



BUET

Semi-Supervised Semantic Depth Estimation using Symbiotic Transformer and NearFarMix Augmentation

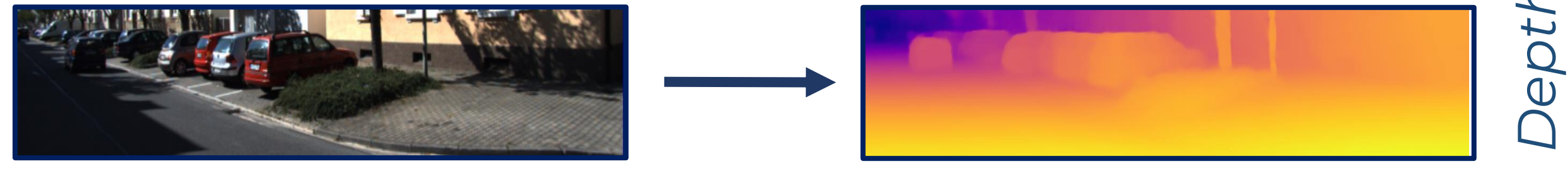
Md Awsafur Rahman, Shaikh Anowarul Fattah

Dept. of EEE, BUET, Bangladesh



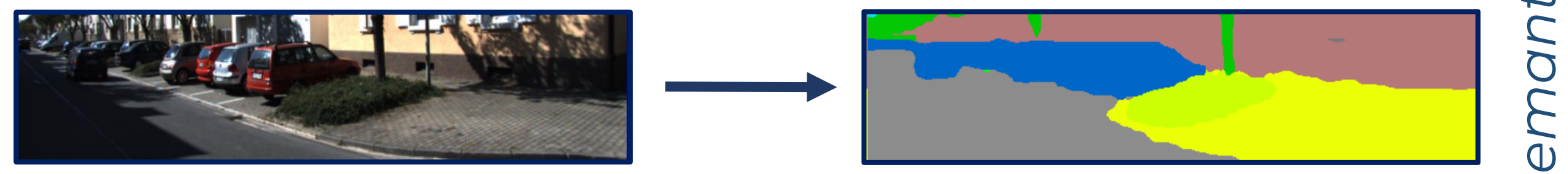
Semantic Depth Estimation ?

- **Semantic Depth**: Estimate both **depth** and **semantics**
- **Depth**: Process of estimating the depth or distance information for each pixel in 2D image.



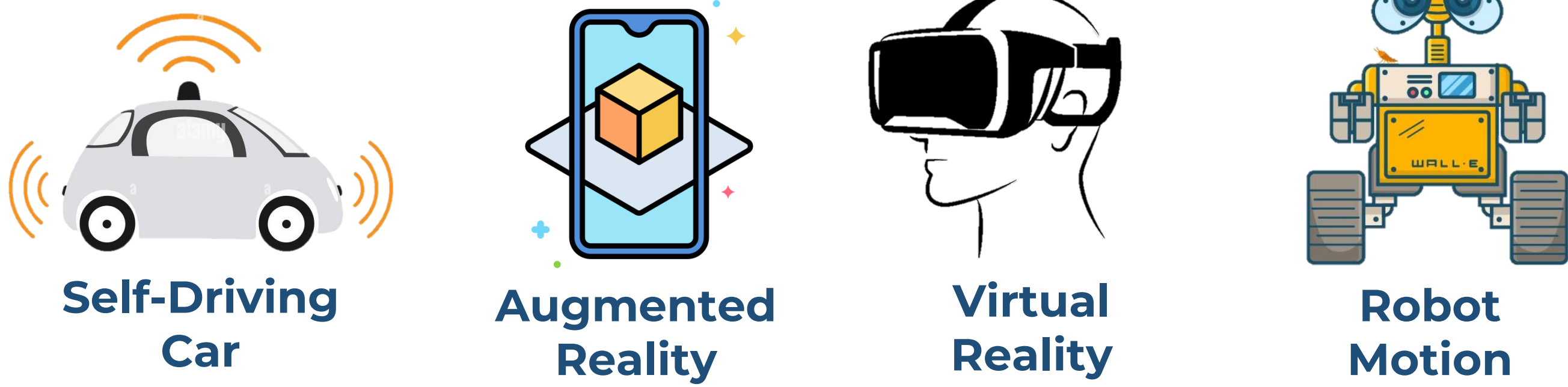
Depth

- **Semantics**: Process of classifying each pixel in 2D image to a specific category.



