

# MMSFormer: Multimodal Transformer for Material and Semantic Segmentation

MD KAYKOBAD REZA <sup>1</sup>, ASHLEY PRATER-BENNETTE <sup>2</sup>, AND M. SALMAN ASIF <sup>1</sup> (Senior Member, IEEE)

<sup>1</sup>University of California, Riverside, CA 92508 USA

<sup>2</sup>Air Force Research Laboratory, Rome, NY 13441 USA

CORRESPONDING AUTHOR: M. SALMAN ASIF (email: sasif@ucr.edu).

This work was supported in part by AFOSR award FA9550-21-1-0330 and ONR award N00014-19-1-2264.  
The code and pretrained models will be made available at <https://github.com/csiplab/MMSFormer>.

**ABSTRACT** Leveraging information across diverse modalities is known to enhance performance on multimodal segmentation tasks. However, effectively fusing information from different modalities remains challenging due to the unique characteristics of each modality. In this paper, we propose a novel fusion strategy that can effectively fuse information from different modality combinations. We also propose a new model named Multi-Modal Segmentation TransFormer (MMSFormer) that incorporates the proposed fusion strategy to perform multimodal material and semantic segmentation tasks. MMSFormer outperforms current state-of-the-art models on three different datasets. As we begin with only one input modality, performance improves progressively as additional modalities are incorporated, showcasing the effectiveness of the fusion block in combining useful information from diverse input modalities. Ablation studies show that different modules in the fusion block are crucial for overall model performance. Furthermore, our ablation studies also highlight the capacity of different input modalities to improve performance in the identification of different types of materials.

**INDEX TERMS** Multimodal fusion, multimodal image segmentation, material segmentation, semantic segmentation, transformer.

## I. INTRODUCTION

Image segmentation [1], [2] methods assign one class label to each pixel in an image. The segmentation map can be used for holistic understanding of objects or context of the scene. Image segmentation can be further divided into different types; examples include semantic segmentation [3], [4], instance segmentation [5], [6], panoptic segmentation [7], [8] and material segmentation [9], [10]. Each of these segmentation tasks are designed to address specific challenges and applications.

Multimodal image segmentation [11], [12] aims to enhance the accuracy and completeness of the task by leveraging diverse sources of information, and potentially leading to a more robust understanding of complex scenes. In contrast to single-modal segmentation [2], the multimodal approach [12] is more complex due to the necessity of effectively integrating heterogeneous data from different modalities. Key challenges arise from variations in data quality and attributes, distinct

traits of each modality, and need to create models capable of accurately and coherently segmenting with the fused information.

Most of the existing multimodal segmentation methods are designed to work with specific modality pairs, such as RGB-Depth [13], [14], [15], RGB-Thermal [16], [17], [18], and RGB-Lidar [19], [20], [21]. As they are designed for specific modality combinations, most of them generally do not work well with modality combinations different from the ones used in the original design. Recently, CMX [22] introduced a technique to fuse information from RGB and one other supplementary modality, but it is incapable of fusing more than two modalities at the same time. Some recent models have proposed techniques to fuse more than two modalities [9], [23], [24]. However, they either use very complex fusion strategies [22], [23] or require additional information like semantic labels [9] for performing underlying tasks.

In this paper, we propose a novel fusion block that can fuse information from diverse combination of modalities. We also propose a new model for multimodal material and semantic segmentation tasks that we call MMSFormer. Our model uses transformer based encoders [25] to capture hierarchical features from different modalities, fuses the extracted features with our novel fusion block and utilizes MLP decoder to perform multimodal material and semantic segmentation. In particular, our proposed fusion block uses parallel convolutions to capture multi-scale features, channel attention to re-calibrate features along the channel dimension and linear layer to combine information across multiple modalities. Such a design provides a simple and computationally efficient fusion block that can handle an arbitrary number of input modalities and combine information effectively from different modality combinations. An illustration of the proposed method is presented in Fig. 1. We compare our fusion block with some of the existing fusion methods in terms of number of parameters and GFLOPs in Table 9.

To evaluate our proposed MMSFormer and fusion block, we focus on multimodal material segmentation on MCubeS [9] dataset and multimodal semantic segmentation on FMB [27] and PST900 [28] datasets. MCubeS dataset consists of four different modalities: RGB, angle of linear polarization (AoLP), degree of linear polarization (DoLP) and near-infrared (NIR). FMB dataset includes RGB and infrared modalities, while PST900 dataset comprises RGB and thermal modalities. We show the overall and per-class performance comparison in Table 1–5 for these datasets. A series of experiments highlight the ability of the proposed fusion block to effectively combine features from different modality combinations, resulting in superior performance compared to current state-of-the-art methods. Ablation studies show that different input modalities assist in identifying different types of material classes as shown in Table 8. Furthermore, as we add new input modalities, overall performance increases gradually highlighting the ability of the fusion block to incorporate useful information from new modalities. We summarize the results in Tables 4 and 6 for FMB and MCubeS datasets respectively.

Main contributions of this paper can be summarized as follows.

- We propose a new multimodal segmentation model called MMSFormer. The model incorporates a novel fusion block that can fuse information from arbitrary (heterogeneous) combinations of modalities.
- Our model achieves new state-of-the-art performance on three different datasets. Furthermore, our method achieves better performance for all modality combinations compared to the current leading models.
- A series of ablation studies show that each module on the fusion block has an important contribution towards the overall model performance and each input modality assists in identifying specific material classes.

Rest of the paper is structured as follows. Section II presents a brief review of related work. We describe our model

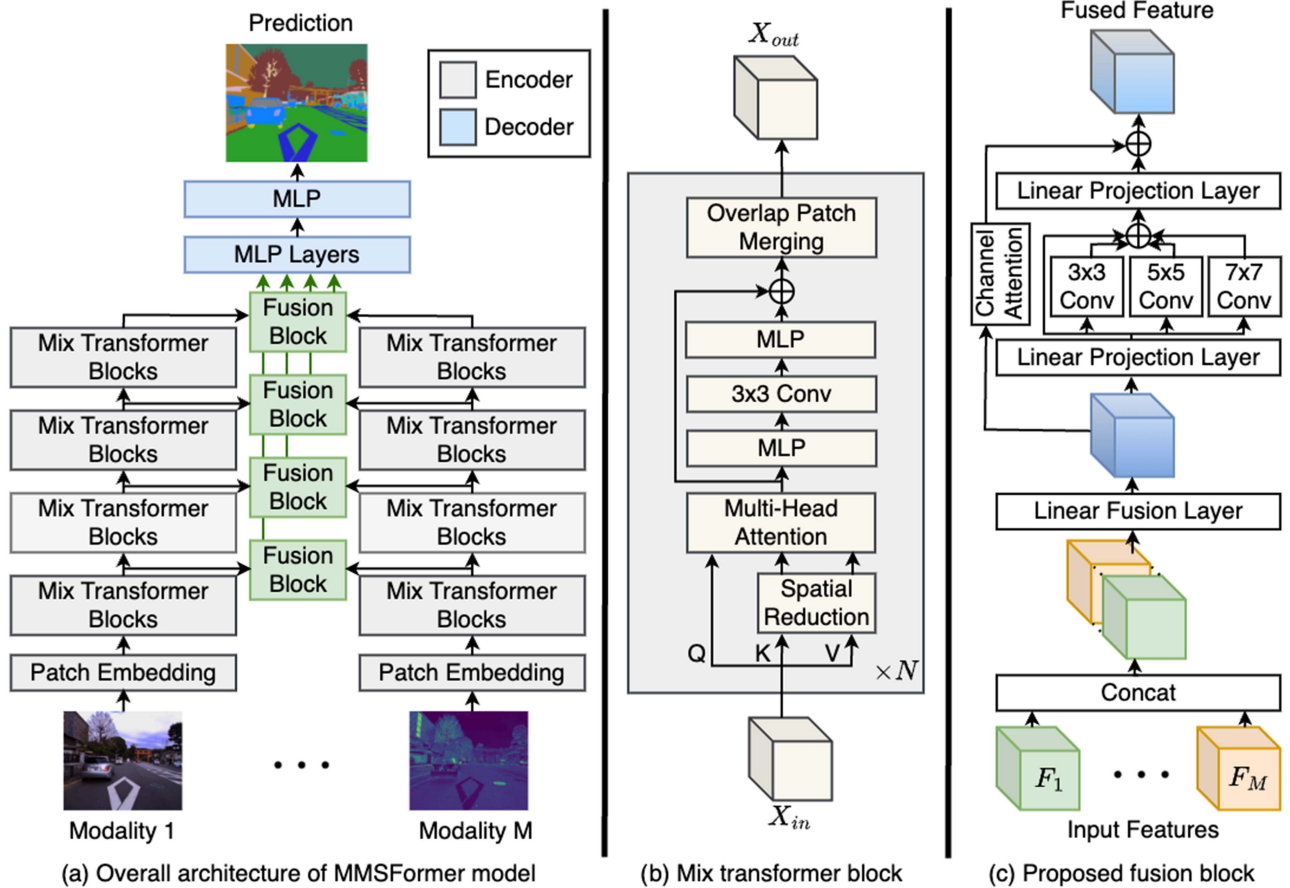
and fusion block in detail in Section III. Section IV presents experimental results and ablation studies on multimodal material and semantic segmentation tasks with qualitative and quantitative analysis.

