

# Supplementary Material

## 1. NearFarMix Augmentation

Figure 1 elegantly illustrates the operation of the proposed NearFarMix augmentation, delineating the four regions—near, pre-far, overlap, and exclusive—that contribute to the final augmented output. Despite the amalgamation of four regions, the pre-far, overlap, and exclusive regions primarily generate far regions, with overlap regions consequently subtracted. Importantly, the augmentation is designed for batch-wise application to enhance execution speed. The step-by-step implementation is comprehensively depicted in Algorithm 3.

Additionally, Fig. 3 and Fig. 2 provide additional examples between the proposed NearFarMix and DepthMix [2].

---

### Algorithm 1 Batchwise Window Partition

---

$x \leftarrow \text{Input}$  ▷ Features  
 $h_w, w_w \leftarrow \text{window\_size}$  ▷ Size of each window

#### # Partition features into local dense windows

$B, H, W, C = \text{shape}(x)$   
 $x = \text{reshape}(x, \text{shape} = [B, \frac{H}{h_w}, h_w, \frac{W}{w_w}, w_w, C])$   
 $x = \text{transpose}(x, \text{permute\_axis} = [0, 1, 3, 2, 4, 5])$   
 $x = \text{reshape}(x, \text{shape} = [B \times \frac{H \times W}{h_w \times w_w}, h_w, w_w, C])$

---



---

### Algorithm 2 Batchwise Grid Partition

---

$x \leftarrow \text{Input}$  ▷ Features  
 $h_g, w_g \leftarrow \text{grid\_size}$  ▷ Size of each grid

#### # Partition features into global sparse grids

$B, H, W, C = \text{shape}(x)$   
 $x = \text{reshape}(x, \text{shape} = [B, h_g, \frac{H}{h_g}, w_g, \frac{W}{w_g}, C])$   
 $x = \text{transpose}(x, \text{permute\_axis} = [0, 1, 3, 2, 4, 5])$   
 $x = \text{reshape}(x, \text{shape} = [B, h_g \times w_g, \frac{H \times W}{h_g \times w_g}, C])$   
 $x = \text{transpose}(x, \text{permute\_axis} = [0, 2, 1, 3])$   
 $x = \text{reshape}(x, \text{shape} = [B \times \frac{H \times W}{h_g \times w_g}, h_g, w_g, C])$

---



---

### Algorithm 3 Batchwise NearFarMix Augmentation

---

$I_1 \leftarrow \text{Images}$  ▷ Input images  
 $D_1 \leftarrow \text{Depths}$  ▷ Input depths  
 $S_1 \leftarrow \text{Semantics}$  ▷ Input semantics  
 $\mathcal{U} \leftarrow \text{random\_uniform}$  ▷ Uniform distribution  
 $D_{min} \leftarrow 20(\text{KITTI}) \text{ or } 1.5(\text{NYUv2})$  ▷ Min depth  
 $D_{max} \leftarrow 60(\text{KITTI}) \text{ or } 6.5(\text{NYUv2})$  ▷ Max depth

#### # Roll batch for fast shuffling

$I_2 = \text{roll}(I_1, \text{shift} = 1, \text{axis} = 0)$  ▷ Roll Images  
 $D_2 = \text{roll}(D_1, \text{shift} = 1, \text{axis} = 0)$  ▷ Roll Depths  
 $S_2 = \text{roll}(S_1, \text{shift} = 1, \text{axis} = 0)$  ▷ Roll Semantics

#### # Depth threshold range for batch

$d_{min} = \max(\min(D_1, \text{axis} = (1, 2, 3)))$  ▷ Min depth  
 $d_{max} = \min(\max(D_1, \text{axis} = (1, 2, 3)))$  ▷ Max depth

#### # Threshold for Near-Far region

$B, H, W, C = \text{shape}(I_1)$   
 $\text{thr}_{min} = \max(D_{min}, d_{min})$  ▷ Clip min depth  
 $\text{thr}_{max} = \min(D_{max}, d_{max})$  ▷ Clip max depth  
 $\text{thrs} = \mathcal{U}(\text{shape} = [B, 1, 1, 1],$  ▷ Random thresholds  
 $\quad \text{min} = \text{thr}_{min},$   
 $\quad \text{max} = \text{thr}_{max})$

