# ReMamber: Referring Image Segmentation with Mamba Twister

Yuhuan Yang*[1], Chaofan Ma*[1], Jiangchao Yao[1], Zhun Zhong[✉2], Ya Zhang[1], and Yanfeng Wang[✉1]

[1] Shanghai Jiao Tong University
[2] University of Nottingham

{yangyuhuan,chaofanma,sunarker,ya_zhang,wangyanfeng622}@sjtu.edu.cn
zhunzhong007@gmail.com

**Abstract.** Referring Image Segmentation (RIS) leveraging transformers has achieved great success on the interpretation of complex visual-language tasks. However, the quadratic computation cost makes it resource-consuming in capturing long-range visual-language dependencies. Fortunately, Mamba addresses this with efficient linear complexity in processing. However, directly applying Mamba to multi-modal interactions presents challenges, primarily due to inadequate channel interactions for the effective fusion of multi-modal data. In this paper, we propose **ReMamber**, a novel RIS architecture that integrates the power of Mamba with a multi-modal *Mamba Twister* block. The Mamba Twister explicitly models image-text interaction, and fuses textual and visual features through its unique channel and spatial *twisting mechanism*. We achieve competitive results on three challenging benchmarks with a simple and efficient architecture. Moreover, we conduct thorough analyses of **ReMamber** and discuss other fusion designs using Mamba. These provide valuable perspectives for future research. The code has been released at: https://github.com/yyh-rain-song/ReMamber.

**Keywords:** Referring Image Segmentation (RIS) · Multi-Modal Understanding · Mamba Architecture

## 1 Introduction

Referring Image Segmentation (RIS) is a crucial yet challenging task in the area of multi-modal understanding [15, 25, 29]. Unlike ordinary image segmentation, RIS involves the identification and segmentation of specific objects in image based on the textual descriptions. This thereby requires the model to be capable of understanding vision-language interactions, which is the core challenge of RIS.

Thanks to the powerful attention mechanisms, it has achieved great success for RIS by leveraging transformers to promote the exact recognition of multi-modal information. For example, existing works have designed transformer decoder [4, 22] or transformer encoder-decoder [24, 35, 36, 51, 55, 56] to comprehensively fuse visual and linguistic features, which have achieved great progress.
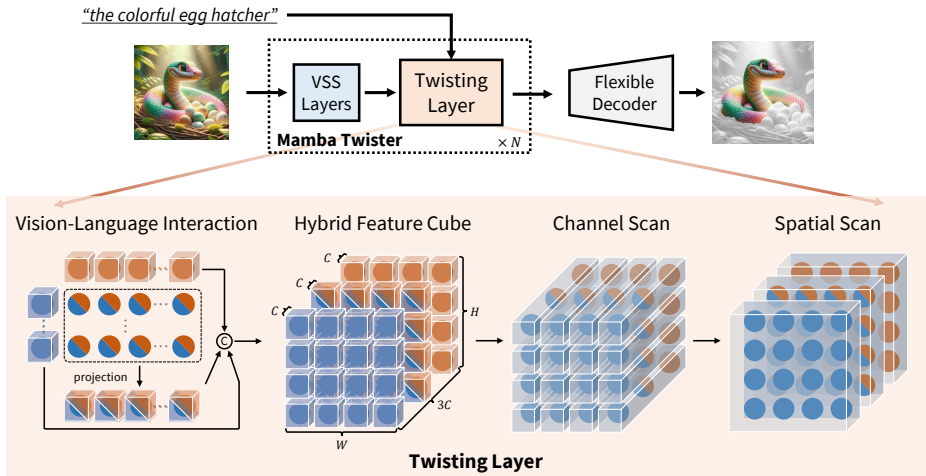
**Fig. 1:** We propose **ReMamber**, a novel referring segmentation architecture with Mamba twister. It consists of several *Mamba Twister block*. Each block contains several *visual state space (VSS) layers* and a *Twisting layer*. The Twisting layer first calculates the interaction between image and text, and then forms a hybrid feature cube. Finally, it "twists" the feature cube using the Channel and Spatial Scan along each dimension.

Nonetheless, it is a quadratic increase in both computation and memory when applying full attention in transformers [5, 49]. This leads to the limitation for resource-intensive scenarios, for example, in capturing long-range visual-language dependencies. And this is particularly important in the context of large-size images with long textual descriptions.

Fortunately, recent advances of State Space Models (SSMs) [7, 8, 11, 46] has emerged as promising architectures for solving the above issue. Specifically, Mamba [8] marks a significant advancement for efficient training and inference with linear complexity, which has been incorporated into various visual tasks. However, current efforts primarily pay attention to *single-modality* settings, such as image classification [33, 62], biomedical image segmentation [31, 37, 44, 52], low-level vision [14, 60], and point cloud analysis [26, 59]. In this paper, we pioneer the exploration of Mamba in *multi-modal* RIS setting, and identify that the prevalent multi-modality token splicing method in transformers [4, 24, 56] is no longer effective. Since in Mamba, a fundamental deficiency exists: the interactions are insufficient between channels of different tokens [8], which adversely affects the fusion of multi-modal information in RIS.

To address this dilemma, we propose **ReMamber**, a novel referring segmentation architecture with Mamba twister. As shown in Fig. 1, **ReMamber** comprises several *Mamba Twister blocks*, which allow the model to be "aware" of the textual context at every spatial location. Each block consists of several *visual state space (VSS) layers* and a *Twisting layer*. Specifically, the VSS layers initially extract visual features, and the Twisting layer injects the textual information into the visual modality. The *Twisting layer* is structured into three critical

components. **(1)** A *vision-language interaction operation* is designed to explicitly capture the fine-grained interactions between modalities. This is achieved by evaluating the similarity between visual and textual tokens then mapping them into a shared feature space, thus generating a multi-modal feature. **(2)** A *hybrid feature cube* is then created by concatenating the visual feature, multi-modal feature, and global textual feature. This ensures that each visual token receives uniform influence from both local and global contexts, preventing the overshadowing of subtle textual cues by predominant visual features. **(3)** Lastly, to address the inadequate interaction within channels in Mamba, a *twisting mechanism* is deployed. This process "twists" the feature cube channel- and spatial-wise, thereby enhancing interaction within and across modalities. It is accomplished through two consecutive SSMs scanning along channel and spatial dimensions, respectively.

To sum up, the key contributions of our research are highlighted as follows:
• We pioneer the exploration of Mamba in referring image segmentation (RIS), demonstrating Mamba's significant potential for multi-modal understanding.
• We design a novel framework, **ReMamber**, that mainly contains several Mamba Twister blocks. This design effectively captures vision-language interactions using the "twisting mechanism".
• We achieve competitive results on three challenging benchmarks. Moreover, we conduct thorough analyses of **ReMamber** and discuss other vision-language fusion designs using Mamba. These provide valuable perspectives for future research.

## 2    Related Work

### 2.1    Referring Image Segmentation

Referring image segmentation (RIS) aims to segment the entities in the image following natural language instruction. Early approaches [15, 25, 29, 39] leveraged RNNs or LSTMs to encode linguistic representations, and CNNs to extract spatial features from the image at varying levels. These disparate modalities were then integrated using the multi-modal LSTM [2, 15, 21, 25, 29, 39], attention mechanism [6, 16, 45, 57], cycle-consistency [3], and graph convolution [17].

Recent trends have shifted towards leveraging transformers for enhanced capturing and fusion of vision-language modalities. MDETR [22] and VLT [4] design a transformer decoder for fusing linguistic and visual features. LAVT [56] adopts Swin Transformer as the visual backbone and incorporate vision-language fusion modules within the visual encoder's final layers. Similar strategies are employed by ReSTR [24] and CRIS [51], which utilize dual transformer encoders for initial modality encoding, followed by feature fusion through a multi-modal transformer encoder or decoder. Other models like PolyFormer [30], SeqTr [61] and [43] also adopt a multi-modal transformer for vision-language fusion, but they output masks as sequences of contour points. Meanwhile, GRES [28] and CGFormer [47] consider transformer queries as region proposals, and regarding

segmentations as proposal-level classification problems. Distinct from all existing works, we introduce a pioneering multi-modal architecture named **ReMamber**, based on Mamba [8]. This novel approach underscores the untapped potential of Mamba in advancing the field of referring image segmentation.