## II. RELATED WORK

Image segmentation has witnessed significant evolution, spurred by advancements in machine learning and computational capabilities. A significant improvement in this evolution came with the inception of fully convolutional networks (FCNs) [29], [30], which enabled pixel-wise predictions through the utilization of hierarchical features within convolutional neural networks (CNNs). This led to the development of a variety of CNN based models for different images segmentation tasks. U-Net [31] is one such model that utilizes skip-connections between the lower-resolution and corresponding higher-resolution feature maps. DeepLabV3+ [32] introduced dilated convolutions (atrous convolutions) into the encoder allowing the expansion of the receptive field without increasing computational complexity significantly. PSPNet [33] introduced global context modules that enable the model to gather information from a wide range of spatial scales, essentially integrating both local and global context into the segmentation process.

Recently, Transformer based models have proven to be very effective in handling complex image segmentation tasks. Some of the notable transformer-based models are Pyramid Vision Transformer (PVT) [34], SegFormer [25], and Mask2Former [35]. PVT [34] utilizes transformer based design for various computer vision tasks. SegFormer [25] utilizes efficient self-attention and lightweight MLP decoder for simple and efficient semantic segmentation. Mask2Former [35] uses masked-attention along with pixel decoder and transformer decoder for any segmentation task. Their success demonstrates the capacity of these models to provide state-of-the-art solutions in various segmentation tasks.

In the context of multimodal image segmentation, fusion of data from diverse sources [12] has gained traction as a means to extract richer information and improve accuracy. A variety of models and fusion strategies have been proposed for RGB-Depth segmentation tasks. FuseNet [14] model integrates depth feature maps into RGB feature maps, while SA-Gate [13] employs Separation-and-Aggregation Gating to mutually filter and recalibrate RGB and depth modalities before fusion. Attention Complementary Module has been proposed by AC-Net [15] that extracts weighted RGB and depth features for fusion. The domain of RGB-Thermal image segmentation has also gained prominence. Recent models include RTFNet [18] that achieves fusion through elementwise addition of thermal features with RGB, RSFNet [16] proposing Residual Spatial Fusion module to blend RGB and Thermal modalities, and EAEFNet [17] utilizing attention interaction and attention complement mechanisms to merge RGB and Thermal features. A number of methods also focus on fusing RGB-Lidar data that include TransFuser [20] that employs Transformer



**FIGURE 1.** (a) Overall architecture of MMSFormer model. Each image passes through a modality-specific encoder where we extract hierarchical features. Then we fuse the extracted features using the proposed fusion block and pass the fused features to the decoder for predicting the segmentation map. (b) Illustration of the mix transformer [25] block. Each block applies a spatial reduction before applying multi-head attention to reduce computational cost. (c) Proposed multimodal fusion block. We first concatenate all the features along the channel dimension and pass it through linear fusion layer to fuse them. Then the feature tensor is fed to linear projection and parallel convolution layers to capture multi-scale features. We use Squeeze and Excitation block [26] as channel attention in the residual connection to dynamically re-calibrate the features along the channel dimension.

blocks, whereas LIF-Seg [21] relies on coarse feature extraction, offset learning, and refinement for effective fusion.

While the previously mentioned studies focus on specific pairs of modalities, some recent research has demonstrated promising results in the fusion of arbitrary modalities. CMX [22] introduces cross-modal feature rectification and fusion modules to merge RGB features with supplementary modalities. For multimodal material segmentation, MCubeSNet [9] model is proposed, which can seamlessly integrate four different modalities to enhance segmentation accuracy. In the context of arbitrary modal semantic segmentation, CMNeXt [24] introduces Self-Query Hub and Parallel Pooling Mixer modules, offering a versatile approach for fusing diverse modalities. Additionally, HRFuser [23] employs multi-window cross-attention to fuse different modalities at various resolutions, thereby enriching model performance.

Though some of these models can fuse different modalities, they either use very complex fusion strategies [22], [23], [24] or requires additional information [9] to perform underlying task. We aim to design a simple fusion module that can handle

arbitrary number of input modalities and able to effectively fuse information from diverse modality combinations.

### III. PROPOSED MODEL

The overall architecture of our proposed MMSFormer model and the fusion block is shown in Fig. 1. The model has three modules: (1) Modality specific encoder; (2) Multimodal fusion block; and (3) Shared MLP decoder. We use mix transformer [25] as the encoder of our model. We choose mix transformer for various reasons. First, it can provide hierarchical features without positional encoding. Second, it uses spatial reduction before attention that reduces the number of parameters significantly [25], [34]. Third, it also works well with simple and lightweight MLP decoder [25].

#### A. OVERALL MODEL ARCHITECTURE

Our overall model architecture is shown in Fig. 1(a). Assume we have  $M$  distinct modalities. Given a set of modalities as input, each modality-specific encoder captures distinctive features from each input modality by mapping the corresponding

image into modality-specific hierarchical features as

$$F_m = \text{Encoder}_m(I_m), \quad (1)$$

where  $I_m \in \mathbb{R}^{H \times W \times 3}$  represents input image for modality  $m \in \{1, 2, \dots, M\}$  and  $\text{Encoder}_m(\cdot)$  denotes the encoder for that modality. The encoder generates four feature maps at  $\{\frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{32}\}$  of the input image resolution. We represent them as  $F_m = \{F_m^1, F_m^2, F_m^3, F_m^4\}$ . For simplicity we denote the shape of the feature map for the  $i$ th encoder stage as  $(H_i \times W_i \times C_i)$  where  $i \in \{1, 2, 3, 4\}$ .

We use four separate fusion blocks, one corresponding to each encoder stage, to fuse the features from each stage of the encoder. We pass the extracted features  $F_m^i$  for all modalities to the  $i$ th fusion block as

$$F^i = \text{FusionBlock}^i(\{F_m^i\}_m). \quad (2)$$

Each fusion block fuses the features extracted from all the modalities to generate a combined feature representation  $F = \{F^1, F^2, F^3, F^4\}$ , where  $F^i$  denotes the fused feature at  $i$ th stage. Finally, we pass the combined features  $F$  to the MLP decoder [25] to predict the segmentation labels.

## B. MODALITY SPECIFIC ENCODER

We use mix transformer encoder [25] to capture hierarchical features from the input modalities. Each input image  $I_m$  goes through patch embedding layer where it is divided into  $4 \times 4$  patches following [25] and then fed to the mix transformer blocks. The design of mix transformer block is shown in Fig. 1(b). We denote the input to any mix transformer block as  $X_{in} \in \mathbb{R}^{H_i \times W_i \times C_i}$  that is reshaped to  $N_i \times C_i$  (with  $N_i = H_i W_i$ ) and used as query  $Q$ , key  $K$ , and value  $V$ .

To reduce the computational overhead, spatial reduction is applied following [34] using a reduction ratio  $R$ .  $K$  and  $V$  are first reshaped into  $\frac{N_i}{R} \times C_i R$  matrices and then mapped to  $\frac{N_i}{R} \times C_i$  matrices via linear projection. A standard multi-head self-attention (MHSA) maps  $Q, K, V$  to intermediate features as

$$\begin{aligned} \text{MHSA}(Q, K, V) &= \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O, \\ \text{head}_j &= \text{Attention}\left(QW_j^Q, KW_j^K, VW_j^V\right), \end{aligned} \quad (3)$$

where  $h$  represents the number of attention heads,  $W_j^Q \in \mathbb{R}^{C_i \times d_K}$ ,  $W_j^K \in \mathbb{R}^{C_i \times d_K}$ ,  $W_j^V \in \mathbb{R}^{C_i \times d_V}$ , and  $W^O \in \mathbb{R}^{hd_V \times C_i}$  are the projection matrices,  $d_K, d_V$  represent dimensions of  $K, V$ , respectively. We can formulate the Attention function as

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)V, \quad (4)$$

where  $Q, K$ , and  $V$  are the input query, key, and value matrices. MHSA is followed by a mixer layer (with two MLP and one  $3 \times 3$  convolution layer). The convolution layer provides sufficient positional encoding into the transformer encoder for optimal segmentation performance [25]. This layer can be

written as

$$\begin{aligned} \hat{X}_{in} &= \text{MHSA}(Q, K, V), \\ X_{out} &= \text{MLP}(\text{GELU}(\text{Conv}_{3 \times 3}(\text{MLP}(\hat{X}_{in})))) + \hat{X}_{in}, \end{aligned} \quad (5)$$

Finally, overlap patch merging is applied to  $X_{out}$  following [25] to generate the final output.