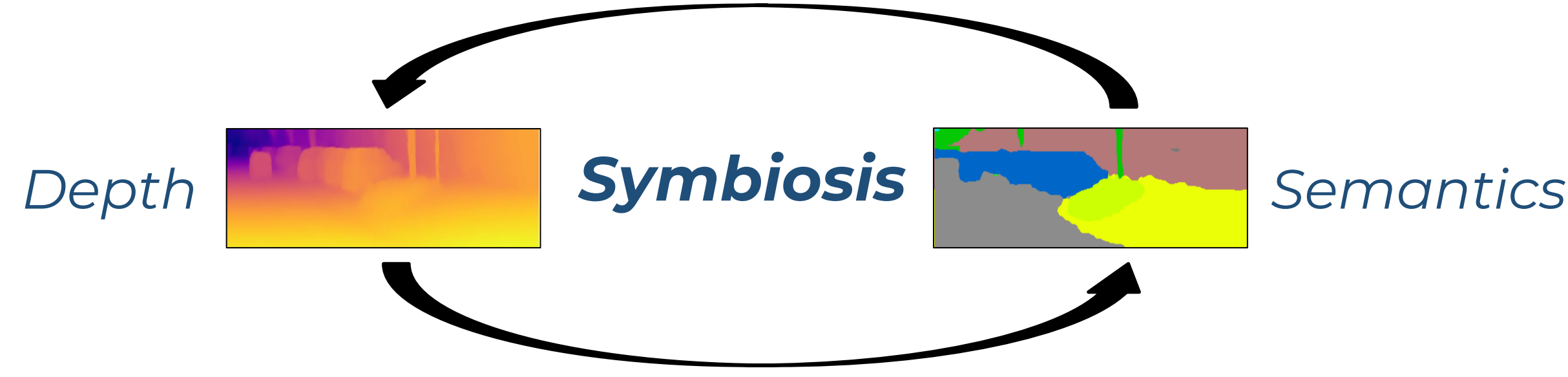
Semantics

Applications



Depth Semantics Symbiosis

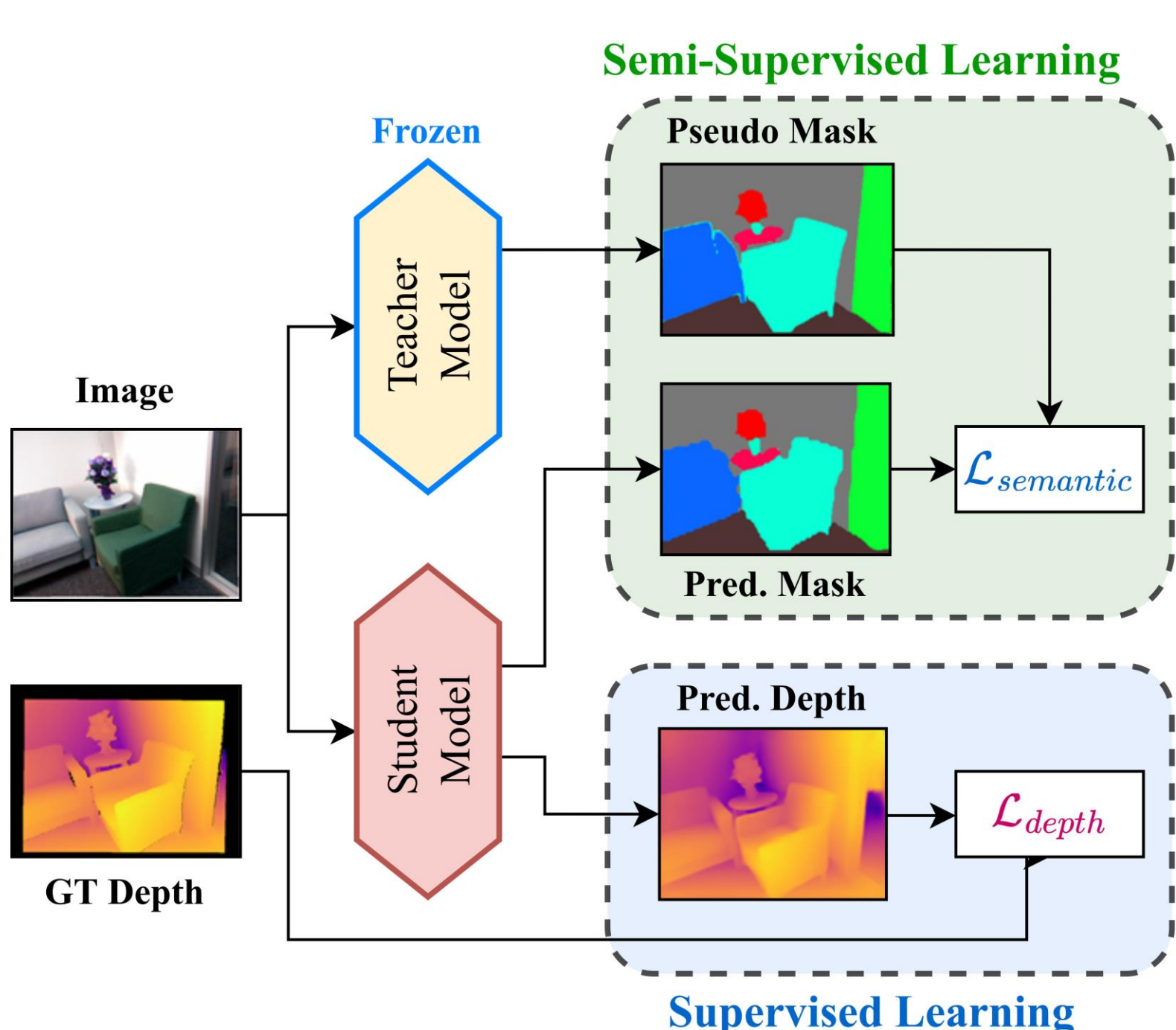
- **Depth** can help **Semantics** by adding a 3D view to scenes, thus clarifying the spatial relations of similar objects at varying distances.
- **Semantics** can help **Depth** by providing object categories and boundaries, essentially facilitating consistent and sharp-edged depth.