#### # Compute binary masks of regions for blending

$M_1 = D_1 \leq \text{thrs}$  ▷ Broadcasted Near region mask  
 $M_2 = D_2 > \text{thrs}$  ▷ Broadcasted pre-Far region mask  
 $M_o = M_1 \odot M_2$  ▷ Overlap region mask  
 $M_e = (1 - M_1) \odot (1 - M_2)$  ▷ Exclusive region mask

#### # Perform blending of regions

$I' = (I_1 \odot M_1)_{near} + ((I_2 \odot M_2) + (I_2 \odot M_e) - (I_2 \odot M_o))_{far}$  ▷ Augmented image  
 $D' = (D_1 \odot M_1)_{near} + ((D_2 \odot M_2) + (D_2 \odot M_e) - (D_2 \odot M_o))_{far}$  ▷ Augmented depth  
 $S' = (S_1 \odot M_1)_{near} + ((S_2 \odot M_2) + (S_2 \odot M_e) - (S_2 \odot M_o))_{far}$  ▷ Augmented semantics

---

## 2. Symbiotic Transformer

**Transformers:** Equation (1) presents the detailed mathematical expression for Symbiotic Transformer, which sym-

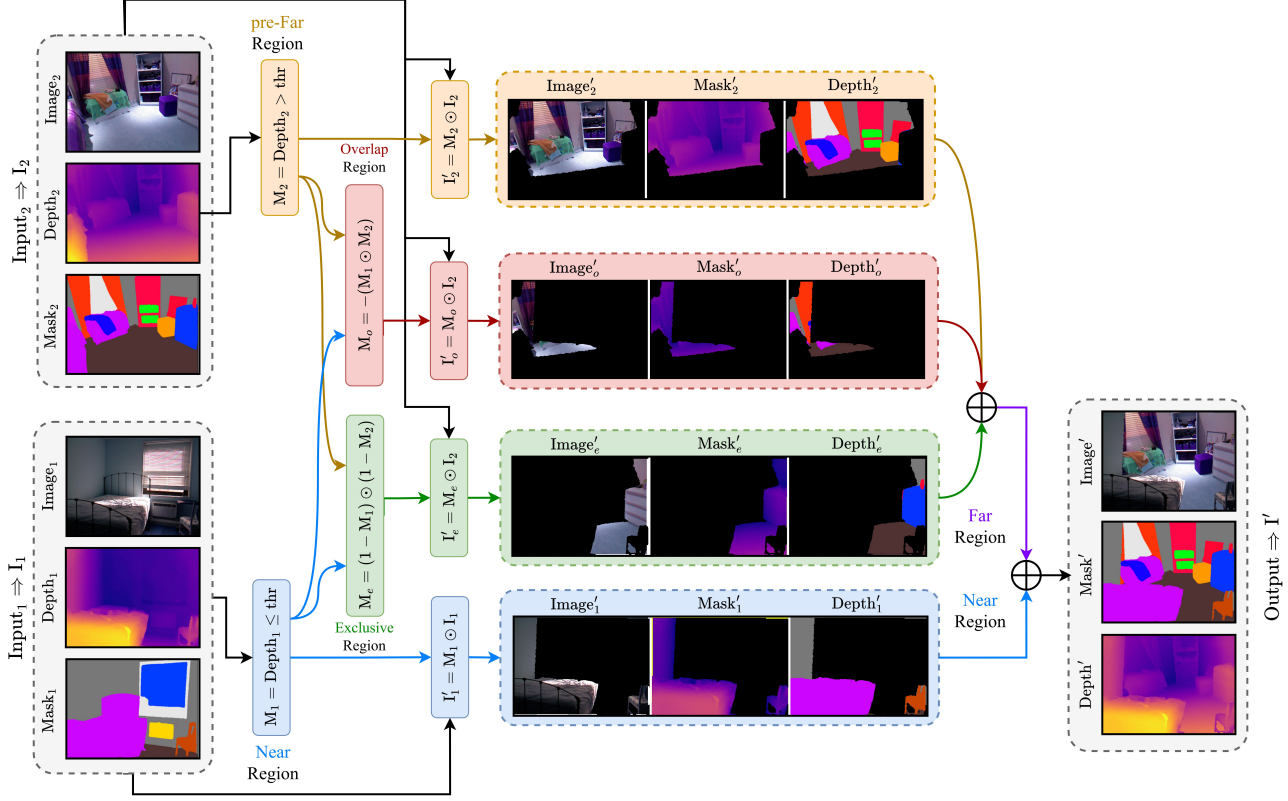


Figure 1. Proposed NearFarMix augmentation. The depth map undergoes thresholding to generate **Near** and **pre-Far** regions. Then, these regions are manipulated to produce **Overlap** and **Exclusive** regions. The **Far** region is then generated by combining the **pre-Far**, **Overlap**, and **Exclusive** regions, with the **Exclusive** region being subtracted and the remaining regions added. Finally, the **Near** and **Far** regions are combined to generate the augmented image.