## C. MULTIMODAL FUSION BLOCK

After extracting hierarchical features, we fuse them using our proposed fusion block. The fusion block shown in Fig. 1(c) is responsible for fusing the features extracted from the modality specific encoders. We have one fusion block for each of the four encoder stages. For the  $i$ th fusion block, let us assume the input feature maps are given as  $F_m^i \in \mathbb{R}^{H_i \times W_i \times C_i} \forall m \in \{1, 2, \dots, M\}$ . First, we concatenate the feature maps from  $M$  modalities along the channel dimension to get the combined feature map  $F^i \in \mathbb{R}^{H_i \times W_i \times MC_i}$ . Then we pass the features through a linear fusion layer that combines the features and reduces the channel dimension to  $C_i$ . Let us denote the resulting features as  $\hat{F}^i \in \mathbb{R}^{H_i \times W_i \times C_i}$ . We represent the operation as

$$\hat{F}^i = \text{Linear}(F_1^i || \dots || F_M^i). \quad (6)$$

Here  $||$  represents concatenation of features along the channel dimension and the linear layer takes an  $MC_i$  dimensional input and generates a  $C_i$  dimensional output.

After the linear fusion layer, we added a module for capturing and mixing multi-scale features. The module consists of two linear projection layers having parallel convolution layers in between them. First we apply a linear transformation on  $\hat{F}^i$  along the channel dimension by passing it through the first linear projection layer. It refines and tunes the features from different channels. Then we apply  $3 \times 3$ ,  $5 \times 5$ , and  $7 \times 7$  convolutions to effectively capture diverse features across different spatial contexts. By employing convolutions with different sizes, the fusion block can attend to local patterns as well as capture larger spatial structures, thereby enhancing its ability to extract meaningful features from the input data. Finally we apply another linear transformation along the channel dimension using the second linear project layer to consolidate the information captured by the parallel convolutions, promoting feature consistency and enhancing the discriminative power of the fused features. These steps can be performed as

$$\tilde{F}^i = \text{Linear}\left(\hat{F}^i\right), \quad (7)$$

$$F^i = \text{Linear}\left(\tilde{F}^i + \sum_{k \in \{3, 5, 7\}} \text{Conv}_{k \times k}(\tilde{F}^i)\right). \quad (8)$$

We found that using 3 parallel convolution layers with sizes  $3 \times 3$ ,  $5 \times 5$ , and  $7 \times 7$  provide optimal performance. Increasing the convolution kernel size reduces performance



**TABLE 1. Performance comparison on Multimodal Material Segmentation (MCubeS) dataset [9]**

Method	Modalities	% mIoU
DRCConv [36]	RGB-A-D-N	34.63
DDF [37]	RGB-A-D-N	36.16
TransFuser [20]	RGB-A-D-N	37.66
DeepLabv3+ [32]	RGB-A-D-N	38.13
MMTM [38]	RGB-A-D-N	39.71
FuseNet [14]	RGB-A-D-N	40.58
MCubeSNet [9]	RGB-A-D-N	42.46
CBAM [39]	RGB-A-D-N	51.32
CMNeXt [24]	RGB-A-D-N	51.54
<b>MMSFormer (Ours)</b>	RGB-A-D-N	<b>53.11</b>

Here A, D, and N represent angle of linear polarization (AoLP), degree of linear polarization (DoLP), and near-infrared (NIR) respectively.

which we show in Table 7. As larger kernels reduce performance, we did not add more than 3 parallel convolutions in our model. We apply Squeeze-and-Excitation block [26] as channel attention in the residual connection. The final fused feature can be represented as

$$F^i = \text{ChannelAttention}(\hat{F}^i) + F^i. \quad (9)$$

Channel attention re-calibrates interdependence between channels and allows the model to select the most relevant features or channels while suppressing less important ones [26]. This leads to more effective feature representations and thus better performance on the underlying task.

#### D. SHARED MLP DECODER

The fused features generated from all the 4 fusion blocks are sent to the shared MLP decoder. We use the decoder design proposed in [25]. The decoder shown in Fig. 1(a) can be represented as the following equations:

$$\begin{aligned} F^i &= \text{Linear}(F^i), \quad \forall i \in \{1, 2, 3, 4\} \\ F^i &= \text{Upsample}(F^i), \quad \forall i \in \{1, 2, 3, 4\} \\ F &= \text{Linear}(F^1 || \dots || F^4), \\ P &= \text{Linear}(F). \end{aligned} \quad (10)$$

The first linear layers take the fused features of different shapes and generate features having the same channel dimension. Then the features are up-sampled to  $\frac{1}{4}$ th of the original input shape, concatenated along the channel dimension and passed through another linear layer to generate the final fused feature  $F$ . Finally  $F$  is passed through the last linear layer to generate the predicted segmentation map  $P$ .

## IV. EXPERIMENTS AND RESULTS

We evaluated our model and proposed fusion block on multiple datasets and with different modality combinations for multimodal semantic and material segmentation tasks. We

**TABLE 2. Performance comparison on FBM [27] dataset**

Methods	Modalities	% mIoU
CBAM [39]	RGB-Infrared	50.1
GMNet [40]	RGB-Infrared	49.2
LASNet [41]	RGB-Infrared	42.5
EGFNet [42]	RGB-Infrared	47.3
FEANet [43]	RGB-Infrared	46.8
DIDFuse [44]	RGB-Infrared	50.6
ReCoNet [45]	RGB-Infrared	50.9
U2Fusion [46]	RGB-Infrared	47.9
TarDAL [47]	RGB-Infrared	48.1
SegMiF [27]	RGB-Infrared	54.8
<b>MMSFormer (Ours)</b>	RGB-Infrared	<b>61.7</b>

We show performance for different methods from already published works.

also compared our methods with existing baseline methods both qualitatively and quantitatively. We report results from already published works whenever possible. \* indicates that we have used the code and pretrained models from the papers to generate the results.

#### A. DATASETS

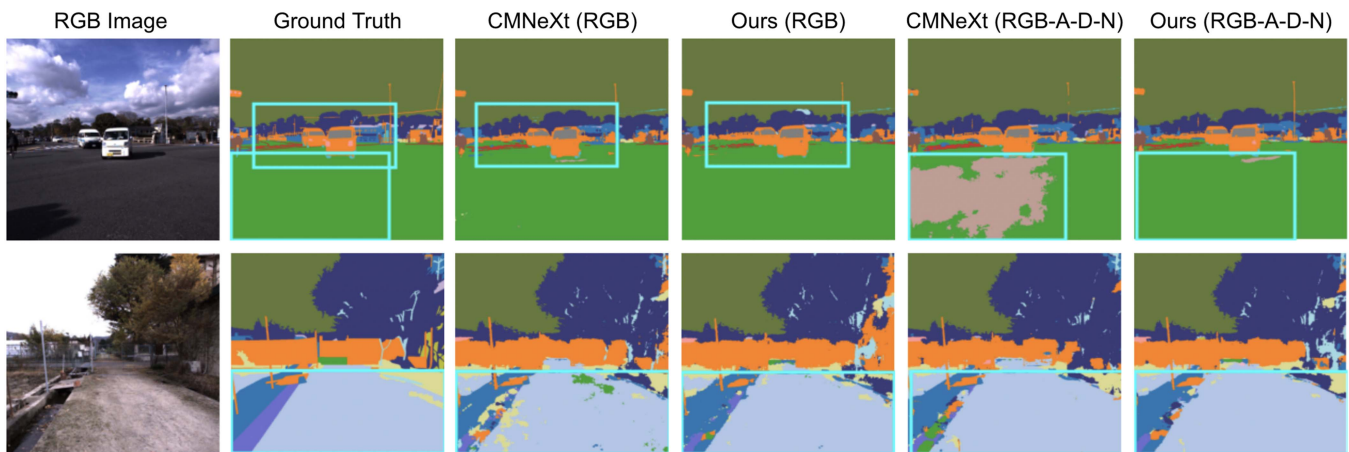
*Multimodal material segmentation (MCubeS) dataset* [9] contains 500 sets of images from 42 street scenes having four modalities: RGB, angle of linear polarization (AoLP), degree of linear polarization (DoLP), and near-infrared (NIR). It provides annotated ground truth labels for both material and semantic segmentation and divided into training set with 302 image sets, validation set with 96 image sets, and test set with 102 image sets. This dataset has 20 class labels corresponding to different materials.