Contributions

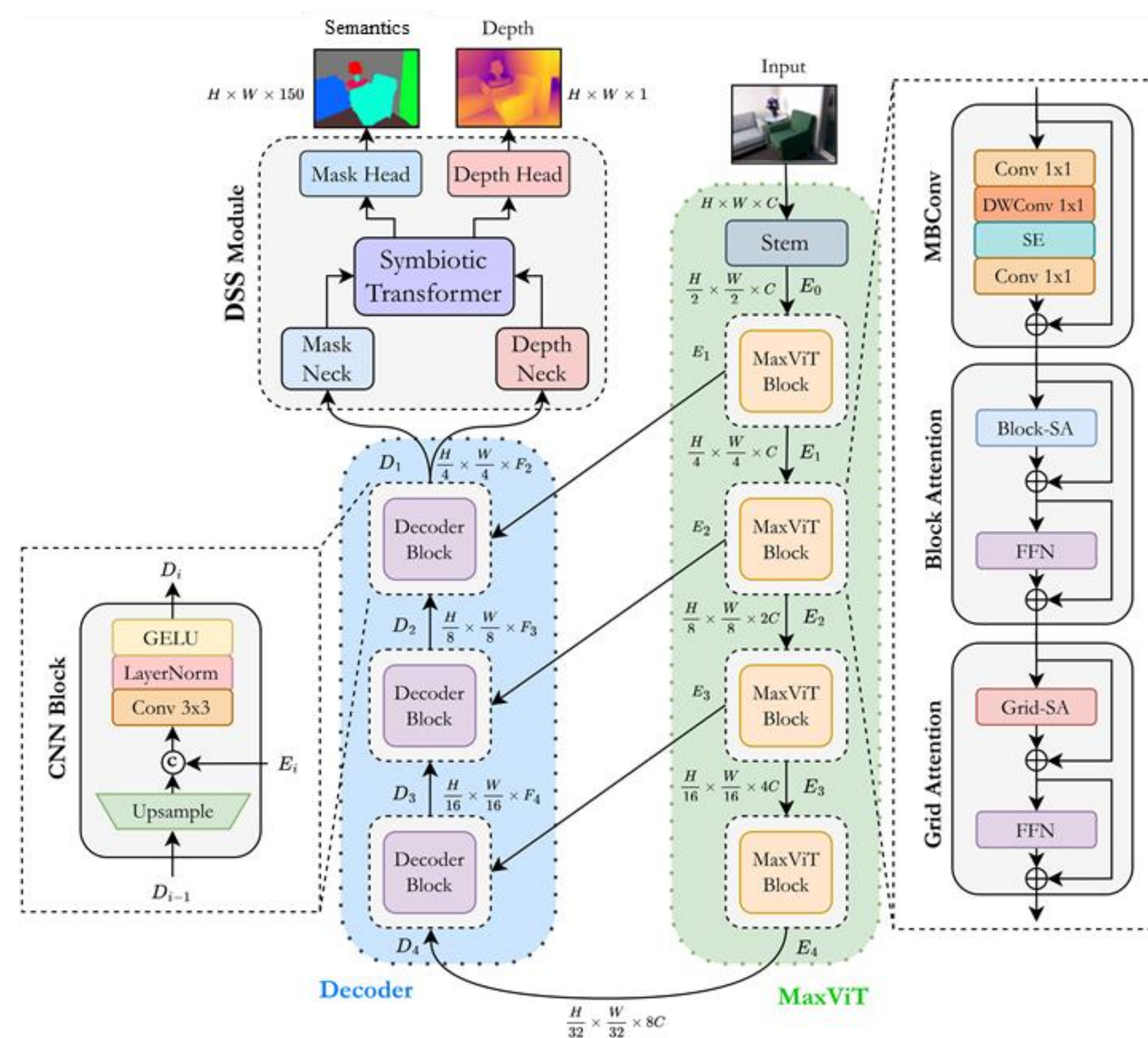
- **Semi-supervised dataset-agnostic** strategy to mitigate semantic label scarcity.
- **Symbiotic Transformer** to resolve limited symbiosis by exchanging information between depth and semantics tasks within both local and global contexts.
- **NearFarMix augmentation** to tackle overfitting in both depth and semantics tasks while solving existing issues such as loss of object integrity, limited diversity, and limited control.