## 2.2   State Space Models and Visual Applications

State space models (SSMs), originally derived from control theory, have been effectively combined with deep learning to model the long-range dependencies. Early works like LSSL [12] show potential when combined with HiPPO [9] initialization. S4 [11] was designed to diminish both computational and memory demands associated with state representations. It scales linearly with sequence length, offering a notable advantage over CNNs and transformers. Building on S4, S5 [46] introduces the efficient parallel scan and MIMO SSM, further refining the approach. Later, H3 [7] expanded on these foundations, achieving competitive performance with transformers in language modeling. Recently, Mamba [8] marked a significant advancement with its linear-time inference and efficient training, incorporating a selection mechanism and hardware-aware algorithms building upon prior works [10, 13, 40].

With the demonstrated success of SSMs like Mamba in language modeling, researchers have begun to investigate their applicability to visual tasks. Works such as ViS4mer [19], Selective S4 [50] and TranS4mer [20] directly use the S4's ability to capture long-range sequences for understanding inter-frame relations in video classification and detection. However, they still employ Vision Transformers for intra-frame feature extraction. More recently, Vim [62] and VMamba [33] have shown promising results as fully Mamba-based vision backbones for image classification, detection, and segmentation. At the same time, several studies have explored Mamba-based architecture on various vision tasks, such as biomedical image segmentation [31, 37, 44, 52], low-level vision [14, 60], and point cloud analysis [26, 59]. Nevertheless, all above efforts primarily focus on *single-modality* vision tasks. In this paper, we pioneer the exploration of Mamba-base architecture in *multi-modality* RIS task.

## 3   Methods

This paper aims to develop a multi-modal Mamba-based architecture for RIS task. In Sec. 3.1, we give the preliminary. From Sec. 3.2 to Sec. 3.5, we describe our architecture from coarse to fine. In Sec. 3.6, we introduce other three variants besides our design. And in Sec. 3.7 we detail the training process.

### 3.1   Preliminary: State Space Model

State Space Model (SSM) is a sequence model inspired by continuous systems. It is designed to capture a mapping relationship between two functions or sequences, expressed as $x(t) \in \mathbb{R} \mapsto y(t) \in \mathbb{R}$, through a hidden state $h(t) \in \mathbb{R}^N$.

The evolution of the hidden state over time is governed by specific parameters $\mathbf{A} \in \mathbb{R}^{N \times N}$, which directs the state evolution, and $\mathbf{B} \in \mathbb{R}^N, \mathbf{C} \in \mathbb{R}^N$, which performs the projection. Formally, SSM can be described by the following equations:

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

Discretizing this system introduces a timescale parameter $\mathbf{\Delta}t$, which transforms the continuous parameters $\mathbf{A}, \mathbf{B}$ into their discrete counterparts $\overline{\mathbf{A}}, \overline{\mathbf{B}}$. A common way for this transformation is the zero-order hold (ZOH) approach, defined as:

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

Then, the discretized version of this system can be represented as:

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

Mamba [8], as a variant of SSM, releases the linearity constraint of the original SSM described in Eq. (1), by making $\mathbf{B}$ and $\mathbf{C}$ to be input-dependent, whereas still maintaining the linear complexity during forward process.

### 3.2   Architecture Overview

The pivotal aspect of Referring Image Segmentation (RIS) is the correspondence between image and textual input. To achieve this, we propose **ReMamber**. Fig. 2 shows the overview of our architecture. The basic block for **ReMamber** is the *Mamba Twister* block. It is a multi-modal fusion block that takes both visual and textual feature as input, and outputs the fused multi-modal feature representation. We extract the intermediate feature after each Mamba Twister block, and then feed it into a flexible decoder to generate the final segmentation mask. The decoder is task-invariant, so any downstream architecture can be applied. Below in Sec. 3.3, we will first introduce the key component: Mamba Twister block.

### 3.3   Mamba Twister Block

As shown in Fig. 2, the Mamba Twister Block consists of several *visual state space (VSS) layers* and a *Twisting layer*. The *VSS layer* is designed to process features in the spatial domain. It treats the input feature as a series of image tokens and aims to discern the spatial relationships among these tokens. The *Twisting layer* aims to inject text condition into the image feature, thus guiding the transformation of image feature. It begins by condensing the textual sequence to compute the cross-correlation between the image and textual modalities. Subsequently, it disseminates this information across each image feature patch via a twisting operation. We'll detail the introduction of VSS layer in Sec. 3.4 and the Twisting layer in Sec. 3.5.
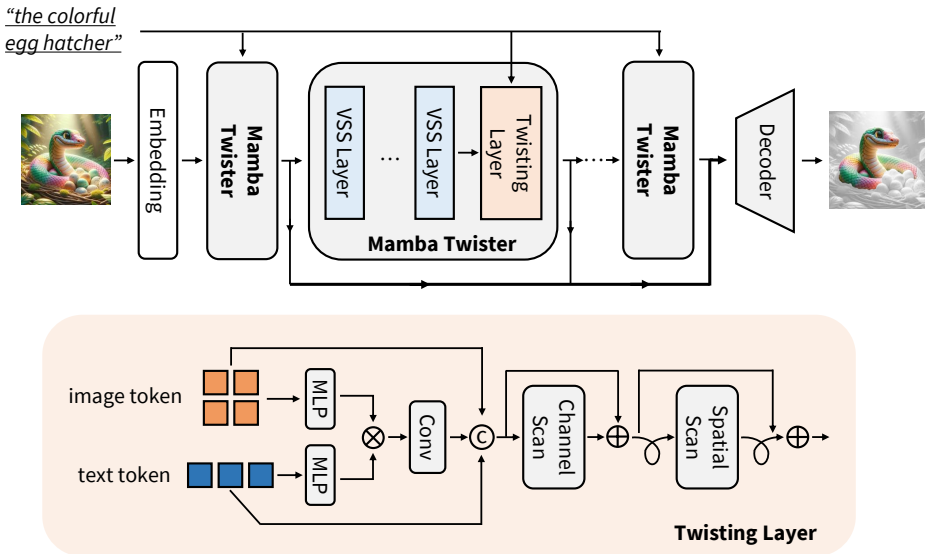
**Fig. 2: Overview architecture of our ReMamber.** The basic block for **ReMamber** is the *Mamba Twister block*. It consists of several visual state space *(VSS) layers* and a *Twisting layer*. The Twisting layer first constructs hybrid feature cube from text, image, and multi-modal features via channel concatenation. Then, it "twists" the cube by Channel Scan and Spatial Scan. We extract intermediate features after each Mamba Twister block, and feed it into a flexible decoder for final segmentation.

### 3.4 VSS Layer for Spatial Data Processing

Since the SSM is initially designed for processing temporal or causal data, which can not effectively process non-causal data types, *i.e.*, 2-dimensional image in our case. To solve this issue, we adopt the Cross-Scan-Module proposed in [33], which unfolds image patches into sequences and then scans them in four distinct directions, ensuring that the information from all pixels is integrated during feature transformation. We replace the scanning operation in vanilla Mamba block with CSM, forming a visual state space (VSS) layer for our spatial feature transformation.

### 3.5 Twisting Layer for Multi-modal Fusion

Here we introduce the detailed design of our *Twisting layer*. Intuitively, the Twisting layer first constructs a feature cube by arranging information from different modalities orderly. Then, pieces of information intertwine with each other when the cube is *twisted* through two SSM layers along different axis, thereby achieving modality fusion. Formally, given image feature $\mathbf{F}_i \in \mathbb{R}^{h \times w \times C_i}$ and text feature $\mathbf{F}_t \in \mathbb{R}^{L \times C_t}$, the Mamba twister aims to learn a function $(\mathbf{F}_i, \mathbf{F}_t) \mapsto \tilde{\mathbf{F}} \in \mathbb{R}^{h \times w \times C_o}$ by first forming a hybrid multi-modal feature cube and then conduct two SSM scans on it. $C_i, C_t$ and $C_o$ are the dimensions of visual, textual, and output feature; $h, w$ and $L$ are height, weight, and length.