biotically enhances both depth and semantics via local-global cross-attention. In the equation,  $F_x^Q$  represents query features of  $x$ ,  $F_y^{KV}$  denotes key-value features of  $y$ , and  $F_y$  signifies the output features contextualized by  $x$ . Specifically, for SGT,  $x = s$  and  $y = d$ , while for DGT,  $x = d$  and  $y = s$ . Moreover,  $DGT = LG-CAT_{DG}$  and  $SGT = LG-CAT_{SG}$  correspond to depth and semantics-guided local-global cross-attention transformers.

**Cross Attentions:** Under the hood, DGT and SGT employ semantics-guided cross attention (SG-CA) and depth-guided cross attention (DG-CA), respectively to contextualize features. SG-CA and DG-CA is mathematically elaborated in Eq. (2) and Eq. (3). In these equations,  $W$  represents the weight of the FFN layer,  $Q_x$  represents query features,  $K_y$  and  $V_y$  represent key and value features,  $Smx$  denotes softmax, and  $d$  is the query/key dimension.  $B$  represents relative positional bias, sampled similar to [1]. The Local-Global Cross-Attention Transformer (LG-CAT), employed by both SGT and DGT, can be implemented using Algorithm 4.

**Partition Operation:** Further, the algorithms for WindowPartition and GridPartition operations which are

used in Block and Grid attention also slightly different from the methods of Max-ViT [4], are detailed in Algorithm 1 and Algorithm 2, respectively. It is noteworthy that Max-ViT implemented these operations using einops [3].