*FMB dataset* [27] is a new and challenging dataset with 1500 pairs of calibrated RGB-Infrared image pairs. The training and test set contains 1220 and 280 image pairs respectively. The dataset covers a wide range of scenes under different lighting and weather conditions (Tyndall effect, rain, fog, and strong light). It also provides per pixel ground truth annotation for 14 different classes.

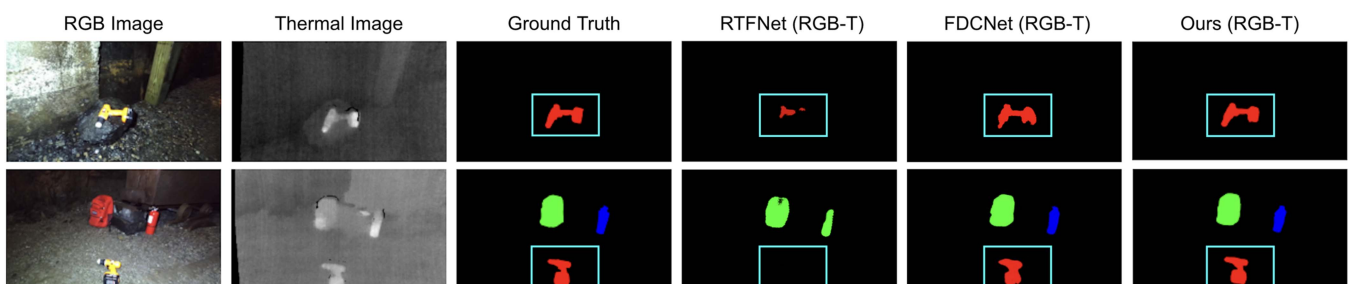
*PST900 dataset* [28] contains 894 pairs of synchronized RGB-Thermal image pairs. The dataset is divided into training and test sets with per pixel ground truth annotation for five different classes.

#### B. IMPLEMENTATION DETAILS

To ensure a fair comparison with prior models, we followed the same data preprocessing and augmentation strategies employed in previous studies [9], [24], [27]. We used the Mix-Transformer (MiT) [25] encoder pretrained on the ImageNet [61] dataset as the backbone for our model to extract features from different modalities. Each modality has a separate encoder. We used a shared MLP decoder introduced in SegFormer [25] and used random initialization for it. We



(a) Visualization of predictions on MCubeS dataset



(b) Visualization of predictions on PST900 dataset

**FIGURE 2.** Visualization of predictions on MCubeS and PST900 datasets. Fig. 2(a) shows RGB and all modalities (RGB-A-D-N) prediction from CMNeXt [24] and our model on MCubeS dataset. For brevity, we only show the RGB image and ground truth material segmentation maps along with the predictions. Fig. 2(b) shows predictions from RTFNet [18], FDCNet [48] and our model for RGB-thermal input modalities on PST900 dataset. Our model shows better predictions on both of the datasets. (a) Visualization of predictions on MCubeS dataset (b) Visualization of predictions on PST900 dataset.

trained and evaluated all our models using two NVIDIA RTX 2080Ti GPUs and used PyTorch for model development.

We utilized a polynomial learning rate scheduler with a power of 0.9 to dynamically adjust the learning rate during training. The first 10 epochs were designated as warm-up epochs with a learning rate of 0.1 times the original rate. For loss computation, we used the cross-entropy loss function. Optimization was performed using the AdamW [62] optimizer with an epsilon value of  $10^{-8}$  and weight decay set to 0.01. For CBAM [39], we use the same encoder, decoder and hyperparameters used in our experiments and replace our fusion block with CBAM<sup>1</sup> module. To be specific, after extracting the feature maps from each input modality using the modality specific encoders, we add them (sum them up) and pass the combined feature map to the CBAM module.

### C. PERFORMANCE COMPARISON WITH EXISTING METHODS

We conducted a rigorous performance evaluation of our model compared to established baseline models for three datasets. The comprehensive results are summarized in Tables 1–6. We

report results for CBAM from our experiments. Other results are taken from published literature.

*Results on MCubeS Dataset:* Table 1 shows the overall performance comparison between our model and existing baseline models for MCubeS dataset. Our model achieves a mean intersection-over-union (mIoU) of 53.11%, surpassing the current state-of-the-art model. It shows 1.57% improvement over CMNeXt [24], 1.79% improvement over CBAM [39] and 10.65% improvement over MCubeSNet [9] models. To further analyze the performance of our model, we conducted a per-class IoU analysis and presented in Table 3. Our model performs better in detecting most of the material classes compared to the current state-of-the-art models. Notably, our model demonstrates a substantial improvement in the detection of plastic (+3.7%), fabric (+3.1%), asphalt (+2.3%), and cobblestone (2.3%) classes while maintaining competitive or better performance in other classes. This led to the overall better performance and sets new state-of-the-art for this dataset.

*Results on FMB Dataset:* Performance Comparison for FMB dataset is shown on Table 2. Our model shows a significant improvement of 6.9% mIoU compared to the current state-of-the-art model. Per-class IoU analysis for this dataset is shown on Table 4. For a fair comparison, we only compare

<sup>1</sup>[Online]. Available: <https://github.com/luuuyi/CBAM.PyTorch>

**TABLE 3. Per-Class % IoU Comparison on Multimodal Material Segmentation (MCubeS) [9] Dataset**

Methods	Asphalt	Concrete	Metal	Road marking	Fabric	Glass	Plaster	Plastic	Rubber	Sand	Gravel	Ceramic	Cobblestone	Brick	Grass	Wood	Leaf	Water	Human	Sky	Mean
MCubeSNet [9]	85.7	42.6	47.0	59.2	12.5	44.3	3.0	10.6	12.7	66.8	67.1	27.8	65.8	36.8	54.8	39.4	73.0	13.3	0.0	94.8	42.9
CBAM [39]	85.7	47.7	55.4	70.4	27.6	54.7	<b>0.9</b>	30.9	26.5	61.6	63.0	28.0	71.1	41.8	58.6	47.4	76.7	<b>56.3</b>	25.9	96.5	51.3
CMNeXt [24] *	84.3	44.9	53.9	<b>74.5</b>	32.3	54.0	0.8	28.3	<b>29.7</b>	<b>67.7</b>	66.5	27.7	68.5	42.9	58.7	<b>49.7</b>	75.4	55.7	18.9	96.5	51.5
<b>MMSFormer (Ours)</b>	<b>88.0</b>	<b>48.3</b>	<b>56.2</b>	72.2	<b>35.4</b>	<b>54.9</b>	0.5	<b>34.6</b>	29.4	67.2	<b>69.0</b>	<b>29.9</b>	<b>73.4</b>	<b>44.7</b>	<b>59.5</b>	47.8	<b>77.1</b>	50.5	<b>26.9</b>	<b>96.6</b>	<b>53.1</b>

Our proposed MMSFormer model shows better performance in detecting most of the classes compared to the current state-of-the-art models. \* indicates that the models. \* indicates that the code and pretrained model from the authors were used to generate the results.