**Forming the Hybrid Feature Cube.** To explicitly build the image-text correspondence, we design a novel *vision-language interaction operation* and then form a *hybrid feature cube*. This allows more effective modality fusion.

To be specific, we formulated the vision-language interaction operation into two manners: **global interaction** and **local interaction**. The **global interaction** treats the text sequence as a whole, which means that all image patches should be aware of the semantic meaning of the text expression. We pool a global representation $\mathbf{F}_t^{\text{CLS}} \in \mathbb{R}^{C_t}$ from the text sequence $\mathbf{F}_t$ to convey this information. The $\mathbf{F}_t^{\text{CLS}}$ is further expanded to the same size as the image feature. Formally:

$$\tilde{\mathbf{F}}_t = \text{Expand}(\mathbf{F}_t^{\text{CLS}}) \in \mathbb{R}^{h \times w \times C_t}, \tag{4}$$

where $\text{Expand}(\cdot)$ denotes the operation that expands the input tensor to the same size as the image feature.

However, a global representation may not be enough to capture the intricate relationships between different modalities. For example, some words such as color, shape or location may be more relevant to certain image patches than others. The **local interaction** aims to capture such correlation in a fine-grained level. Formally, we calculate the local interaction map $\mathbf{F}_c$ using matrix multiplication:

$$\mathbf{F}_c = \mathbf{F}_i \mathbf{W}_i \cdot (\mathbf{F}_t \mathbf{W}_t)^T \in \mathbb{R}^{h \times w \times L}, \tag{5}$$

where $\mathbf{W}_i \in \mathbb{R}^{C_i \times C_c}$ and $\mathbf{W}_t \in \mathbb{R}^{C_t \times C_c}$ are learnable parameters. To enhance feature processing in high dimension, we use a single convolutional layer to transform $\mathbf{F}_c$ into $\tilde{\mathbf{F}}_c \in \mathbb{R}^{h \times w \times C_c}$. We then concatenate $\mathbf{F}_i$, $\tilde{\mathbf{F}}_t$, and $\tilde{\mathbf{F}}_c$ along the channel dimension to form the hybrid feature cube. Formally:

$$\mathbf{F}_{\text{cube}} = [\mathbf{F}_i, \tilde{\mathbf{F}}_t, \tilde{\mathbf{F}}_c] \in \mathbb{R}^{h \times w \times (C_i + C_t + C_c)}. \tag{6}$$

**Twisting the Hybrid Feature Cube.** Generally, the scanning operation for vanilla SSM is almost independent across channels [8]. This is not enough for feature communication, especially when information from different modalities is arranged along the channel dimension, *i.e.*, our hybrid feature cube.

To foster the interactions **within** and **across** modalities, here we design a **twisting mechanism** to "twist" the hybrid feature cube along different axis. Specifically, it is composed of the Channel Scan and Spacial Scan successively. The **Channel Scan** first treats the channel as an ordered sequence, and learns to communicate cross channels, thus foster the modality fusion. Then, the **Spatial Scan** operates on the spatial dimension, and learns to communicate information cross patches within each modality separately. Formally, this process can be written as:

$$\mathbf{F}_{\text{out}} = \text{SSM}_{\text{spatial}}\left(\text{SSM}_{\text{channel}}\left(\mathbf{F}_{\text{cube}}\right)\right), \tag{7}$$

where $\text{SSM}_{\text{channel}}$ refers to the Channel Scan. It is an SSM layer that treats the concatenated feature as an ordered sequence through channel dimension, and presents 1-D selective scan along the channel. While $\text{SSM}_{\text{spatial}}$ denotes the Spatial Scan. It is an VSS layer and presents 2D selective scan along the two spatial dimension. $\mathbf{F}_{\text{out}}$ is then feed into the next layer for further process.
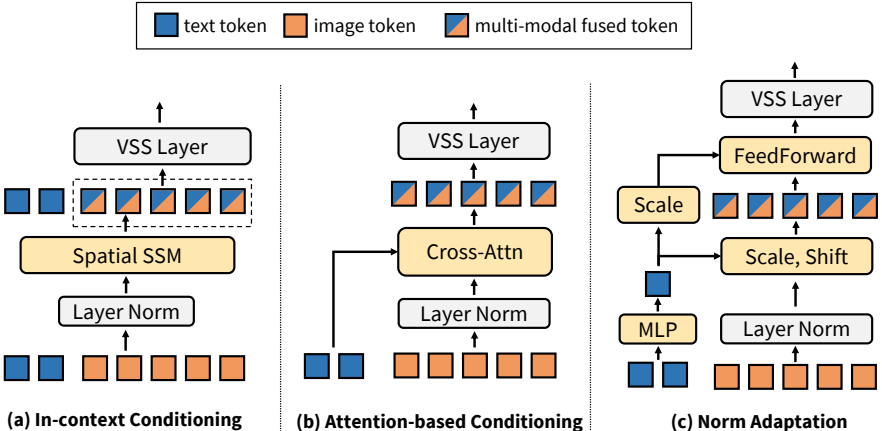
Fig. 3: **Other multi-modal fusion designs.** (a) **In-context Conditioning** appends text tokens ahead of image tokens. (b) **Attention-based Conditioning** utilizes cross-attention mechanism for modality fusion. (c) **Norm Adaptation** learns a scale and bias for the model's normalization layers.

## 3.6    Discussion: Other Variants of Modality Fusion Designs

Besides the proposed Twisting layer, we also explore other variants of multi-modal fusion designs. Fig. 3 shows the structure of three other different variants. We here discuss these variants, and later in Sec. 4.3, we analyze the results for these variants.

**In-context Conditioning.** As shown in Fig. 3 (a), in this variant, we simply append the text sequence before the image feature, forming a longer sequence than originally being processed. In this way, the model is able to fuse image features with the previous text as a context. This is the most straightforward way to allow the model aware of the text condition.

**Attention-based Conditioning.**   Attention mechanism is usually a strong baseline for sequence modeling. Here we utilize the cross-attention mechanism for two modality fusion, as shown in Fig. 3 (b), To be specific, we use the image feature as query, and the text feature as key and value, hoping to capture the correlation between image and text.

**Norm Adaptation.**   In this variant, we integrate the text input by adapting the scale and bias after norm layer using FiLM [42]. We first pool a global representation from text sequence, and then use a linear projection to transform it into the scale and shift in layer normalization.

## 3.7    Model Training

The proposed **ReMamber** is a general framework for multi-modal fusion. Therefore, there is no restriction on the downstream decoder design nor the loss function. Here we use a simple convolution-based decoder for simplicity. The entire network is trained in an end-to-end manner.

## 4    Experiments

### 4.1    Datasets and Metrics

**Dataset.** To assess the efficacy of our proposed approach, we executed experiments across three widely recognized datasets tailored for the referring image segmentation (RIS) task: RefCOCO, RefCOCO+ [23], and G-Ref [23,38,41]. RefCOCO [23] is a prominent dataset in the RIS domain, comprising 19,994 images and 142,210 referring expressions associated with 50,000 objects derived from the MSCOCO [27] dataset through a two-player game. RefCOCO+ enhances the challenge by excluding expressions with absolute location references, featuring 141,564 expressions for 49,856 objects across 19,992 images. G-Ref [38,41] enriches the dataset diversity with 104,560 referring expressions for 54,882 objects in 26,711 images, showcasing an average sentence length of 8.4 words with a heightened focus on location and appearance descriptions.

**Evaluation Metrics.** In line with preceding studies [48], we adopt mIoU and oIoU as the main metrics. Furthermore, we incorporate the metric Precision@X (X$\in \{50, 60, 70, 90\}$) for a more comprehensive evaluation, whereas Percision@X means the percentage of test images with an IoU score high than X%.

