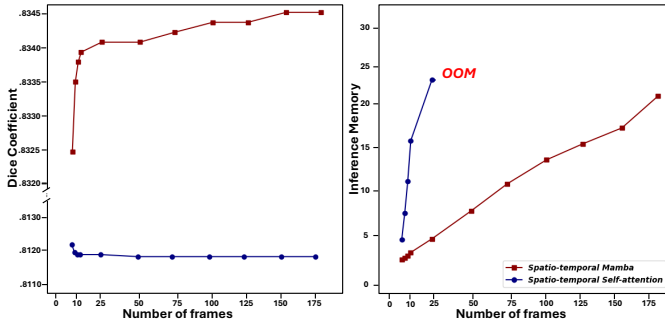| Methods | Core Module | Input Size | TM (M) | IM (M) | Run-time (s) | Is Global |
|---|---|---|---|---|---|---|
| M1 | Spatio-temporal self-attention | $32 \times 256^2$ | OOM | - | - | ✓ |
| M2 | Spatio-temporal Window self-attention | $32 \times 256^2$ | 25,861 | 7,795 | 0.142 | ✗ |
| M3 | Spatio-temporal Factorized self-attention | $32 \times 256^2$ | 29,110 | 9,288 | 0.156 | ✗ |
| **Our method** | Spatio-temporal Mamba | $32 \times 256^2$ | **19,216** | **5,112** | **0.121** | ✓ |



Fig. 7. ST-Mamba performs better in effectiveness and efficiency when addressing long sequence modeling. (a) More reference frames can help improve the segmentation performance of ST-Mamba, but it is not applicable for spatio-temporal self-attention. (b) Vivim has a lighter memory burden than traditional attention-based methods when increasing the sequence length.

and temporal information simultaneously, faces challenges due to the memory constraints when processing a video clip of 32 frames at a resolution of $256 \times 256$. In contrast, M2 and M3 compromise on the receptive field to ensure that spatio-temporal modeling can be conducted within the available memory capacity. Instead, our approach introduces an efficient global modeling module based on Mamba, leading to superior performance in terms of training memory, inference memory, and average run-time when compared to the other model variants.

Fig. 7 displays the Dice coefficient and memory costs with an increasing number of frames in one video clip at the inference stage. We evaluate M1 and our method, *i.e.*, spatio-temporal self-attention and spatio-temporal Mamba, to verify the efficiency of our model. As observed, M1 tends to maintain and even degrade the segmentation performance when referring to more neighboring frames. Instead, our method obtains an about 0.2% improvement in Dice coefficient. Furthermore, increasing the number of frames introduces an explosive growth in memory costs for M1. Our method can infer a video clip of over 150 frames using a single RTX4090. This provides a great foundation for longer medical video segmentation within the limited memory capacity.

## V. DISCUSSION

Our Vivim is designed around the diagnostic process of radiologists that obtain the holistic understanding by ST-Mamba, and preserve local details by CNN-based decoder aggregating multiscale features and boundary constraint. Although con-

volutional neural networks (CNNs) and Transformers have achieved impressive performance for many ultrasound video segmentation tasks [27], [28], there remains significant potential for enhancements in both efficiency and effectiveness. A critical challenge limiting the broader application of CNNs and Transformers in medical video analysis is the trade-off between receptive field and computational complexity. This issue arises from the inherently local processing nature of CNNs and the high computational complexity associated with Transformers. State Space Models (SSMs), *e.g.*, Mamba, offer a more efficient technique for global dependency modeling compared to Transformers, facilitating more dynamic references within ultrasound videos. Ultrasound experts typically acquire the complete appearance of target tissues by utilizing both transverse and longitudinal views, which are captured by high-frame-rate devices. This leads to the requirement of long sequence modeling in ultrasound videos. Unlike the self-attention mechanism in Transformers [7], [9], which scales quadratically with video sequence length, the computational complexity of SSMs scales linearly. This linear scalability makes SSMs well-suited for spatio-temporal joint modeling, allowing them to operate within the constraints of limited memory and computational resources, a feat challenging for Transformers, especially with longer video sequences.