**TABLE 4. Per-Class % IoU Comparison on FMB [27] Dataset for Both RGB Only and RGB-Infrared Modalities**

Methods	Modalities	Car	Person	Truck	T-Lamp	T-Sign	Building	Vegetation	Pole	% mIoU
SegMiF [27]	RGB	78.3	46.6	<b>43.4</b>	23.7	64.0	77.8	82.1	41.8	50.5
<b>MMSFormer (Ours)</b>	RGB	<b>80.3</b>	<b>56.7</b>	42.1	<b>31.6</b>	<b>77.8</b>	<b>77.9</b>	<b>85.4</b>	<b>48.1</b>	<b>57.2</b>
CBAM [39]	RGB-Infrared	71.9	49.3	20.9	25.8	67.1	75.8	80.9	19.7	50.1
GMNet [40]	RGB-Infrared	79.3	60.1	22.2	21.6	69.0	79.1	83.8	39.8	49.2
LASNet [41]	RGB-Infrared	72.6	48.6	14.8	2.9	59.0	75.4	81.6	36.7	42.5
EGFNet [42]	RGB-Infrared	77.4	63.0	17.1	25.2	66.6	77.2	83.5	41.5	47.3
FEANet [43]	RGB-Infrared	73.9	60.7	32.3	13.5	55.6	79.4	81.2	36.8	46.8
DIDFuse [44]	RGB-Infrared	77.7	64.4	28.8	29.2	64.4	78.4	82.4	41.8	50.6
ReCoNet [45]	RGB-Infrared	75.9	65.8	14.9	34.7	66.6	79.2	81.3	44.9	50.9
U2Fusion [46]	RGB-Infrared	76.6	61.9	14.4	28.3	68.9	78.8	82.2	42.2	47.9
TarDAL [47]	RGB-Infrared	74.2	56.0	18.8	29.6	66.5	79.1	81.7	41.9	48.1
SegMiF [27]	RGB-Infrared	78.3	65.4	<b>47.3</b>	43.1	74.8	82.0	85.0	49.8	54.8
<b>MMSFormer (Ours)</b>	RGB-Infrared	<b>82.6</b>	<b>69.8</b>	44.6	<b>45.2</b>	<b>79.7</b>	<b>83.0</b>	<b>87.3</b>	<b>51.4</b>	<b>61.7</b>

We show the comparison for 8 classes (out of 14) that are published. T-Lamp and T-Sign stand for Traffic Lamp and Traffic Sign respectively. Our model outperforms all the methods for all the classes except for the truck class.

the performance on 8 classes (out of 14) that are published in literature. T-Lamp and T-Sign represent Traffic Lamp and Traffic Sign, respectively. Our model shows an overall performance improvement of 6.7% mIoU for RGB only predictions compared to the most recent SegMiF [27] model. Alongside this, our model also shows superior performance in detecting all of the classes except for the truck class for both RGB only and RGB-Infrared semantic segmentation tasks. Performance on RGB-Infrared input modalities is much better than RGB only performance for all the classes, which demonstrates the ability of the fusion block to effectively fuse information from the input modalities.

*Results on PST900 Dataset:* We also tested our model on PST900 [28] dataset and summarized the result in Table 5. Experiments show that our model outperforms existing baseline models for RGB-Thermal semantic segmentation on this dataset. It outperforms the most recent CACFNet [60] model by 0.89% mIoU. Our model also shows better performance in detecting 3 out of the 5 classes available in the dataset and competitive performance in other two classes.

#### D. PERFORMANCE COMPARISON FOR INCREMENTAL MODALITY INTEGRATION

A critical aspect of this work involves evaluating the effectiveness of our proposed fusion block in combining valuable information from diverse modalities. To analyze this effect, we trained our model with various combinations of modalities on the MCubeS dataset. The results are presented in Table 6. Our model exclusively trained on RGB data provided an mIoU score of 50.44%, which is 2.28% greater than the current state-of-the-art model. We observe progressive improvement in performance as we incorporated additional modalities: AoLP, DoLP, and NIR. The integration led to incremental performance gains, with the mIoU increasing from 50.44% to 51.30%, then to 52.03%, and ultimately reaching to 53.11%. These findings serve as a compelling evidence that our fusion approach effectively leverages and fuses valuable information from different combination of modalities, resulting in a notable enhancement in segmentation performance.

Furthermore, our model consistently outperforms the current state-of-the-art benchmark across all modality

**TABLE 5.** Performance Comparison on PST900 [28] Dataset

Methods	Modalities	Background	Fire-Extinguisher	Backpack	Hand-Drill	Survivor	% mIoU
ACNet [15]	RGB-Thermal	99.25	59.95	83.19	51.46	65.19	71.81
CCNet [49]	RGB-Thermal	99.05	51.84	66.42	32.27	57.50	61.42
Efficient FCN [50]	RGB-Thermal	98.63	39.96	58.15	30.12	28.00	50.98
RTFNet [18]	RGB-Thermal	99.02	51.93	74.17	7.07	70.11	60.46
PSTNet [28]	RGB-Thermal	98.85	70.12	69.20	53.60	50.03	68.36
EGFNet [51]	RGB-Thermal	99.26	71.29	83.05	64.67	74.30	78.51
MTANet [52]	RGB-Thermal	99.33	64.95	87.50	62.05	79.14	78.60
MFFNet [53]	RGB-Thermal	99.40	66.38	81.02	72.50	75.60	78.98
GMNet [40]	RGB-Thermal	99.44	73.79	83.82	85.17	78.36	84.12
CGFNet [54]	RGB-Thermal	99.30	71.71	82.00	59.72	77.42	78.03
GCNet [55]	RGB-Thermal	99.35	77.68	79.37	82.92	73.58	82.58
GEBNet [56]	RGB-Thermal	99.39	73.07	85.93	67.14	80.21	81.15
GCGLNet [57]	RGB-Thermal	99.39	77.57	81.01	81.90	76.31	83.24
DHFNet [58]	RGB-Thermal	99.44	78.15	87.38	71.18	74.81	82.19
MDRNet+ [59]	RGB-Thermal	99.07	63.04	76.27	63.47	71.26	74.62
FDCNet [48]	RGB-Thermal	99.15	71.52	72.17	70.36	72.36	77.11
CBAM [39]	RGB-Thermal	99.43	73.81	82.75	80.00	69.60	81.12
EGFNet [42]	RGB-Thermal	99.55	79.97	90.62	76.08	<b>80.88</b>	85.42
CACFNet [60]	RGB-Thermal	99.57	<b>82.08</b>	89.49	80.90	80.76	86.56
<b>MMSFormer (Ours)</b>	RGB-Thermal	<b>99.60</b>	81.45	<b>89.86</b>	<b>89.65</b>	76.68	<b>87.45</b>

We show per-class % IoU as well as % mIoU for all the classes.

**TABLE 6.** Performance Comparison (% mIoU) on Multimodal Material Segmentation (MCubeS) [9] Dataset for Different Modality Combinations

Modalities	MCubeSNet [9]	CMNeXt [24]	MMSFormer (Ours)
RGB	33.70	48.16	<b>50.44</b>
RGB-A	39.10	48.42	<b>51.30</b>
RGB-A-D	42.00	49.48	<b>52.03</b>
RGB-A-D-N	42.86	51.54	<b>53.11</b>

Here A, D, and N represent angle of linear polarization (AoLP), degree of linear polarization (DoLP), and near-infrared (NIR) respectively.

combinations. This consistent superiority underscores the robustness and versatility of our fusion block, demonstrating its ability to adapt and excel regardless of the specific modality combination provided.

### E. QUALITATIVE ANALYSIS OF THE PREDICTIONS

Apart from quantitative analysis, we also perform qualitative analysis of the predicted segmentation maps. We show material segmentation results predicted by CMNeXt [24] model and our proposed MMSFormer model in Fig. 2(a). For brevity, we only show RGB images and ground truth material segmentation maps in the illustrations. We show RGB only predictions and all modalities (RGB-A-D-N) predictions for both of the models. As highlighted in the rectangular bounding boxes, our proposed MMSFormer model identifies asphalt, sand and water with greater accuracy than CMNeXt [24] model for both RGB only and all modalities (RGB-A-D-N) predictions.

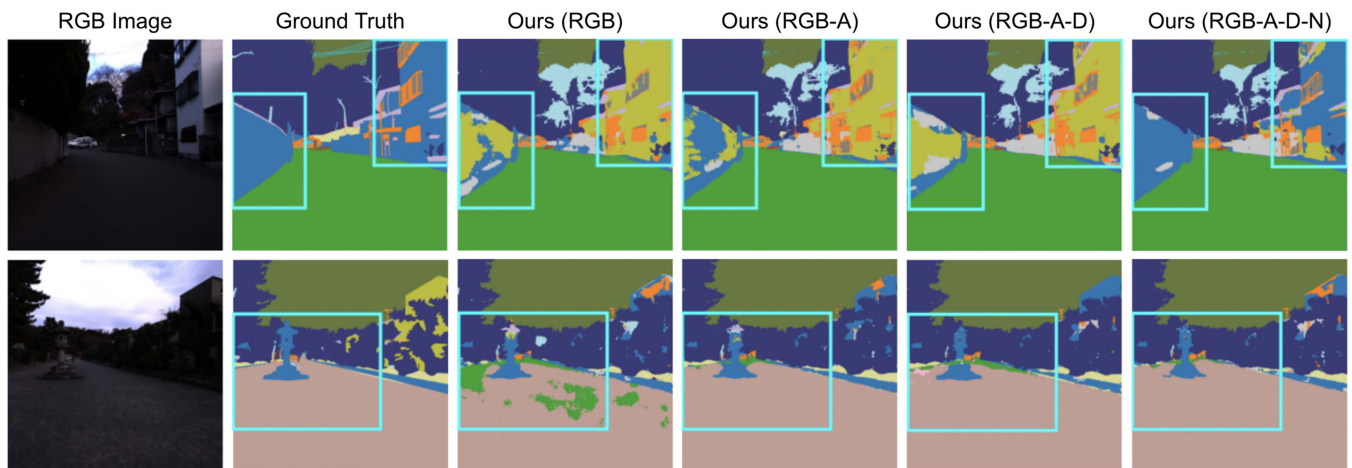
We also compare our prediction on PST900 [28] dataset with RTFNet [18] and FDCNet [48] on Fig. 2(b). We show the input RGB image, thermal images, ground truth segmentation maps and prediction from the models. As highlighted by the rectangular bounding boxes, our model shows better accuracy in detecting objects with more precise contours compared to the other two methods.