### 4.2    Implementation Details

Our model architecture is developed upon the foundations of Mamba [8] and VMamba [33]. We adopt ImageNet pre-trained weights as initialization, and train the model in an end-to-end manner. We set the input image resolution to 480 when compared with the state-of-the-art methods in Tab. 2, following [48,56] for fair comparison. For other experiments (Tabs. 1 and 3 to 4), we set the image resolution to 256 for faster training and ablation study without changing its property. For details about the training settings, please refer to our code.[1]

### 4.3    Evaluation on Other Variants of Modality Fusion Designs

As discussed in Sec. 3.6, besides our Mamba Twister, we also provide three other variant architectures for modality fusion, including **In-context Conditioning**, **Attention-based Conditioning** and **Norm Adaptation**. To make a fair comparison, for all architectures, we initialize their common part (VMamba-based parameters) with the same checkpoint. The difference is that the Twisting layer is replaced by the corresponding fusion methods. The parameters of fusion modules are nearly equal, *i.e.*, {103, 113, 118, 116} KB for {In-Context, Attention, Norm Adaptation, Twister}. Here we provide a comprehensive analysis between these architectures, showing the advantages and disadvantages of each variant.

Tab. 1 shows the comparison result. The Mamba Twister consistently outperforms the other variants across all metrics and datasets, indicating its superior capability in capturing and integrating contextual information for more accurate segmentation.

---

[1] https://github.com/yyh-rain-song/ReMamber

**Table 1: Comparison with other modality fusion variants.** "Attention" means for "Attention-based Conditioning", "In-Context" for "In-Context Conditioning", and "Adaptation" for "Norm Adaptation". Mamba Twister steadily outperforms other variants across all metrics and datasets, indicating its superior capability in capturing and integrating contextual information for more accurate segmentation.

Dataset: RefCOCO

| | val | | | | testA | | | | testB | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Variants | mIoU | oIoU | Pr@50 | Pr@70 | mIoU | oIoU | Pr@50 | Pr@70 | mIoU | oIoU | Pr@50 | Pr@70 |
| Attention | 65.3 | 62.3 | 75.1 | 60.6 | 67.5 | 64.4 | 79.1 | 65.4 | 61.7 | 58.0 | 69.6 | 53.4 |
| In-Context | 69.1 | 65.9 | 79.7 | 67.5 | 71.2 | 68.9 | 82.2 | 71.8 | 66.2 | 62.8 | 75.0 | 61.7 |
| Adaptation | 70.2 | 67.0 | 80.7 | 69.6 | 72.3 | 70.2 | 83.1 | 73.2 | 66.8 | 62.8 | 75.1 | 63.1 |
| **Mamba Twister** | **71.6** | **68.4** | **82.1** | **70.9** | **73.3** | **71.5** | **84.8** | **74.5** | **68.4** | **64.5** | **77.1** | **64.9** |

Dataset: RefCOCO+

| | val | | | | testA | | | | testB | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Variants | mIoU | oIoU | Pr@50 | Pr@70 | mIoU | oIoU | Pr@50 | Pr@70 | mIoU | oIoU | Pr@50 | Pr@70 |
| Attention | 54.0 | 49.7 | 60.7 | 45.5 | 58.7 | 55.8 | 68.5 | 52.9 | 46.9 | 41.9 | 50.6 | 36.6 |
| In-Context | 58.4 | 54.0 | 66.3 | 53.9 | 63.2 | 59.4 | 72.2 | 60.5 | 51.5 | 47.5 | 56.4 | 44.5 |
| Adaptation | 60.3 | 55.2 | 68.2 | 57.3 | 64.2 | 59.9 | 73.4 | 62.9 | 53.9 | 48.7 | 59.0 | 47.2 |
| **Mamba Twister** | **61.6** | **57.3** | **70.0** | **58.5** | **65.8** | **62.1** | **75.5** | **64.9** | **54.0** | **49.9** | **59.4** | **47.8** |

Dataset: G-Ref

| | val | | | | | | test | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Variants | mIoU | oIoU | Pr@50 | Pr@60 | Pr@70 | Pr@90 | mIoU | oIoU | Pr@50 | Pr@60 | Pr@70 | Pr@90 |
| Attention | 50.5 | 49.2 | 54.2 | 46.0 | 36.7 | 8.1 | 50.8 | 50.8 | 57.3 | 49.8 | 38.6 | 6.4 |
| In-Context | 54.8 | 53.4 | 59.5 | 52.3 | 43.4 | 13.3 | 55.5 | 54.8 | 60.7 | 53.6 | 45.3 | 13.3 |
| Adaptation | 59.3 | 56.7 | 66.7 | 61.0 | 53.7 | 18.4 | 58.6 | 56.6 | 65.5 | 59.8 | 52.8 | 18.2 |
| **Mamba Twister** | **61.1** | **58.0** | **68.4** | **62.8** | **55.2** | **18.5** | **61.2** | **59.0** | **69.0** | **63.4** | **55.5** | **18.7** |

Surprisingly, despite widely used in transformers, the **Attention-based Conditioning** performs poorly in our task. This suggests that the cross-attention mechanism is not inherently suitable for Mamba architecture. This may be due to *a fundamental discrepancy* between the two systems. To be specific, **(1)** the Mamba model is predicted on ordered sequences that exhibit *strict sequential dependencies*, where the state at a given time point $t + 1$ is determined by the preceding time point $t$ and a hidden state; **(2)** In contrast, cross-attention mechanisms treats all tokens within a sequence *equally*. This discrepancy may undermine the Mamba model's ability to structurally model sequences, as the cross-attention mechanism does not preserve the sequential integrity and the hierarchical dependencies essential for the model's operation.

Serving as a straightforward baseline, **In-context Conditioning** performs sub-optimal. This may because it models image-text interaction in an implicit way. In our situation, the length of the image feature is much larger than the textual side, textual information may be diluted during the forward process. Besides, this operation does not differentiate the textual information between image patches, and would be insufficient for multi-modal fusion.

Finally, **Norm Adaptation** with explicit image-text interaction modeling appears to be a strong baseline. It achieves better performance than other two variants. However, when calculating scale and bias, it uses only a pooled vector as

textual representation instead of the entire sequence. This potentially results in the loss of information, making it less effective compared to our Mamba Twister.

**Case-study Using Attention Maps.**    In Fig. 4, we further visualized the attention maps in Attention-based Conditioning as well as the local interaction map in Mamba Twister defined in Eq. (5). The four images in Fig. 4(b) are arranged from left to right with the network going from shallow to deep. It can be observed that in the shallower layers, Mamba Twister focuses more on the lower-level features of the images such as edges. As the layer goes deeper, the cost map mostly concentrating on the target object, indicating that Mamba Twister can progressively guide the image feature towards the target described by the language. In contrast, though the Attention-based Conditioning variant is able to locate the target, its attention map performs poorly. This may suggest that the cross-attention mechanism struggles to accurately capture the correct context information, and the two modalities are not truly fused together.
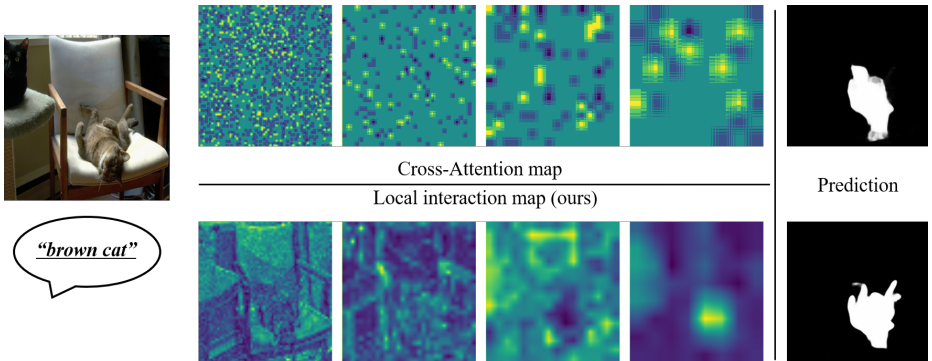


**Fig. 4: Cross-Attention map (up) and our local interaction map (down) comparison.**    Though both methods are able to predict target correctly, the cross-attention maps don't show correct image-text correlation, while ours are able to capture this relationship accurately, indicating that Mamba Twister is able to gradually fusing the two modality.