Dataset-Agnostic Semi-Supervised Strategy

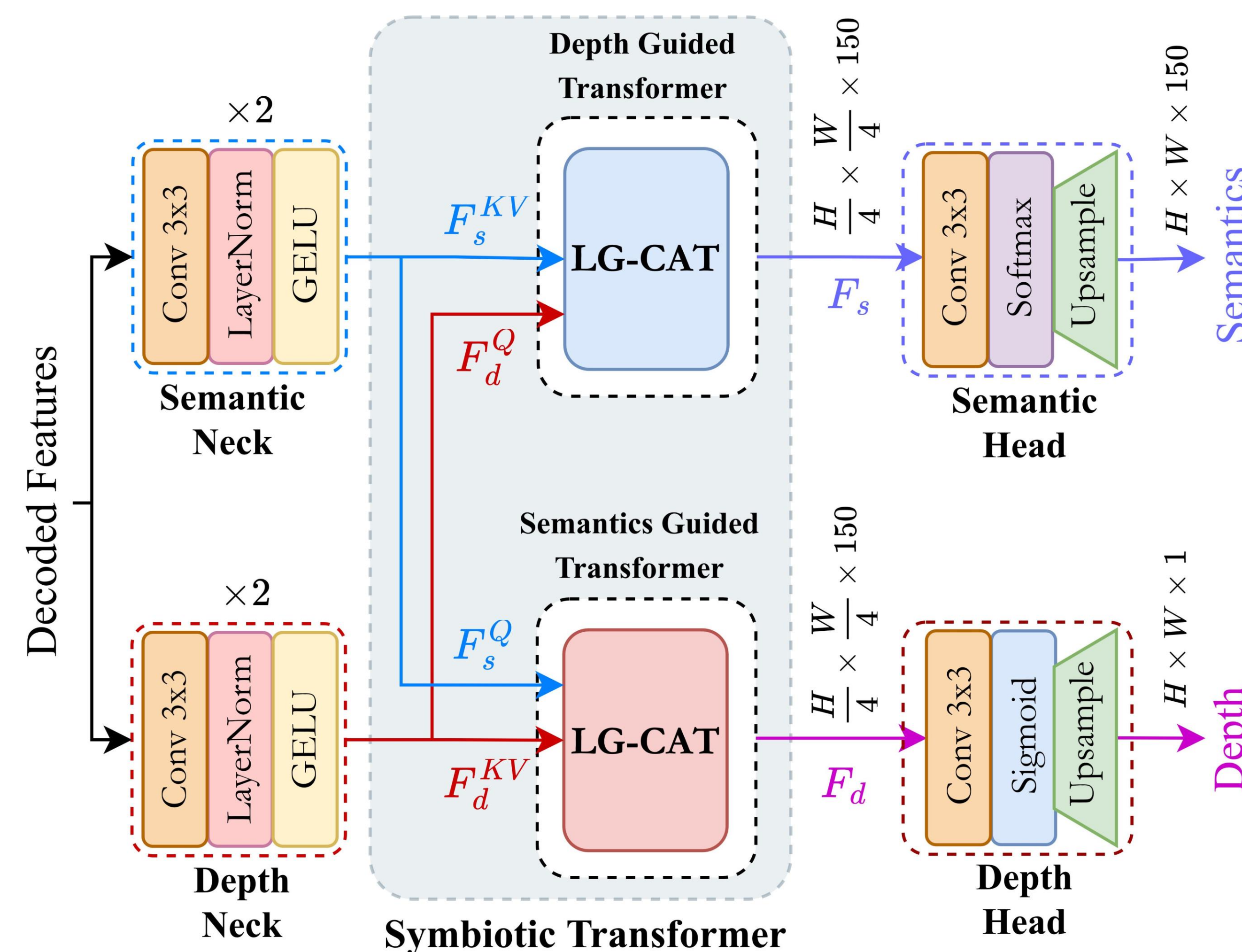


- The teacher model maintains a fixed number of classes in semantic labels across datasets, ensuring a **dataset-invariant architecture**.
- The semi-supervised nature enables its application to datasets with only depth information, **without requiring semantic labels**.