### F. ABLATION STUDY ON THE FUSION BLOCK

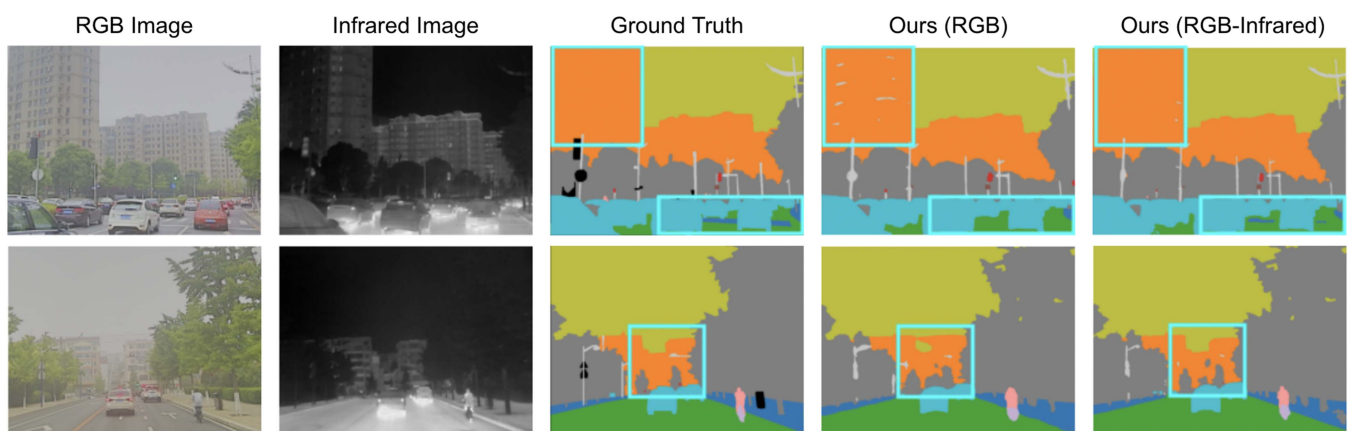
We conducted a number of ablation studies aimed at investigating the contributions of individual components within the fusion block to the overall model performance. The findings, as detailed in Table 7, shed light on the critical importance of these components. We used both RGB and infrared modalities of the FMB dataset during training and testing in these experiments. First, we observed that the absence of channel attention in the residual connection had a negative impact, resulting in a reduction in performance by 3.36%. This indicates that feature calibration along channel dimension plays an important role in capturing and leveraging crucial information effectively. Additionally, while comparing larger convolution kernel sizes ( $3 \times 3$ ,  $7 \times 7$ , and  $11 \times 11$ ) to the originally employed ( $3 \times 3$ ,  $5 \times 5$ , and  $7 \times 7$ ), we noted a decrease in performance by 5.36%. This result underscores the significance of the carefully chosen convolution kernel sizes within the fusion block.

Furthermore, completely removing the parallel convolutions from the block led to a performance decline of 4.51%, emphasizing their substantial contribution in capturing multi-scale features and overall model performance. Finally, if we





(a) Visualization of predictions on MCubeS dataset for different modality combinations



(b) Visualization of predictions on FMB dataset for different modality combinations

**FIGURE 3.** Visualization of predicted segmentation maps for different modality combinations on MCubeS [9] and FMB [27] datasets. Both figures show that prediction accuracy increases as we incrementally add new modalities. They also illustrate the fusion block’s ability to effectively combine information from different modality combinations. (a) Visualization of predictions on MCubeS dataset for different modality combinations (b) Visualization of predictions on FMB dataset for different modality combinations.

only use the linear fusion layer to fuse the features and remove the parallel convolutions and channel attention from the fusion block, performance drop significantly by 9.25%. These studies demonstrate that multi-scale feature capturing via parallel convolutions and channel-wise feature calibration using channel attention is extremely important in learning better feature representation and thus crucial to overall model performance. These comprehensive ablation studies collectively underscore the significance of every component within the fusion block, revealing that each module plays a distinct and vital role in achieving the overall performance of our model.

### G. ABLATION STUDY ON DIFFERENT MODALITY COMBINATIONS

To analyze the contributions of different modalities in the identification of distinct materials, we conducted a series of ablation studies, focusing on per-class IoU for different modality combinations. The insights are summarized in

Table 8. As we progressively integrate new modalities, performance gradually increases for specific classes, which include grass, leaf, asphalt, cobblestone and plastic classes. Particularly noteworthy is the assistance provided by NIR data in classifying asphalt, concrete, plastic, cobblestone, and human categories, leading to significant performance gains in these classes as NIR was added as an additional modality.

Conversely, certain classes such as water and brick exhibited a gradual performance decline as we introduced additional modalities. This suggests that RGB data alone suffices for accurately identifying these classes, and the inclusion of more modalities potentially introduces noise or redundancy that negatively impacts performance. Moreover, AoLP appears to be helpful in enhancing the recognition of materials like road markings, glass and wood. Similarly, DoLP improved performance for classes like plaster, rubber, sand, gravel and ceramic. These findings underscore the relationship between different imaging modalities and the unique characteristics of different types of materials, demonstrating

**TABLE 7. Ablation Study of the Fusion Block on FMB [27] Dataset**

Structure	Parameter Count (M)	% mIoU (Change)
MMSFormer	61.26	61.68
- without channel attention	61.21	58.32 (-3.36)
- without parallel convolutions	61.17	57.17 (-4.51)
- with 3x3, 7x7 and 11x11 convolutions	61.36	56.32 (-5.36)
- only linear fusion	59.57	52.43 (-9.25)

Both RGB and infrared input modalities were used during training and testing. The table shows the contribution of different modules in fusion block in overall model performance.

**TABLE 8. Per Class % IoU Comparison on Multimodal Material Segmentation (MCubeS) [9] Dataset for Different Modality Combinations**

Modalities	Asphalt	Concrete	Metal	Road marking	Fabric	Glass	Plaster	Plastic	Rubber	Sand	Gravel	Ceramic	Cobblestone	Brick	Grass	Wood	Leaf	Water	Human	Sky	Mean
RGB	83.2	44.2	52.1	70.4	31.0	51.6	1.3	26.2	21.8	65.0	61.8	31.3	72.5	<b>45.0</b>	55.4	46.0	74.7	<b>56.0</b>	22.7	96.4	50.4
RGB-A	86.5	46.5	55.9	<b>73.0</b>	35.3	<b>56.0</b>	0.8	27.3	27.8	66.2	67.0	28.6	69.6	43.0	57.6	<b>49.6</b>	76.4	53.8	8.4	96.6	51.3
RGB-A-D	86.0	44.0	55.5	68.1	31.9	54.8	<b>2.3</b>	30.0	<b>29.7</b>	<b>69.4</b>	<b>73.7</b>	<b>32.2</b>	69.4	41.4	59.2	48.3	76.6	50.6	20.9	<b>96.7</b>	52.0
RGB-A-D-N	<b>88.0</b>	<b>48.3</b>	<b>56.2</b>	72.2	<b>35.4</b>	54.9	0.5	<b>34.6</b>	29.4	67.2	69.0	29.9	<b>73.4</b>	44.7	<b>59.5</b>	47.8	<b>77.1</b>	50.5	<b>26.9</b>	96.6	<b>53.1</b>

As we add modalities incrementally, overall performance increases gradually. This table also shows that specific modality combinations assist in identifying specific types of materials better.

that specific modalities excel in detecting particular classes based on their distinctive traits.

In Fig. 3(a), we presents some examples to show how adding different modalities help improve performance of segmentation. We show predictions for RGB, RGB-A, RGB-A-D and RGB-A-D-N inputs from our proposed MMSFormer model. As we add new modalities, the predictions become more accurate as shown in the bounding boxes. The illustrations show that the identification of concrete and gravel becomes more accurate with additional modalities. Fig. 3(b) shows predictions for RGB and RGB-Infrared from FMB dataset. As highlighted by the bounding boxes, adding new modality helps improve performance in detecting building, road and sidewalks. This also illustrates the capability of the fusion block to effectively fuse information from different modality combinations.