### 4.4    Comparison with the State-of-the-Arts

Tab. 2 shows the comparison of our **ReMamber** with state-of-the-art methods on RefCOCO, RefCOCO+, and G-Ref datasets. **ReMamber** outperforms all other methods across all datasets, demonstrating its superior capability in capturing and integrating contextual information for more accurate segmentation. In particular, Mamba outperforms previous Swin-based methods such as LAVT [56] by a considerable margin, indicating the remarkable capability of the Mamba-based architecture in segmentation tasks. and the efficiency of the newly introduced Mamba Twister.

**Table 2: Comparison with state-of-the-art methods on RefCOCO, Ref-COCO+, and G-Ref.** The metric of oIoU is reported. Best results are in bold.

| Model | Backbone | RefCOCO | | | RefCOCO+ | | | G-Ref | |
|---|---|---|---|---|---|---|---|---|---|
| | | val | testA | testB | val | testA | testB | valU | testU |
| PCAN [1] | ResNet-50 | 69.51 | 71.64 | 64.18 | 58.25 | 63.68 | 48.89 | 59.98 | 60.8 |
| MAttNet [58] | MaskRCNN ResNet-101 | 56.51 | 62.37 | 51.7 | 46.67 | 52.39 | 40.08 | 47.64 | 48.61 |
| RMI [29] | DeepLab ResNet-101 | 45.18 | 45.69 | 45.57 | 29.86 | 30.48 | 29.50 | - | - |
| RRN [25] | DeepLab ResNet-101 | 55.33 | 57.26 | 53.95 | 39.75 | 42.15 | 36.11 | - | - |
| CMSA [57] | DeepLab ResNet-101 | 58.32 | 60.61 | 55.09 | 43.76 | 47.6 | 37.89 | - | - |
| CAC [3] | DeepLab ResNet-101 | 58.90 | 61.77 | 53.81 | - | - | - | 46.37 | 46.95 |
| STEP [2] | DeepLab ResNet-101 | 60.04 | 63.46 | 57.97 | 48.19 | 52.33 | 40.41 | - | - |
| BRINet [16] | DeepLab ResNet-101 | 60.98 | 62.99 | 59.21 | 48.17 | 52.32 | 42.11 | - | - |
| CMPC [17] | DeepLab ResNet-101 | 61.36 | 64.53 | 59.64 | 49.56 | 53.44 | 43.23 | - | - |
| LSCM [18] | DeepLab ResNet-101 | 61.47 | 64.99 | 59.55 | 49.34 | 53.12 | 43.50 | - | - |
| CMPC+ [32] | DeepLab ResNet-101 | 62.47 | 65.08 | 60.82 | 50.25 | 54.04 | 43.47 | - | - |
| BUSNet [54] | DeepLab ResNet-101 | 63.27 | 66.41 | 61.39 | 51.76 | 56.87 | 44.13 | - | - |
| CGAN [34] | DeepLab ResNet-101 | 64.86 | 68.04 | 62.07 | 51.03 | 55.51 | 44.06 | 51.01 | 51.69 |
| EFN [6] | Wide ResNet-101 | 62.76 | 65.69 | 59.67 | 51.50 | 55.24 | 43.01 | - | - |
| ETRIS [53] | CLIP ResNet-101 | 71.06 | 74.11 | 66.66 | 62.23 | 68.51 | 52.79 | 60.28 | 60.42 |
| ReSTR [24] | ViT-B-16 | 67.22 | 69.3 | 64.45 | 55.78 | 60.44 | 48.27 | 54.48 | - |
| ETRIS [53] | ViT-B-16 | 70.51 | 73.51 | 66.63 | 60.10 | 66.89 | 50.17 | 59.82 | 59.91 |
| CRIS [51] | CLIP ResNet-50 | 69.52 | 72.72 | 64.7 | 61.39 | 67.1 | 52.48 | 59.87 | 60.36 |
| CRIS [51] | CLIP ResNet-101 | 70.47 | 73.18 | 66.10 | 62.27 | 68.08 | 53.68 | 59.87 | 60.36 |
| LAVT [56] | Swin-B | 72.73 | 75.82 | 68.79 | 62.14 | 68.38 | 55.10 | 61.24 | 62.10 |
| **ReMamber** (Ours) | Mamba-B | **74.54** | **76.74** | **70.89** | **65.00** | **70.78** | **57.53** | **63.9** | **64.0** |

**Table 3: Ablation on the combination of two scans.** We ablated each scan separately, as well as combining them parallel or swap the order. The occurence of Spatial Scan affect most to the performance. The three variants performs similarly, with the combination of Spatial-Channel slightly better. Results are evaluated on RefCOCO.

| Channel Scan | Spatial Scan | val | | | | testA | | | | testB | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | mIoU | oIoU | Pr@50 | Pr@70 | mIoU | oIoU | Pr@50 | Pr@70 | mIoU | oIoU | Pr@50 | Pr@70 |
| ✓ | - | 62.3 | 60.1 | 70.7 | 57.2 | 64.9 | 63.2 | 74.3 | 61.9 | 59.7 | 57.5 | 66.3 | 51.9 |
| - | ✓ | 70.0 | 65.8 | 80.5 | 69.1 | 72.3 | 71.0 | 83.7 | 73.4 | 67.5 | 63.5 | 75.2 | 62.7 |
| Parallel | | 71.0 | 68.8 | 81.9 | 70.0 | 73.1 | **71.8** | 84.7 | 74.1 | 67.9 | 64.6 | 76.9 | 64.5 |
| Spatial-Channel | | 71.4 | **68.8** | **82.4** | 70.8 | 73.3 | 71.4 | **84.9** | 74.2 | 67.9 | **64.7** | 76.5 | 64.4 |
| Channel-Spatial | | **71.6** | 68.4 | 82.1 | **70.9** | **73.3** | 71.5 | 84.8 | **74.5** | **68.4** | 64.5 | **77.1** | **64.9** |

## 4.5   Ablation Study

**Effects of Two Scans.**   Tab. 3 presents a comparative analysis of different scan variants. We provide outcomes of conducting Channel Scan and Spatial Scan independently. Additionally, results from various combinations of these two scans are discussed. The term "Parallel" denotes the concurrent execution of both scans, followed by adding the independent output for fusion. "Channel-Spatial" refers to the original scanning order in Mamba Twister shown in Fig. 2, where the Channel Scan is executed first, followed by the Spatial Scan. "Spatial-Channel" refers to reversing the sequence of the two scans, starting with Spatial Scan and then proceeding to Channel Scan.

Tab. 3 indicates that utilizing any single scan on its own yields suboptimal results, with the use of Channel Scan alone experiencing the most significant drop in effectiveness. This decline may indicate that a pure Channel Scan alters the data distribution, adversely affecting network stability. Among the three

hybrid scanning policies, each has its strengths and weaknesses. Overall, the combination of Channel-Spatial Scan appears to offer a considerable advantage.

**Distribution Visualization.**    The visualization of the distribution of the multi-modal input data is depicted in Fig. 5, where we use PCA to map the feature to a 3D space. The **image** and **text** data are indicated in red and blue respectively.