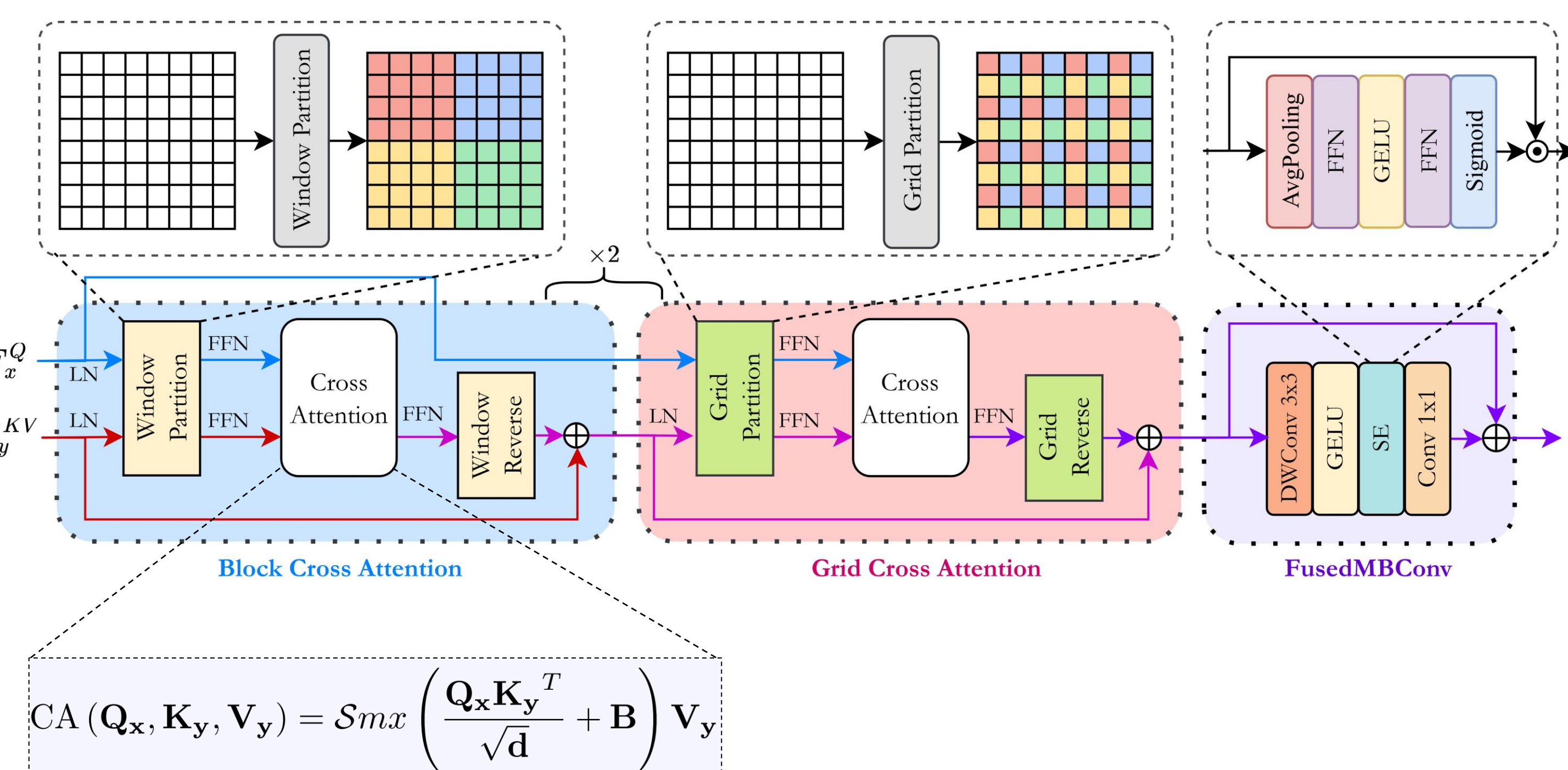
Architecture Overview



Depth Semantics Symbiosis (DSS) Module



Local Global Cross Attention