The causal nature of SSMs is particularly well-aligned with tasks in Natural Language Processing (NLP) and video processing, where understanding the context in textual and temporal data is crucial [11]. A key challenge in adapting the Mamba model for medical video tasks lies in designing selective scan directions that effectively preserve non-causal spatial details of lesions and tissues while exploring temporal dependencies. Therefore, our approach goes beyond a straightforward application of Mamba; we establish a baseline that combines Transformers for spatial modeling with SSMs for spatio-temporal modeling. We introduce a tri-directional scan mechanism that simultaneously operates along temporal forward, temporal backward, and spatial forward directions, carefully balancing cross-frame coherence with single-frame spatial integrity. In the future, our focus will be on further investigating a fully SSM-based solution for efficient medical video segmentation and developing more innovative selective scan mechanisms, tailored to medical video segmentation tasks for better lesion location. We will strive to better preserve the spatial correlation between neighboring patches of the target lesion when conducting 1-D sequence interaction.

## VI. CONCLUSION

In this paper, we present a Mamba-based framework Vivim to address the challenges of medical video segmentation, especially in modeling long-range temporal dependencies due to the inherent locality of CNNs and the high computational complexity of the self-attention mechanism. The main idea of Vivim is to introduce the structured state space models with spatiotemporal selective scan, ST-Mamba, into the standard hierarchical Transformer architecture. This facilitates the exploration of single-frame spatial coherence and cross-frame coherence in a computationally cheaper way than using the self-attention mechanism. An improved boundary-aware constraint at the training stage is proposed to mitigate the ambiguous prediction of our model. We also contribute a video thyroid ultrasound segmentation dataset VTUS with 100 videos and 9342 annotated frames. Experimental results on our collected VTUS dataset, ultrasound breast lesion videos and polyp colonoscopy videos reveal that Vivim outperforms state-of-the-art segmentation networks. Ablation studies also validate the superior efficiency of ST-Mamba to other spatiotemporal Transformer-based methods.

## REFERENCES

[1] Q. Huang, Y. Huang, Y. Luo, F. Yuan, and X. Li, "Segmentation of breast ultrasound image with semantic classification of superpixels," *Medical Image Analysis*, vol. 61, p. 101657, 2020.

[2] Z. Lin, J. Lin, L. Zhu, H. Fu, J. Qin, and L. Wang, "A new dataset and a baseline model for breast lesion detection in ultrasound videos," in *MICCAI*. Springer, 2022, pp. 614–623.

[3] Z. Zhou, M. M. Rahman Siddiquee, N. Tajbakhsh, and J. Liang, "Unet++: A nested u-net architecture for medical image segmentation," in *Deep learning in medical image analysis and multimodal learning for clinical decision support*. Springer, 2018, pp. 3–11.

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

[5] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.

[6] Y. Yang, S. Wang, L. Zhu, and L. Yu, "Hcdg: A hierarchical consistency framework for domain generalization on medical image segmentation," *arXiv preprint arXiv:2109.05742*, 2021.

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

[8] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr *et al.*, "Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 6881–6890.

[9] A. Arnab, M. Dehghani, G. Heigold, C. Sun, M. Lučić, and C. Schmid, "Vivit: A video vision transformer," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 6836–6846.

[10] Y. Liang, X. Li, N. Jafari, and J. Chen, "Video object segmentation with adaptive feature bank and uncertain-region refinement," *Advances in Neural Information Processing Systems*, vol. 33, pp. 3430–3441, 2020.

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

[12] R. E. Kalman, "A new approach to linear filtering and prediction problems," 1960.

[13] J. Ma, F. Li, and B. Wang, "U-mamba: Enhancing long-range dependency for biomedical image segmentation," *arXiv preprint arXiv:2401.04722*, 2024.

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

[15] K. He, C. Gan, Z. Li, I. Rekik, Z. Yin, W. Ji, Y. Gao, Q. Wang, J. Zhang, and D. Shen, "Transformers in medical image analysis," *Intelligent Medicine*, vol. 3, no. 1, pp. 59–78, 2023.