Fig. 5a illustrates the distribution of the input data. We can observe that the text data forms a linear arrangement within the 3D space, as it is duplicated $H \times W$ times to align with the image size. This suggests that the text instruction should be applied for each pixel within the image. Fig. 5b reveals the data distribution following the application of a Channel Scan. Notably, this process appears to lean towards aggregating different modalities towards a distribution pattern similar with text. Fig. 5c demonstrates the data distribution after a Spatial Scan. Here, the features of the textual and image data disperse, with the modalities intermixed. The Spatial Scan thus reintegrates the previously aligned modalities, distributing them in a manner that reflects their combined influence.
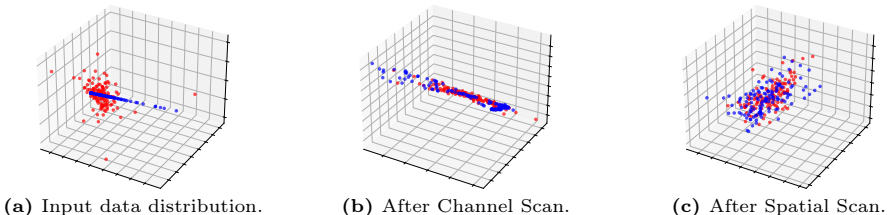


(a) Input data distribution.          (b) After Channel Scan.          (c) After Spatial Scan.

Fig. 5: **Data distribution after Channel Scan and Spacial Scan. Image** in red and **text** data in blue. The Channel Scan tends to aggregate different modalities towards the distribution of textual side. The Spatial Scan reintegrates the previously aligned modalities, distributing them in a manner that reflects their combined influence.

Table 4: **Ablation study on global and local interaction.** Both global and local interactions are crucial for modality fusion. Results are evaluated on RefCOCO.

| Image | Global | Local | val | | | | testA | | | | testB | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | mIoU | oIoU | Pr@50 | Pr@70 | mIoU | oIoU | Pr@50 | Pr@70 | mIoU | oIoU | Pr@50 | Pr@70 |
| ✓ | ✓ | - | 69.1 | 66.6 | 79.5 | 66.3 | 71.4 | 69.8 | 82.3 | 70.9 | 66.3 | 63.7 | 74.7 | 61.0 |
| ✓ | - | ✓ | 69.9 | 67.9 | 80.6 | 68.2 | 72.2 | 70.6 | 83.9 | 73.1 | 66.4 | 63.8 | 75.3 | 61.7 |
| ✓ | ✓ | ✓ | **71.6** | **68.4** | **82.1** | **70.9** | **73.3** | **71.5** | **84.8** | **74.5** | **68.4** | **64.5** | **77.1** | **64.9** |

**Effects of Global and Local Features.** We also ablate the effect of global and local interactions when forming the hybrid feature cube, formally, the $\tilde{\mathbf{F}}_t$ and $\tilde{\mathbf{F}}_c$ in Eq. (6). The results in Tab. 4 illustrate that the integration of both global and local features significantly enhances the performance.
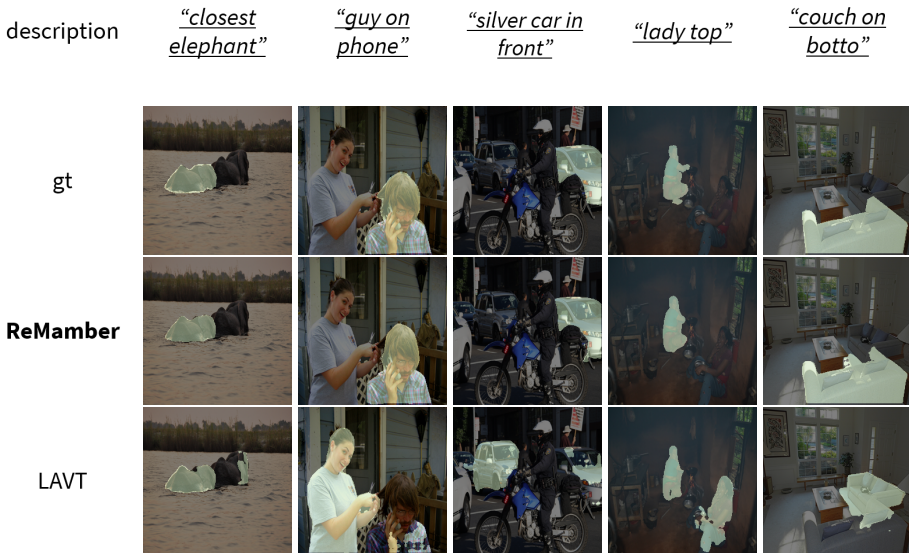
# 5    Visualization Results



**Fig. 6:** Visualization results of our **ReMamber** and the baseline model LAVT. Our model is able to predict more accurate masks.

Fig. 6 presents the visualization results of our method compared with the baseline method LAVT [56]. Fig. 7 presents the visualization outcomes of our method alongside three other variants. Our **ReMamber** is capable of producing segmentation results with higher accuracy. In contrast, the other three variants occasionally encounter issues with inaccurate segmentation masks or are misled towards incorrect objects.

# 6    Conclusion

In this study, we introduce **ReMamber**, a novel architecture utilizing the Mamba framework in Referring Image Segmentation (RIS). It marks a significant advancement in multi-modal understanding. By integrating visual and textual information through innovative Mamba Twister blocks, our approach sets new benchmarks in image-text modality fusion. Achieving competitive results across multiple RIS datasets, our research highlights the potential of Mamba architecture in enhancing the scalability and performance of multi-modal tasks, offering promising directions for future exploration in the field.

**Limitations and Future Works.** In our current architecture, the segmentation decoder is constructed by only a few convolutional layers. As shown
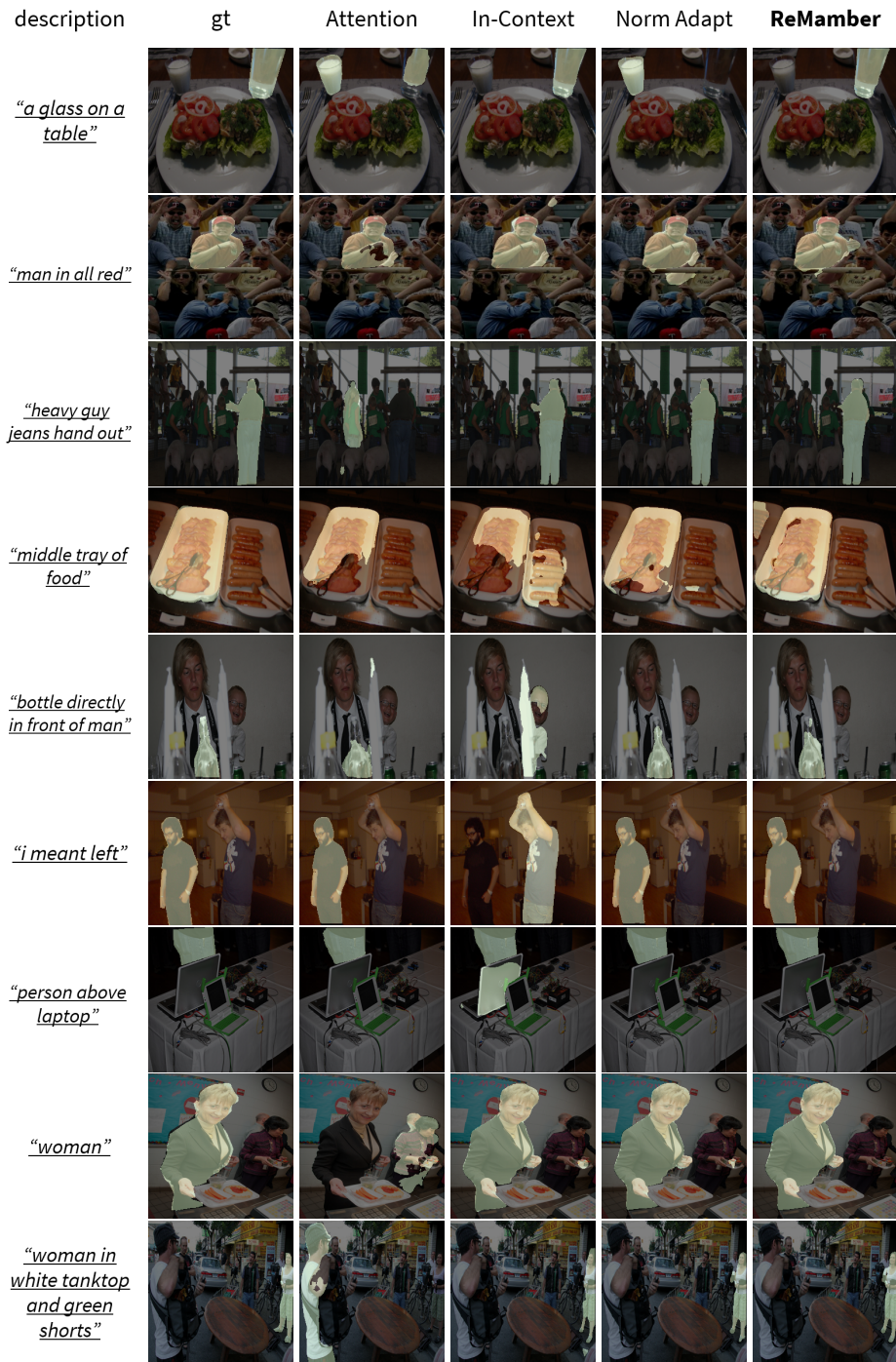
**Fig. 7: Visualization results.** Our **ReMamber** is capable of producing segmentation results with higher accuracy. While other three variants occasionally encounter issues with inaccurate segmentation masks or are misled towards incorrect objects.