Transformer (LG-CAT)

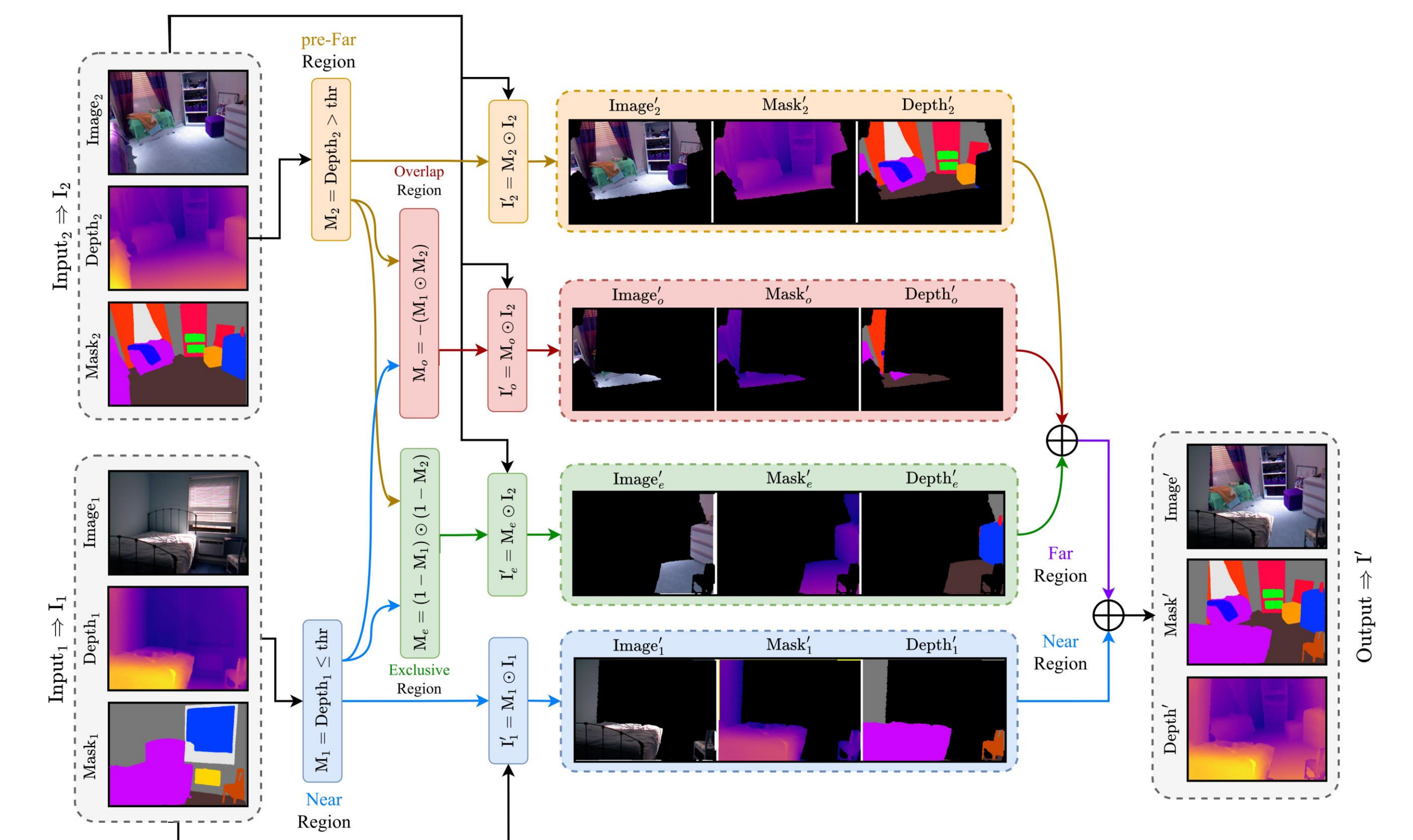


The Symbiotic Transformer block can be mathematically expressed as:

$$\begin{cases} i\mathbf{F}_y^{\text{KV}} = \text{Block-Cross-Attention}(\mathbf{F}_x^{\text{Q}}, i-1\mathbf{F}_y^{\text{KV}}) \\ i+1\mathbf{F}_y^{\text{KV}} = \text{Grid-Cross-Attention}(\mathbf{F}_x^{\text{Q}}, i\mathbf{F}_y^{\text{KV}}) \end{cases} \times N_s$$

$$\mathbf{F}_y = \text{FusedMBConv}(N+1\mathbf{F}_y^{\text{KV}})$$

NearFarMix Augmentation



NearFarMix Algorithm

- **Compute binary masks of regions for blending**
- $M_1 = D_1 \leq thr_s \triangleright$ Broadcasted Near region mask
- $M_2 = D_2 > thr_s \triangleright$ Broadcasted pre-Far region mask
- $M_o = M_1 \odot M_2 \triangleright$ Overlap region mask
- $M_e = (1 - M_1) \odot (1 - M_2) \triangleright$ Exclusive region mask

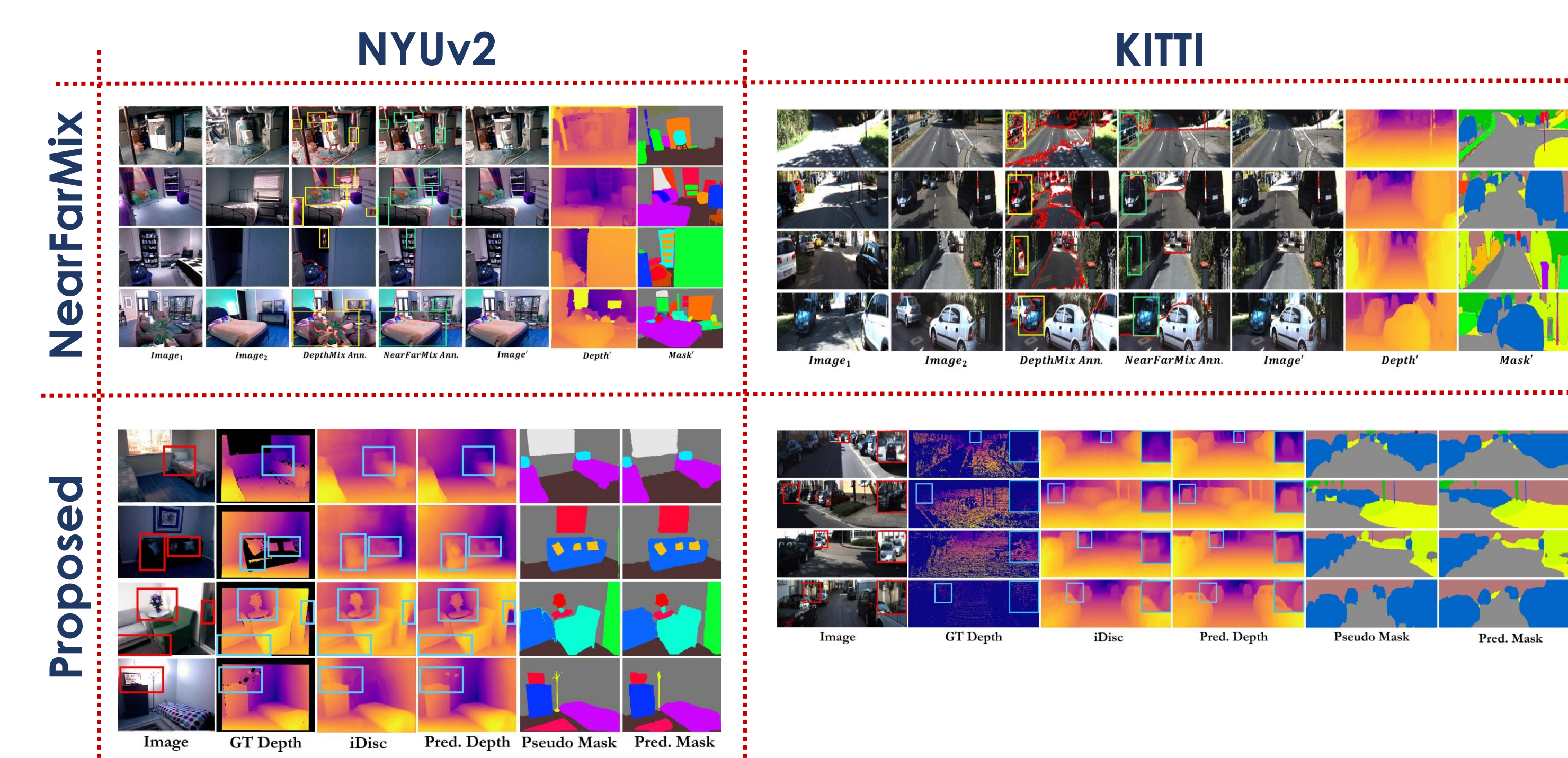
Perform blending of regions

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

Ablation Study

DSS					NearFarMix							
Dataset	Method	RMS ↓	Abs Rel ↓	$\delta_1 \uparrow$	mIoU ↑	Dataset	Method	RMS ↓	Abs Rel ↓	$\delta_1 \uparrow$	mIoU ↑	
KITTI	Baseline	2.295	0.056	0.966	0.615	KITTI	Baseline	2.045	0.051	0.977	0.663	
	SW-Map	2.275	0.055	0.967	0.623		CutMix	2.441	0.059	0.962	0.671	
	SIG	2.105	0.053	0.969	0.650		ClassMix	2.302	0.056	0.967	0.693	
	AG-MMD	2.096	0.052	0.973	0.677		AffineMix	2.102	0.053	0.969	0.711	
	CCAM	2.080	0.050	0.976	0.695		DepthMix	2.104	0.054	0.968	0.708	
	DSS*	1.984	0.048	0.979	0.731		CutDepth	2.968	0.059	0.973	0.640	