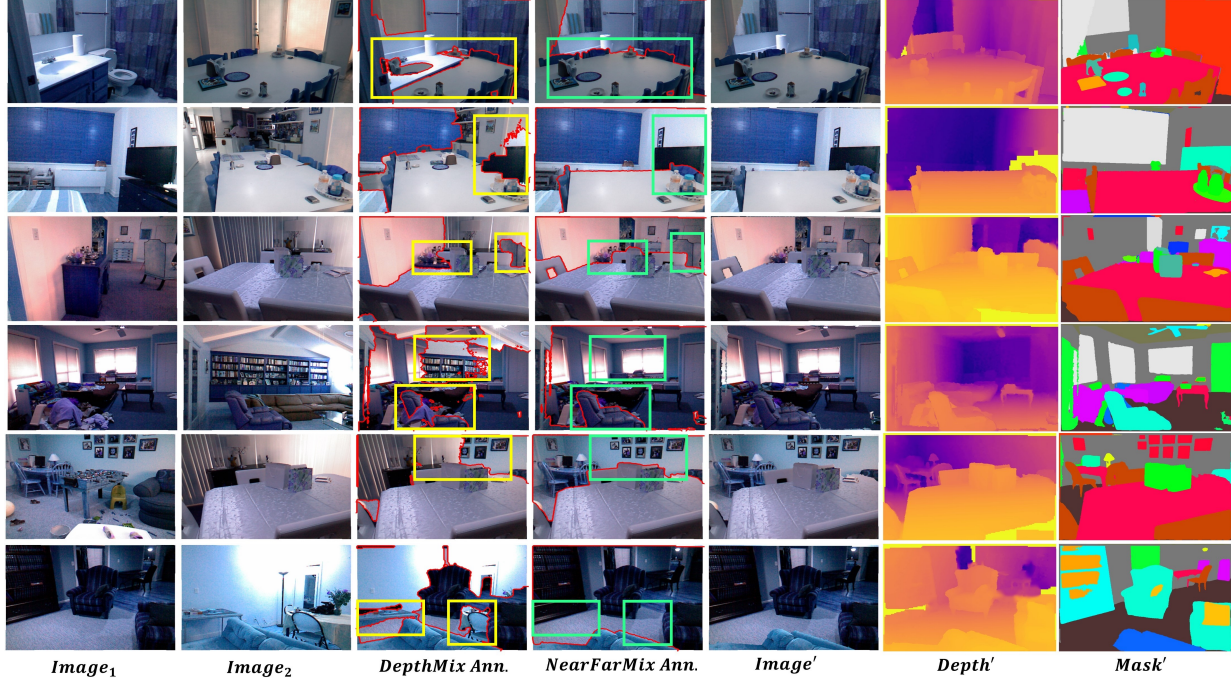


Figure 2. Additional qualitative comparisons between proposed NearFarMix and DepthMix augmentation on NYUv2 dataset.

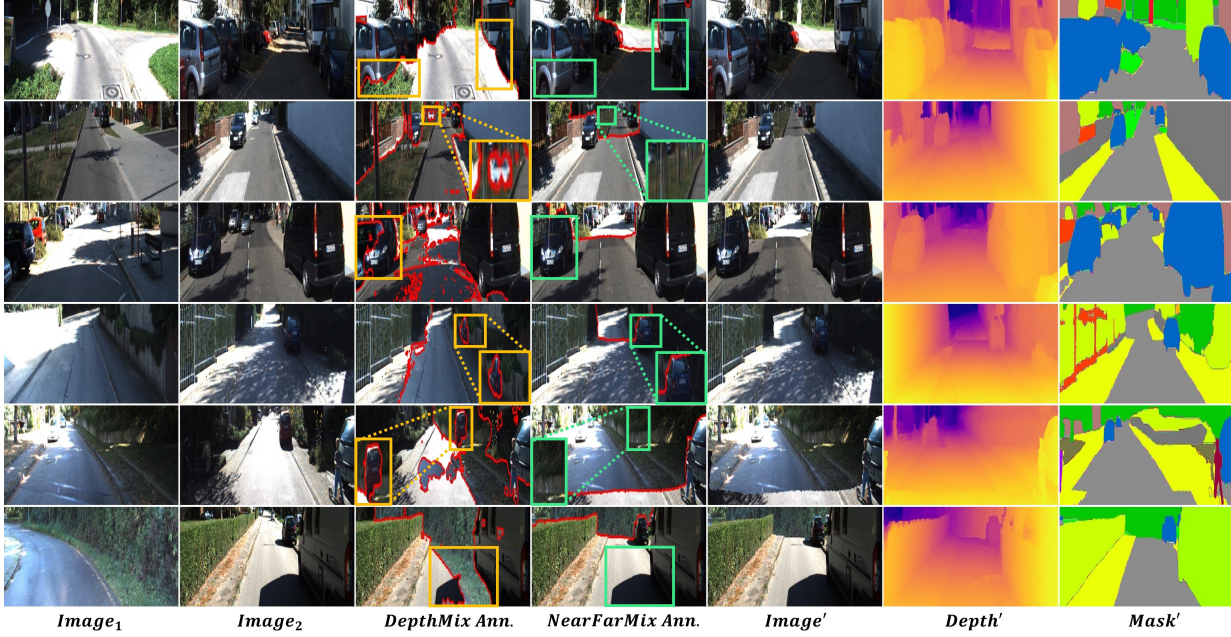


Figure 3. Additional qualitative comparisons between proposed NearFarMix and DepthMix augmentation on KITTI dataset.

$$\text{Sym-T}(\mathbf{F}'_d, \mathbf{F}'_s) = \begin{bmatrix} \mathbf{F}'_d = \mathbf{F}_d^{\text{Q}} = \mathbf{F}_d^{\text{KV}} = \mathbf{F}'_d \\ \mathbf{F}'_s = \mathbf{F}_s^{\text{KV}} = \mathbf{F}'_s \\ \mathbf{F}_s = \text{LG-CAT}_{DG}(\mathbf{F}_d^{\text{Q}}, \mathbf{F}_s^{\text{KV}}) \\ \mathbf{F}_d = \text{LG-CAT}_{SG}(\mathbf{F}'_s, \mathbf{F}_d^{\text{KV}}) \end{bmatrix} \quad (1)$$