in Sec. 4.3, the integration of the Cross-Attention mechanism within the Mamba-based architecture demonstrates sub-optimal compatibility, undermining the overall efficacy of our method. In light of these findings, future endeavors will be directed towards the investigation of more sophisticated multi-modal segmentation decoders which best fit Mamba architecture.

# Appendix

## A    Speed Analysis

Here, we supplement two experiments *w.r.t.* model FPS and memory cost in Fig. 8 under different resolutions. ReMamber is consistently faster and requires fewer memory cost than LAVT, especially with large resolution (*e.g.*, 1,024).



**Fig. 8:** Comparison of inference FPS (left) and training GPU memory (right) between LAVT and our **ReMamber**.

## B    Implementation Details

Here we provide more details about the implementation of our method, including detailed architectural structure, training settings and other baseline implementation.

### B.1    Architecture Details

The whole **ReMamber** architecture consists of an encoder and a decoder. The encoder is made up by a patch-embedding layer with patch-size 4 and hidden dimension 128, followed by 4 Mamba Twister blocks. Each Twister block consists of several VSS Layers and a Twisting Layer. The VSS Layer number configuration is set as 2-2-15-2, with hidden dimension 128-256-512-1024, respectively.

For the decoder part, we provide two variants in our code implementation: convolution-based decoder (`ReMamber_Conv`) and Mamba-based decoder

(`ReMamber_Mamba`). `ReMamber_Conv` uses a progressive upsampling architecture with 4 residual blocks, 2 convolutional layers in each. `ReMamber_Mamba` is similar with `ReMamber_Conv`, but uses VSS layers instead of convolutional layers. This variant is slightly faster.

## B.2   Implementation for Other Three Variants

Here, we detail the implementation of other three architecture variants in our paper, *i.e.*, In-context Conditioning, Attention-based Conditioning and Norm Adaptation.

**In-context Conditioning.**   To enable the model to distinguish between two modalities, we add learnable positional embeddings to the image and text tokens separately at each layer before the Spatial SSM.

**Attention-based Conditioning.**   In this variant, we also incorporate learnable positional embeddings. For the cross-attention block, we use image tokens as the query, and text tokens as the key and value.

**Norm Adaptation.**   Norm Adaptation learns a global scale and bias. First, we use an MLP layer to pool a global vector from the text. This vector is then used to scale and bias the image tokens. An additional feed-forward layer is added after adjusting the scale and bias to maintain parameter size comparable to other variants.

# References

1. Chen, B., Hu, Z., Ji, Z., Bai, J., Zuo, W.: Position-aware contrastive alignment for referring image segmentation. arXiv preprint arXiv:2212.13419 (2022) 12
2. Chen, D.J., Jia, S., Lo, Y.C., Chen, H.T., Liu, T.L.: See-through-text grouping for referring image segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7454–7463 (2019) 3, 12
3. Chen, Y.W., Tsai, Y.H., Wang, T., Lin, Y.Y., Yang, M.H.: Referring expression object segmentation with caption-aware consistency. arXiv preprint arXiv:1910.04748 (2019) 3, 12
4. Ding, H., Liu, C., Wang, S., Jiang, X.: Vision-language transformer and query generation for referring segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 16321–16330 (2021) 1, 2, 3
5. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020) 2
6. Feng, G., Hu, Z., Zhang, L., Lu, H.: Encoder fusion network with co-attention embedding for referring image segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 15506–15515 (2021) 3, 12
7. Fu, D.Y., Dao, T., Saab, K.K., Thomas, A.W., Rudra, A., Ré, C.: Hungry hungry hippos: Towards language modeling with state space models. arXiv preprint arXiv:2212.14052 (2022) 2, 4
8. Gu, A., Dao, T.: Mamba: Linear-time sequence modeling with selective state spaces. arXiv preprint arXiv:2312.00752 (2023) 2, 4, 5, 7, 9
9. Gu, A., Dao, T., Ermon, S., Rudra, A., Ré, C.: Hippo: Recurrent memory with optimal polynomial projections. Advances in neural information processing systems **33**, 1474–1487 (2020) 4
10. Gu, A., Goel, K., Gupta, A., Ré, C.: On the parameterization and initialization of diagonal state space models. Advances in Neural Information Processing Systems **35**, 35971–35983 (2022) 4
11. Gu, A., Goel, K., Ré, C.: Efficiently modeling long sequences with structured state spaces. arXiv preprint arXiv:2111.00396 (2021) 2, 4
12. Gu, A., Johnson, I., Goel, K., Saab, K., Dao, T., Rudra, A., Ré, C.: Combining recurrent, convolutional, and continuous-time models with linear state space layers. Advances in neural information processing systems **34**, 572–585 (2021) 4
13. Gupta, A., Gu, A., Berant, J.: Diagonal state spaces are as effective as structured state spaces. Advances in Neural Information Processing Systems **35**, 22982–22994 (2022) 4
14. He, X., Cao, K., Yan, K., Li, R., Xie, C., Zhang, J., Zhou, M.: Pan-mamba: Effective pan-sharpening with state space model. arXiv preprint arXiv:2402.12192 (2024) 2, 4
15. Hu, R., Rohrbach, M., Darrell, T.: Segmentation from natural language expressions. In: Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14. pp. 108–124. Springer (2016) 1, 3
16. Hu, Z., Feng, G., Sun, J., Zhang, L., Lu, H.: Bi-directional relationship inferring network for referring image segmentation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4424–4433 (2020) 3, 12

17. Huang, S., Hui, T., Liu, S., Li, G., Wei, Y., Han, J., Liu, L., Li, B.: Referring image segmentation via cross-modal progressive comprehension. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10488–10497 (2020) 3, 12

18. Hui, T., Liu, S., Huang, S., Li, G., Yu, S., Zhang, F., Han, J.: Linguistic structure guided context modeling for referring image segmentation. In: European Conference on Computer Vision. pp. 59–75. Springer (2020) 12

19. Islam, M.M., Bertasius, G.: Long movie clip classification with state-space video models. In: European Conference on Computer Vision. pp. 87–104. Springer (2022) 4

20. Islam, M.M., Hasan, M., Athrey, K.S., Braskich, T., Bertasius, G.: Efficient movie scene detection using state-space transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18749–18758 (2023) 4

21. Jing, Y., Kong, T., Wang, W., Wang, L., Li, L., Tan, T.: Locate then segment: A strong pipeline for referring image segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9858–9867 (2021) 3