[16] Y. He, V. Nath, D. Yang, Y. Tang, A. Myronenko, and D. Xu, "Swinunetr-v2: Stronger swin transformers with stagewise convolutions for 3d medical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2023, pp. 416–426.

[17] A. Hatamizadeh, Y. Tang, V. Nath, D. Yang, A. Myronenko, B. Landman, H. R. Roth, and D. Xu, "Unetr: Transformers for 3d medical image segmentation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 574–584.

[18] 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," *arXiv preprint arXiv:2102.04306*, 2021.

[19] Y. Wang, Z. Deng, X. Hu, L. Zhu, X. Yang, X. Xu, P.-A. Heng, and D. Ni, "Deep attentional features for prostate segmentation in ultrasound," in *Medical Image Computing and Computer Assisted Intervention–MICCAI 2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part IV 11*. Springer, 2018, pp. 523–530.

[20] S. Ali, Y. Espinel, Y. Jin, P. Liu, B. Güttner, X. Zhang, L. Zhang, T. Dowrick, M. J. Clarkson, S. Xiao *et al.*, "An objective comparison of methods for augmented reality in laparoscopic liver resection by preoperative-to-intraoperative image fusion," *arXiv preprint arXiv:2401.15753*, 2024.

[21] J. M. Webb, D. D. Meixner, S. A. Adusei, E. C. Polley, M. Fatemi, and A. Alizad, "Automatic deep learning semantic segmentation of ultrasound thyroid cineclips using recurrent fully convolutional networks," *IEEE Access*, vol. 9, pp. 5119–5127, 2020.

[22] L. Ma, G. Tan, H. Luo, Q. Liao, S. Li, and K. Li, "A novel deep learning framework for automatic recognition of thyroid gland and tissues of neck in ultrasound image," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 9, pp. 6113–6124, 2022.

[23] J. Chi, Z. Li, Z. Sun, X. Yu, and H. Wang, "Hybrid transformer unet for thyroid segmentation from ultrasound scans," *Computers in Biology and Medicine*, vol. 153, p. 106453, 2023.

[24] S. W. Oh, J.-Y. Lee, N. Xu, and S. J. Kim, "Video object segmentation using space-time memory networks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9226–9235.

[25] H. K. Cheng, Y.-W. Tai, and C.-K. Tang, "Rethinking space-time networks with improved memory coverage for efficient video object segmentation," *Advances in Neural Information Processing Systems*, vol. 34, pp. 11781–11794, 2021.

[26] K. Park, S. Woo, S. W. Oh, I. S. Kweon, and J.-Y. Lee, "Per-clip video object segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 1352–1361.

[27] J. Li, Q. Zheng, M. Li, P. Liu, Q. Wang, L. Sun, and L. Zhu, "Rethinking breast lesion segmentation in ultrasound: A new video dataset and a baseline network," in *MICCAI*. Springer, 2022, pp. 391–400.

[28] J. Lin, Q. Dai, L. Zhu, H. Fu, Q. Wang, W. Li, W. Rao, X. Huang, and L. Wang, "Shifting more attention to breast lesion segmentation in ultrasound videos," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2023, pp. 497–507.

[29] A. Gu, I. Johnson, K. Goel, K. Saab, T. Dao, A. Rudra, and C. Ré, "Combining recurrent, convolutional, and continuous-time models with linear state space layers," *Advances in neural information processing systems*, vol. 34, pp. 572–585, 2021.

[30] A. Gu, K. Goel, and C. Ré, "Efficiently modeling long sequences with structured state spaces," *arXiv preprint arXiv:2111.00396*, 2021.

[31] E. Nguyen, K. Goel, A. Gu, G. Downs, P. Shah, T. Dao, S. Baccus, and C. Ré, "S4nd: Modeling images and videos as multidimensional signals with state spaces," *Advances in neural information processing systems*, vol. 35, pp. 2846–2861, 2022.