#### H. COMPUTATIONAL COMPLEXITY OF THE FUSION BLOCK

In addition to better performance, our fusion block is also computationally efficient compared to most of the fusion blocks proposed for these datasets. We show a comparison in terms of the number of parameters in the fusion block and GFLOPs for some recent models on Table 9 for MCubeS dataset having 4 input modalities with an input shape of  $(3 \times 512 \times 512)$  for each modality. As observed from the table, our proposed fusion strategy demonstrates significantly lower complexity in terms of both the number of parameters and GFLOPs compared to existing methods. HRFuser [23] has a lower parameter count than ours but it requires more than  $7 \times$  GFLOPs. Other methods require significantly more

**TABLE 9. Comparison of Number of Parameters in the Fusion Block and GFLOPs for MCubeS Dataset Having 4 Input Modalities With an Input Shape of  $(3 \times 512 \times 512)$  for Each Modality. Our Fusion Block Shows Significantly Lower Complexity Compared to Existing Methods**

Methods	Fusion Block Parameters (M)	GFLOPs
CMNeXt [24]	16.63	6.47
MCubeSNet [9]	7.41	12.10
HRFuser [23]	<b>1.72</b>	17.50
CMX [22]	16.59	6.41
DDF (Resnet-101) [37]	28.10	4.10
<b>MMSFormer (Ours)</b>	3.23	<b>2.47</b>

Our fusion block shows significantly lower complexity compared to existing methods.

parameters ( $2.3 \times -8.7 \times$ ) and GFLOPs ( $1.6 \times -7 \times$ ) compared to our fusion strategy. Our comparison only includes models for which these results are available in the published literature.

#### V. CONCLUSION

In this article, we introduce a novel fusion module designed to combine useful information from various modality combinations. We also propose a new model called MMSFormer that integrates the proposed fusion block to accomplish multimodal material and semantic segmentation tasks. Experimental results illustrate the model's capability to efficiently fuse information from different combination of modalities, leading to new state-of-the-art performance on three different datasets. Experiments also show that the fusion block can extract useful

information from different modality combinations that helps the model to consistently outperform current state-of-the-art models. Starting from only one input modality, performance increases gradually as we add new modalities. Several ablation studies further highlight how different components of the fusion block contribute to the overall model performance. Ablation studies also reveal that different modalities assist in identifying different types of material classes. However, one limitation of the proposed model is the use of modality specific encoders and the number of encoders grows with number of modalities. Future work will include exploring the possibility and effectiveness of using a shared encoder for all the modalities, investigating and extending the model's capability with other modalities and multimodal tasks.

## REFERENCES

- [1] H.-D. Cheng, X. H. Jiang, Y. Sun, and J. Wang, "Color image segmentation: Advances and prospects," *Pattern Recognit.*, vol. 34, no. 12, pp. 2259–2281, 2001.
- [2] S. Minaee, Y. Boykov, F. Porikli, A. Plaza, N. Kehtarnavaz, and D. Terzopoulos, "Image segmentation using deep learning: A survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 7, pp. 3523–3542, Jul. 2022.
- [3] Y. Guo, Y. Liu, T. Georgiou, and M. S. Lew, "A review of semantic segmentation using deep neural networks," *Int. J. Multimedia Inf. Retrieval*, vol. 7, pp. 87–93, 2018.
- [4] P. Wang et al., "Understanding convolution for semantic segmentation," in *Proc. IEEE Winter Conf. Appl. Comput. Vis.*, 2018, pp. 1451–1460.
- [5] W. Gu, S. Bai, and L. Kong, "A review on 2 d instance segmentation based on deep neural networks," *Image Vis. Comput.*, vol. 120, 2022, Art. no. 104401.
- [6] A. M. Hafiz and G. M. Bhat, "A survey on instance segmentation: State of the art," *Int. J. Multimedia Inf. Retrieval*, vol. 9, no. 3, pp. 171–189, 2020.
- [7] O. Elharrouss, S. A. Al-Maadeed, N. Subramanian, N. Ottakath, N. Almaadeed, and Y. Himeur, "Panoptic segmentation: A review," 2021, *arXiv:2111.10250*.
- [8] A. Kirillov, K. He, R. Girshick, C. Rother, and P. Dollár, "Panoptic segmentation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 9404–9413.
- [9] Y. Liang, R. Wakaki, S. Nobuhara, and K. Nishino, "Multimodal material segmentation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 19800–19808.
- [10] P. Upchurch and R. Niu, "A dense material segmentation dataset for indoor and outdoor scene parsing," in *Proc. Eur. Conf. Comput. Vis.*, 2022, pp. 450–466.
- [11] Z. Guo, X. Li, H. Huang, N. Guo, and Q. Li, "Deep learning-based image segmentation on multimodal medical imaging," *IEEE Trans. Radiat. Plasma Med. Sci.*, vol. 3, no. 2, pp. 162–169, Mar. 2019.
- [12] Y. Zhang, D. Sidibé, O. Morel, and F. Mériaudeau, "Deep multimodal fusion for semantic image segmentation: A survey," *Image Vis. Comput.*, vol. 105, 2021, Art. no. 104042.
- [13] X. Chen et al., "Bi-directional cross-modality feature propagation with separation-and-aggregation gate for RGB-D semantic segmentation," in *Proc. Eur. Conf. Comput. Vis.*, 2020, pp. 561–577.
- [14] C. Hazirbas, L. Ma, C. Domokos, and D. Cremers, "Fusenet: Incorporating depth into semantic segmentation via fusion-based CNN architecture," in *Proc. Asian Conf. Comput. Vis.*, 2016, pp. 213–228.
- [15] X. Hu, K. Yang, L. Fei, and K. Wang, "ACNet: Attention based network to exploit complementary features for RGBD semantic segmentation," in *Proc. IEEE Int. Conf. Image Process.*, 2019, pp. 1440–1444.
- [16] P. Li, J. Chen, B. Lin, and X. Xu, "Residual spatial fusion network for RGB-thermal semantic segmentation," 2023, *arXiv:2306.10364*.
- [17] M. Liang, J. Hu, C. Bao, H. Feng, F. Deng, and T. L. Lam, "Explicit attention-enhanced fusion for RGB-thermal perception tasks," *IEEE Robot. Autom. Lett.*, vol. 8, no. 7, pp. 4060–4067, Jul. 2023.
- [18] Y. Sun, W. Zuo, and M. Liu, "RTFNet: RGB-Thermal fusion network for semantic segmentation of urban scenes," *IEEE Robot. Autom. Lett.*, vol. 4, no. 3, pp. 2576–2583, Jul. 2019.
- [19] J. Li, H. Dai, H. Han, and Y. Ding, "MSEG3D: Multi-modal 3 D semantic segmentation for autonomous driving," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2023, pp. 21694–21704.
- [20] A. Prakash, K. Chitta, and A. Geiger, "Multi-modal fusion transformer for end-to-end autonomous driving," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 7077–7087.
- [21] L. Zhao, H. Zhou, X. Zhu, X. Song, H. Li, and W. Tao, "LIF-SEG: LiDAR and camera image fusion for 3 D LiDAR semantic segmentation," *IEEE Trans. Multimedia*, vol. 26, pp. 1158–1168, 2024.
- [22] J. Zhang, H. Liu, K. Yang, X. Hu, R. Liu, and R. Stiefelhagen, "CMX: Cross-modal fusion for RGB-X semantic segmentation with transformers," *IEEE Trans. Intell. Transp. Syst.*, vol. 24, no. 12, pp. 14679–14694, Dec. 2023.
- [23] T. Broedermann, C. Sakaridis, D. Dai, and L. Van Gool, "HRFuser: A multi-resolution sensor fusion architecture for 2D object detection," in *Proc. IEEE Int. Conf. Intell. Transp. Syst.*, 2023, pp. 4159–4166.
- [24] J. Zhang et al., "Delivering arbitrary-modal semantic segmentation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2023, pp. 1136–1147.
- [25] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, "SegFormer: Simple and efficient design for semantic segmentation with transformers," in *Proc. Neural Inf. Process. Syst.*, 2021, pp. 12077–12090.
- [26] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 7132–7141.
- [27] J. Liu et al., "Multi-interactive feature learning and a full-time multi-modality benchmark for image fusion and segmentation," in *Proc. Int. Conf. Comput. Vis.*, 2023, pp. 8115–8124.
- [28] S. S. Shivakumar, N. Rodrigues, A. Zhou, I. D. Miller, V. Kumar, and C. J. Taylor, "PST900: RGB-thermal calibration, dataset, and segmentation network," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2020, pp. 9441–9447.
- [29] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 3431–3440.
- [30] H. Wu, J. Zhang, K. Huang, K. Liang, and Y. Yu, "FastFCN: Rethinking dilated convolution in the backbone for semantic segmentation," 2019, *arXiv:1903.11816*.
- [31] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *Proc. Med. Image Comput. Comput.-Assist. Interv.: 18th Int. Conf.*, 2015, pp. 234–241.
- [32] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 801–818.
- [33] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 2881–2890.
- [34] W. Wang et al., "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2021, pp. 568–578.
- [35] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, and R. Girdhar, "Masked-attention mask transformer for universal image segmentation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 1290–1299.
- [36] J. Chen, X. Wang, Z. Guo, X. Zhang, and J. Sun, "Dynamic region-aware convolution," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 8064–8073.
- [37] J. Zhou, V. Jampani, Z. Pi, Q. Liu, and M.-H. Yang, "Decoupled dynamic filter networks," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 6647–6656.
- [38] H. R. V. Joze, A. Shaban, M. L. Iuzzolino, and K. Koishida, "MMTM: Multimodal transfer module for CNN fusion," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 13289–13299.
- [39] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "CBAM: Convolutional block attention module," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 3–19.
- [40] W. Zhou, J. Liu, J. Lei, L. Yu, and J.-N. Hwang, "GMNet: Graded-feature multilabel-learning network for RGB-thermal urban scene semantic segmentation," *IEEE Trans. Image Process.*, vol. 30, pp. 7790–7802, 2021.
- [41] G. Li, Y. Wang, Z. Liu, X. Zhang, and D. Zeng, "RGB-T semantic segmentation with location, activation, and sharpening," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 33, no. 3, pp. 1223–1235, Mar. 2023.