NYUv2	Baseline	0.323	0.089	0.934	0.519	NYUv2	V-CutDepth	2.965	0.050	0.975	0.645	
	SW-Map	0.314	0.087	0.933	0.522			NearFarMix*	1.984	0.048	0.979	0.731
	SIG	0.309	0.086	0.937	0.555		Baseline	0.296	0.082	0.945	0.568	
	AG-MMD	0.302	0.085	0.939	0.571		CutMix	0.331	0.094	0.927	0.572	
	CCAM	0.300	0.083	0.943	0.589		ClassMix	0.320	0.090	0.930	0.585	
	DSS*	0.289	0.080	0.948	0.620	AffineMix	0.314	0.086	0.941	0.603		
						DepthMix	0.305	0.084	0.943	0.601		
						CutDepth	0.294	0.082	0.946	0.542		
						V-CutDepth	0.290	0.081	0.946	0.546		
						NearFarMix*	0.289	0.080	0.948	0.620		

Qualitative Results



Quantitative Result

NYUv2					KITTI				
Method	Abs Rel ↓	RMS ↓	$\log_{10} \downarrow$	$\delta_1 \uparrow$	Method	Abs Rel ↓	RMS ↓	$RMS_{log} \downarrow$	$\delta_1 \uparrow$
Eigen et al.	0.158	0.641	-	0.769	Eigen et al.	0.203	6.307	0.270	0.702
DORN	0.115	0.509	0.051	0.828	DORN	0.072	2.727	0.120	0.932
BTS	0.110	0.392	0.047	0.885	BTS	0.059	2.756	0.090	0.956
TransDepth	0.106	0.365	0.045	0.900	TransDepth	0.064	2.755	0.098	0.956
DPT	0.110	0.367	0.045	0.904	Adabins	0.058	2.360	0.088	0.964
Adabins	0.103	0.364	0.044	0.903	DPT	0.060	2.573	0.088	0.959
P3Depth	0.104	0.356	0.043	0.898	NeWCRFs	0.052	2.129	0.079	0.974
NeWCRFs	0.095	0.334	0.041	0.922	PixelFormer	0.051	2.081	0.077	0.976
PixelFormer	0.090	0.322	0.039	0.929	iDisc	0.050	2.067	0.077	0.977
iDisc	0.086	0.313	0.037	0.940	Proposed	0.048	1.984	0.075	0.979
Proposed	0.080	0.289	0.034	0.948					