22. Kamath, A., Singh, M., LeCun, Y., Synnaeve, G., Misra, I., Carion, N.: Mdetr-modulated detection for end-to-end multi-modal understanding. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1780–1790 (2021) 1, 3

23. Kazemzadeh, S., Ordonez, V., Matten, M., Berg, T.: ReferItGame: Referring to objects in photographs of natural scenes. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP) (2014) 9

24. Kim, N., Kim, D., Lan, C., Zeng, W., Kwak, S.: Restr: Convolution-free referring image segmentation using transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18145–18154 (2022) 1, 2, 3, 12

25. Li, R., Li, K., Kuo, Y.C., Shu, M., Qi, X., Shen, X., Jia, J.: Referring image segmentation via recurrent refinement networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5745–5753 (2018) 1, 3, 12

26. Liang, D., Zhou, X., Wang, X., Zhu, X., Xu, W., Zou, Z., Ye, X., Bai, X.: Point-mamba: A simple state space model for point cloud analysis. arXiv preprint arXiv:2402.10739 (2024) 2, 4

27. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: Common Objects in Context. ECCV (2014) 9

28. Liu, C., Ding, H., Jiang, X.: Gres: Generalized referring expression segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 23592–23601 (2023) 3

29. Liu, C., Lin, Z., Shen, X., Yang, J., Lu, X., Yuille, A.: Recurrent multimodal interaction for referring image segmentation. In: Proceedings of the IEEE international conference on computer vision. pp. 1271–1280 (2017) 1, 3, 12

30. Liu, J., Ding, H., Cai, Z., Zhang, Y., Satzoda, R.K., Mahadevan, V., Manmatha, R.: Polyformer: Referring image segmentation as sequential polygon generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18653–18663 (2023) 3

31. Liu, J., Yang, H., Zhou, H.Y., Xi, Y., Yu, L., Yu, Y., Liang, Y., Shi, G., Zhang, S., Zheng, H., et al.: Swin-umamba: Mamba-based unet with imagenet-based pre-training. arXiv preprint arXiv:2402.03302 (2024) 2, 4

32. Liu, S., Hui, T., Huang, S., Wei, Y., Li, B., Li, G.: Cross-modal progressive comprehension for referring segmentation. IEEE Transactions on Pattern Analysis and Machine Intelligence **44**(9), 4761–4775 (2021) 12

33. Liu, Y., Tian, Y., Zhao, Y., Yu, H., Xie, L., Wang, Y., Ye, Q., Liu, Y.: Vmamba: Visual state space model. arXiv preprint arXiv:2401.10166 (2024) 2, 4, 6, 9

34. Luo, G., Zhou, Y., Ji, R., Sun, X., Su, J., Lin, C.W., Tian, Q.: Cascade grouped attention network for referring expression segmentation. Proceedings of the 28th ACM International Conference on Multimedia (2020) 12

35. Ma, C., Yang, Y., Ju, C., Zhang, F., Zhang, Y., Wang, Y.: Attrseg: Open-vocabulary semantic segmentation via attribute decomposition-aggregation. NeurIPS (2023) 1

36. Ma, C., Yang, Y., Wang, Y., Zhang, Y., Xie, W.: Open-vocabulary semantic segmentation with frozen vision-language models. In: British Machine Vision Conference (2022) 1

37. Ma, J., Li, F., Wang, B.: U-mamba: Enhancing long-range dependency for biomedical image segmentation. arXiv preprint arXiv:2401.04722 (2024) 2, 4

38. Mao, J., Huang, J., Toshev, A., Camburu, O.M., Yuille, A., Murphy, K.: Generation and comprehension of unambiguous object descriptions. CVPR (2015) 9

39. Margffoy-Tuay, E., Pérez, J.C., Botero, E., Arbeláez, P.: Dynamic multimodal instance segmentation guided by natural language queries. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 630–645 (2018) 3

40. Mehta, H., Gupta, A., Cutkosky, A., Neyshabur, B.: Long range language modeling via gated state spaces. arXiv preprint arXiv:2206.13947 (2022) 4

41. Nagaraja, V.K., Morariu, V.I., Davis, L.S.: Modeling context between objects for referring expression understanding. In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) ECCV (2016) 9

42. Perez, E., Strub, F., de Vries, H., Dumoulin, V., Courville, A.C.: Film: Visual reasoning with a general conditioning layer. In: AAAI (2018) 8

43. Qu, M., Wu, Y., Wei, Y., Liu, W., Liang, X., Zhao, Y.: Learning to segment every referring object point by point. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2023) 3

44. Ruan, J., Xiang, S.: Vm-unet: Vision mamba unet for medical image segmentation. arXiv preprint arXiv:2402.02491 (2024) 2, 4

45. Shi, H., Li, H., Meng, F., Wu, Q.: Key-word-aware network for referring expression image segmentation. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 38–54 (2018) 3

46. Smith, J.T., Warrington, A., Linderman, S.W.: Simplified state space layers for sequence modeling. arXiv preprint arXiv:2208.04933 (2022) 2, 4

47. Tang, J., Zheng, G., Shi, C., Yang, S.: Contrastive grouping with transformer for referring image segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 23570–23580 (2023) 3

48. Tang, J., Zheng, G., Shi, C., Yang, S.: Contrastive grouping with transformer for referring image segmentation. In: CVPR (2023) 9

49. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems **30** (2017) 2

50. Wang, J., Zhu, W., Wang, P., Yu, X., Liu, L., Omar, M., Hamid, R.: Selective structured state-spaces for long-form video understanding. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6387–6397 (2023) 4

51. Wang, Z., Lu, Y., Li, Q., Tao, X., Guo, Y., Gong, M., Liu, T.: Cris: Clip-driven referring image segmentation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 11686–11695 (2022) 1, 3, 12
52. Xing, Z., Ye, T., Yang, Y., Liu, G., Zhu, L.: Segmamba: Long-range sequential modeling mamba for 3d medical image segmentation. arXiv preprint arXiv:2401.13560 (2024) 2, 4
53. Xu, Z., Chen, Z., Zhang, Y., Song, Y., Wan, X., Li, G.: Bridging vision and language encoders: Parameter-efficient tuning for referring image segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 17503–17512 (2023) 12
54. Yang, S., Xia, M., Li, G., Zhou, H.Y., Yu, Y.: Bottom-up shift and reasoning for referring image segmentation. In: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 11261–11270 (2021) 12
55. Yang, Y., Ma, C., Ju, C., Zhang, Y., Wang, Y.: Multi-modal prototypes for open-set semantic segmentation. IJCV (2024) 1
56. Yang, Z., Wang, J., Tang, Y., Chen, K., Zhao, H., Torr, P.H.: Lavt: Language-aware vision transformer for referring image segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18155–18165 (2022) 1, 2, 3, 9, 11, 12, 14
57. Ye, L., Rochan, M., Liu, Z., Wang, Y.: Cross-modal self-attention network for referring image segmentation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10502–10511 (2019) 3, 12
58. Yu, L., Lin, Z., Shen, X., Yang, J., Lu, X., Bansal, M., Berg, T.L.: Mattnet: Modular attention network for referring expression comprehension. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018) 12
59. Zhang, T., Li, X., Yuan, H., Ji, S., Yan, S.: Point could mamba: Point cloud learning via state space model. arXiv preprint arXiv: 2403.00762 (2024) 2, 4
60. Zheng, Z., Wu, C.: U-shaped vision mamba for single image dehazing. arXiv preprint arXiv:2402.04139 (2024) 2, 4
61. Zhu, C., Zhou, Y., Shen, Y., Luo, G., Pan, X., Lin, M., Chen, C., Cao, L., Sun, X., Ji, R.: Seqtr: A simple yet universal network for visual grounding. In: European Conference on Computer Vision. pp. 598–615. Springer (2022) 3
62. Zhu, L., Liao, B., Zhang, Q., Wang, X., Liu, W., Wang, X.: Vision mamba: Efficient visual representation learning with bidirectional state space model. arXiv preprint arXiv:2401.09417 (2024) 2, 4