[32] Y. Liu, Y. Tian, Y. Zhao, H. Yu, L. Xie, Y. Wang, Q. Ye, and Y. Liu, "Vmamba: Visual state space model," *arXiv preprint arXiv:2401.10166*, 2024.

[33] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, "Segformer: Simple and efficient design for semantic segmentation with transformers," *Advances in Neural Information Processing Systems*, vol. 34, pp. 12077–12090, 2021.

[34] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 568–578.

[35] S. Borse, Y. Wang, Y. Zhang, and F. Porikli, "Inverseform: A loss function for structured boundary-aware segmentation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 5901–5911.

[36] J. Canny, "A computational approach to edge detection," *IEEE Transactions on pattern analysis and machine intelligence*, no. 6, pp. 679–698, 1986.

[37] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015, pp. 234–241.

[38] F. Perazzi, A. Khoreva, R. Benenson, B. Schiele, and A. Sorkine-Hornung, "Learning video object segmentation from static images," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2663–2672.

[39] J. Zhou, Z. Pang, and Y.-X. Wang, "Rmem: Restricted memory banks improve video object segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 18 602–18 611.

[40] D. Jha, P. H. Smedsrud, M. A. Riegler, P. Halvorsen, T. de Lange, D. Johansen, and H. D. Johansen, "Kvasir-seg: A segmented polyp dataset," in *MultiMedia Modeling: 26th International Conference, MMM 2020, Daejeon, South Korea, January 5–8, 2020, Proceedings, Part II 26*. Springer, 2020, pp. 451–462.

[41] J. Bernal, J. Sánchez, and F. Vilarino, "Towards automatic polyp detection with a polyp appearance model," *Pattern Recognition*, vol. 45, no. 9, pp. 3166–3182, 2012.

[42] J. Bernal, F. J. Sánchez, G. Fernández-Esparrach, D. Gil, C. Rodríguez, and F. Vilariño, "Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians," *Computerized medical imaging and graphics*, vol. 43, pp. 99–111, 2015.

[43] N. Tajbakhsh, S. R. Gurudu, and J. Liang, "Automated polyp detection in colonoscopy videos using shape and context information," *IEEE transactions on medical imaging*, vol. 35, no. 2, pp. 630–644, 2015.

[44] G.-P. Ji, Y.-C. Chou, D.-P. Fan, G. Chen, H. Fu, D. Jha, and L. Shao, "Progressively normalized self-attention network for video polyp segmentation," in *MICCAI*. Springer, 2021, pp. 142–152.

[45] D. Jha, P. H. Smedsrud, M. A. Riegler, D. Johansen, T. De Lange, P. Halvorsen, and H. D. Johansen, "Resunet++: An advanced architecture for medical image segmentation," in *2019 IEEE international symposium on multimedia (ISM)*. IEEE, 2019, pp. 225–2255.

[46] R. Zhang, G. Li, Z. Li, S. Cui, D. Qian, and Y. Yu, "Adaptive context selection for polyp segmentation," in *MICCAI*. Springer, 2020, pp. 253–262.

[47] D.-P. Fan, G.-P. Ji, T. Zhou, G. Chen, H. Fu, J. Shen, and L. Shao, "Pranet: Parallel reverse attention network for polyp segmentation," in *MICCAI*. Springer, 2020, pp. 263–273.

[48] R. Zhang, P. Lai, X. Wan, D.-J. Fan, F. Gao, X.-J. Wu, and G. Li, "Lesion-aware dynamic kernel for polyp segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2022, pp. 99–109.

[49] D.-P. Fan, M.-M. Cheng, Y. Liu, T. Li, and A. Borji, "Structure-measure: A new way to evaluate foreground maps," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 4548–4557.

[50] D.-P. Fan, G.-P. Ji, X. Qin, and M.-M. Cheng, "Cognitive vision inspired object segmentation metric and loss function," *Scientia Sinica Informationis*, vol. 6, no. 6, 2021.

[51] Z. Liu, J. Ning, Y. Cao, Y. Wei, Z. Zhang, S. Lin, and H. Hu, "Video swin transformer," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 3202–3211.