$$\begin{aligned} \text{SG-CA} &= \text{CA}(\mathbf{Q}_s, \mathbf{K}_d, \mathbf{V}_d) \\ &= \text{Smx} \left( \frac{\mathbf{Q}_s \mathbf{K}_d^T}{\sqrt{d}} + \mathbf{B} \right) \mathbf{V}_d \\ &= \text{Smx} \left( \frac{(\mathbf{F}'_s \cdot \mathbf{W}_s^{\text{Q}})(\mathbf{F}_d^{\text{KV}} \cdot \mathbf{W}_d^{\text{K}})^T}{\sqrt{d}} + \mathbf{B} \right) (\mathbf{F}_d^{\text{KV}} \cdot \mathbf{W}_d^{\text{V}}) \end{aligned} \quad (2)$$



---

**Algorithm 4** Local-Global Cross-Attention Transformer (LG-CAT)
 

---

$F_x^Q, F_y^{KV} \leftarrow inputs$  ▷ Input features  
 $x \leftarrow F_x^Q$  ▷ Query features of depth/semantics  
 $y \leftarrow F_y^{KV}$  ▷ Key-value features of semantics/depth  
 $i \leftarrow 0$  ▷ Initialize counter

**while**  $i \neq 2$  **do**

**# Block Cross Attention**

$x_1 = \text{layer\_norm}(x)$   
 $y_1 = \text{layer\_norm}(y)$   
 $x_{1,q} = \text{FFN}(\text{window\_partition}(x_1))$  ▷ Query gen.  
 $y_{1,k}, y_{1,v} = \text{FFN}(\text{window\_partition}(y_1))$  ▷ Key Value  
 $y_2 = \text{CA}(x_{1,q}, y_{1,k}, y_{1,v})$  ▷ Apply cross-attention  
 $y_2 = y_1 + \text{window\_reverse}(\text{FFN}(y_2))$  ▷ Residual

**# Grid Cross Attention**

$y_2 = \text{layer\_norm}(y_2)$   
 $x_{2,q} = \text{FFN}(\text{grid\_partition}(x_1))$  ▷ Query  
 $y_{2,k}, y_{2,v} = \text{FFN}(\text{grid\_partition}(y_2))$  ▷ Key Value  
 $y_3 = \text{CA}(x_{2,q}, y_{2,k}, y_{2,v})$  ▷ Apply cross-attention  
 $y_3 = y_2 + \text{grid\_reverse}(\text{FFN}(y_3))$  ▷ Residual  
 $y = y_3$  ▷ Reset variable for loop  
 $i = i + 1$  ▷ Increment counter

**end while**

**# FusedMBConv - Channel Attention**

$\hat{y} = \text{DWConv}3 \times 3(y)$  ▷ Depthwise convolution  
 $\hat{y} = \text{GELU}(\hat{y})$  ▷ Apply activation  
 $\hat{y} = \text{SE}(\hat{y})$  ▷ Squeeze-Excitation  
 $\hat{y} = \text{Conv}1 \times 1(\hat{y})$  ▷ Convolution  
 $\hat{y} = y + \hat{y}$  ▷ Residual

$output \leftarrow \hat{y}$

---

## References

- [1] Ali Hatamizadeh, Hongxu Yin, Jan Kautz, and Pavlo Molchanov. Global context vision transformers. *arXiv preprint arXiv:2206.09959*, 2022. 2
- [2] Lukas Hoyer, Dengxin Dai, Yuhua Chen, Adrian Koring, Suman Saha, and Luc Van Gool. Three ways to improve semantic segmentation with self-supervised depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11130–11140, 2021. 1
- [3] Alex Rogozhnikov. Einops: Clear and reliable tensor manipulations with einstein-like notation. In *International Conference on Learning Representations*, 2021. 2
- [4] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. Maxvit: Multi-axis vision transformer. *arXiv preprint arXiv:2204.01697*, 2022. 2

$$\begin{aligned}
 \text{DG-CA} &= \text{CA}(\mathbf{Q}_d, \mathbf{K}_s, \mathbf{V}_s) \\
 &= \text{Smax} \left( \frac{\mathbf{Q}_d \mathbf{K}_s^T}{\sqrt{d}} + \mathbf{B} \right) \mathbf{V}_s \\
 &= \text{Smax} \left( \frac{(\mathbf{F}_d^Q \cdot \mathbf{W}_Q^d)(\mathbf{F}_s^{KV} \cdot \mathbf{W}_K^s)^T}{\sqrt{d}} + \mathbf{B} \right) (\mathbf{F}_s^{KV} \cdot \mathbf{W}_V^s)
 \end{aligned} \tag{3}$$