- [42] S. Dong, W. Zhou, C. Xu, and W. Yan, "EGFNet: Edge-aware guidance fusion network for RGB-thermal urban scene parsing," *IEEE Trans. Intell. Transp. Syst.*, vol. 25, no. 1, pp. 657–669, Jan. 2024.
- [43] F. Deng et al., "FEAnet: Feature-enhanced attention network for RGB-thermal real-time semantic segmentation," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2021, pp. 4467–4473.
- [44] Z. Zhao, S. Xu, C. Zhang, J. Liu, J. Zhang, and P. Li, "DIDFuse: Deep image decomposition for infrared and visible image fusion," in *Proc. Int. Joint Conf. Artif. Intell.*, 2020, pp. 970–976.
- [45] Z. Huang, J. Liu, X. Fan, R. Liu, W. Zhong, and Z. Luo, "RECONet: Recurrent correction network for fast and efficient multi-modality image fusion," in *Proc. Eur. Conf. Comput. Vis.*, 2022, pp. 539–555.
- [46] H. Xu, J. Ma, J. Jiang, X. Guo, and H. Ling, "U2Fusion: A unified unsupervised image fusion network," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 1, pp. 502–518, Jan. 2022.
- [47] J. Liu et al., "Target-aware dual adversarial learning and a multi-scenario multi-modality benchmark to fuse infrared and visible for object detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 5802–5811.
- [48] S. Zhao and Q. Zhang, "A feature divide-and-conquer network for RGB-T semantic segmentation," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 33, no. 6, pp. 2892–2905, Jun. 2023.
- [49] Z. Huang, X. Wang, L. Huang, C. Huang, Y. Wei, and W. Liu, "CCNet: Criss-cross attention for semantic segmentation," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2019, pp. 603–612.
- [50] J. Liu, J. He, J. Zhang, J. S. Ren, and H. Li, "EfficientFCN: Holistically-guided decoding for semantic segmentation," in *Proc. Eur. Conf. Comput. Vis.*, 2020, pp. 1–17.
- [51] W. Zhou, S. Dong, C. Xu, and Y. Qian, "Edge-aware guidance fusion network for RGB-thermal scene parsing," in *Proc. AAAI Conf. Artif. Intell.*, 2022, pp. 3571–3579.
- [52] W. Zhou, S. Dong, J. Lei, and L. Yu, "MTANet: Multitask-aware network with hierarchical multimodal fusion for RGB-T urban scene understanding," *IEEE Trans. Intell. Veh.*, vol. 8, no. 1, pp. 48–58, Jan. 2023.
- [53] W. Zhou, X. Lin, J. Lei, L. Yu, and J.-N. Hwang, "MFFENet: Multiscale feature fusion and enhancement network for RGB-thermal urban road scene parsing," *IEEE Trans. Multimedia*, vol. 24, pp. 2526–2538, 2022.
- [54] J. Wang, K. Song, Y. Bao, L. Huang, and Y. Yan, "CGFNet: Cross-guided fusion network for RGB-T salient object detection," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 32, no. 5, pp. 2949–2961, May 2022.
- [55] J. Liu, W. Zhou, Y. Cui, L. Yu, and T. Luo, "GCNet: Grid-like context-aware network for RGB-thermal semantic segmentation," *Neurocomput.*, vol. 506, no. C., pp. 60–67, Sep. 2022.
- [56] S. Dong, W. Zhou, X. Qian, and L. Yu, "GEBNet: Graph-enhancement branch network for RGB-T scene parsing," *IEEE Signal Process. Lett.*, vol. 29, pp. 2273–2277, 2022.
- [57] T. Gong, W. Zhou, X. Qian, J. Lei, and L. Yu, "Global contextually guided lightweight network for RGB-thermal urban scene understanding," *Eng. Appl. Artif. Intell.*, vol. 117, 2023, Art. no. 105510.
- [58] Y. Cai, W. Zhou, L. Zhang, L. Yu, and T. Luo, "DHFNet: Dual-decoding hierarchical fusion network for RGB-thermal semantic segmentation," *Vis. Comput.*, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:256320037>
- [59] S. Zhao, Y. Liu, Q. Jiao, Q. Zhang, and J. Han, "Mitigating modality discrepancies for RGB-T semantic segmentation," *IEEE Trans. Neural Netw. Learn. Syst.*, early access, Jan. 06, 2023, doi: [10.1109/TNNLS.2022.3233089](https://doi.org/10.1109/TNNLS.2022.3233089).
- [60] W. Zhou, S. Dong, M. Fang, and L. Yu, "CACFNet: Cross-modal attention cascaded fusion network for RGB-T urban scene parsing," *IEEE Trans. Intell. Veh.*, vol. 9, no. 1, pp. 1919–1929, Jan. 2024.
- [61] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2009, pp. 248–255.
- [62] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," in *Proc. Int. Conf. Learn. Representations*, 2019.



**MD KAYKOBAD REZA** received the B.Sc. degree from the Department of Computer Science and Engineering, Bangladesh University of Engineering and Technology, Dhaka, Bangladesh, in 2019. He is currently working toward the Ph.D. degree with the University of California, Riverside, Riverside, CA, USA, under the supervision of Prof. M. Salman Asif. His research interests include multimodal machine learning, computer vision, and natural language processing.



**ASHLEY PRATER-BENNETTE** received the Ph.D. degree in mathematics from Syracuse University, Syracuse, NY, USA, in 2013. She is currently a Senior Research Mathematician with the Air Force Research Laboratory, Information Directorate in Rome NY, where she leads the Autonomy, Command & Control, and Decision Support core technical competency research portfolio for the Air Force. Her research interests include mathematical foundations of machine learning, optimization, tensor decompositions, and sparse and low-rank data representations, and applications of the aforementioned to improved speed, scale, and quality of US Air Force operational level decision making.



**M. SALMAN ASIF** (Senior Member, IEEE) received the B.Sc. degree from the University of Engineering and Technology, Lahore, Pakistan, and the M.S. and Ph.D. degrees from the Georgia Institute of Technology, Atlanta, GA, USA. He was a Postdoctoral Researcher with Rice University, Houston, TX, USA, and a Senior Research Engineer with Samsung Research America, Dallas. He is currently an Associate Professor with the University of California Riverside, Riverside, CA, USA. His research interests include computational imaging, signal/image processing, computer vision, and machine learning. He was the recipient of the NSF CAREER Award, Google Faculty Award, Hershel M. Rich Outstanding Invention Award, and UC Regents Faculty Fellowship and Development Awards. He is an Associate Editor for IEEE TRANSACTIONS ON COMPUTATIONAL IMAGING and a Member of IEEE Signal Processing Society's Computational Imaging Technical Committee.