## 3. Architecture Details

The architectural specifications, encompassing input, output, layer name, and layer details, are succinctly laid out in Table 1. Here,  $E$  and  $D$  represent the input/output of the encoder/decoder, while  $ST$  and  $N$  denote the Symbiotic Transformer and the Neck.

Table 1. Architectural Specifications of proposed method where  $h, w$  signify attention heads and window size;  $Conv$  denotes 2D convolution with  $k, s, c$  as kernel size, stride size, and output channels;  $act$  and  $norm$  represent activation and normalization types;  $sc$  indicates upscale size.

Input Size: $H \times W \times 3$				
Layer Name	Input	Output	Output Size	Architecture
Stem	$Image$	$E_0$	$\frac{H}{4} \times \frac{W}{2} \times 128$	$Conv(c=128, k=3, s=2)$ $Conv(c=128, k=3, s=1)$
Encoder Stage 1	$E_0$	$E_1$	$\frac{H}{8} \times \frac{W}{4} \times 128$	$[Max-ViT-Block(h=4, w=7)] \times 2$
Encoder Stage 2	$E_1$	$E_2$	$\frac{H}{16} \times \frac{W}{8} \times 256$	$[Max-ViT-Block(h=8, w=7)] \times 6$
Encoder Stage 3	$E_2$	$E_3$	$\frac{H}{32} \times \frac{W}{16} \times 512$	$[Max-ViT-Block(h=16, w=7)] \times 14$
Encoder Stage 4	$E_3$	$E_4/D_4$	$\frac{H}{32} \times \frac{W}{32} \times 1024$	$[Max-ViT-Block(h=16, w=7)] \times 2$
Decoder Stage 3	$(D_4, E_3)$	$D_3$	$\frac{H}{32} \times \frac{W}{32} \times 512$	Upsample( $sc=2$ )
				Concat( $[E_3, D_4]$ , axis=-1)
				Conv( $c=512, k=3, s=1, norm='layer', act='gelu'$ )
Decoder Stage 3	$(D_3, E_2)$	$D_2$	$\frac{H}{16} \times \frac{W}{16} \times 256$	Upsample( $sc=2$ )
				Concat( $[E_2, D_3]$ , axis=-1)
				Conv( $c=256, k=3, s=1, norm='layer', act='gelu'$ )
Decoder Stage 1	$(D_2, E_1)$	$D_1$	$\frac{H}{4} \times \frac{W}{4} \times 128$	Upsample( $sc=2$ )
				Concat( $[E_1, D_2]$ , axis=-1)
				Conv( $c=128, k=3, s=1, norm='layer', act='gelu'$ )
Neck	$D_1$	$(N_d, N_s)$	$(\frac{H}{4} \times \frac{W}{4} \times 150,$ $\frac{H}{4} \times \frac{W}{4} \times 150)$	$([Conv(c=150, k=3, s=1, norm='layer', act='gelu')] \times 2,$ $[Conv(c=150, k=3, s=1, norm='layer', act='gelu')] \times 2)$
Symbiotic Transformer	$(N_d, N_s)$	$(ST_d, ST_s)$	$(\frac{H}{4} \times \frac{W}{4} \times 150,$ $\frac{H}{4} \times \frac{W}{4} \times 150)$	$[Block-Cross-Attention(h=4, w=7)] \times 2$ $[Grid-Cross-Attention(h=4, w=7)] \times 2$
				$[FusedMBCConv] \times 1$
Head	$(ST_d, ST_s)$	$(Depth, Semantics)$	$(H \times W \times 1,$ $H \times W \times 150)$	$([Conv(k=3, s=1, act='sigmoid')]$ Upsample( $sc=4$ ) $[Conv(k=3, s=1, act='softmax')]$ Upsample( $sc=4$ )) $\times 2$

$(Depth: H \times W \times 1, Semantics: H \times W \times 150)$