ISSN: 2278-1862



Journal of Applicable Chemistry

2025, 14 (2): 395-470 (International Peer Reviewed Journal)





CNN - 66C

IAM (Intelligence Augmented /Assisted Method(s)) Transformers --architectures & Fits (2025)

Information Source	sciencedirect.com;	
S. Narasinga Rao M D	K. Somasekhara Rao, Ih D	R. Sambasiva Rao, Ih D
Associate Professor,	Dept. of Chemistry,	Dept. of Chemistry,
Emergency Medicine dept.,	Acharya Nagarjuna Univ.,	Andhra University,
Andhra Medical College,	a Medical College, Dr. M.R.Appa Rao Campus,	
King George Hospital	Nuzvid-521 201, India	India
Visakhapatnam, A.P., India		
snrnaveen007@gmail.com	sr_kaza1947@yahoo.com	rsr.chem@gmail.com
(+91 98 48 13 67 04)	<u>(+91 98 48 94 26 18)</u>	(+91 99 85 86 01 82)

Conspectus: In the year 2017, Ashish Vaswani et al. published a paper entitled "Attention is All You Need". It revolutionised sequence data modelling. The new model gained popularity as a Transformer net (TransF Net) or Transformer neural network (TransF NN). The two important modules are attention layer and MLP-NN to carry out Natural Language processing (NLP). The evolution of architecture of TransF NN, attention mechanism, and hybridization with other methods, during these few years, revolutionized computational modelling paradigm. This made a niche in Data Science dealing with multi-modal data (viz.Text, numerical time-series, sound (speech), and image/video sequence) with local and global inter-dependencies.

Natural Language processing (NLP) : Some of earlier models in use were word2vector (2013), MLP, RecNN (1997-2015), capsuleNN, LSTM, GRU and Transformer (2017-). The advances in modules with LSTM are Bidirectional-LSTM, LSTM+attention, LSTM+seq2seq model, LSTM+Reinforced Lrning, LSTM+self sup Lrning, LSTM+Tranformer, peep- LSTM, and Hierarchical LSTM.

The language models (LMs) are also categorized as Large language models (LLMs), Small language models (SLMs) and Large/Small language models (LSLMs) based on training data size and number of parameters.

Transformer models: Attention and MLP NN are the two basic modules of a Transformer invoked in 2017 by Vaswani. Bahadaname (2014) and Luong (2015) invoked the concept of attention in pre-transformer era. The Transformer model with self-attention layers achieved state-of-the-art results in machine translation and completely replaced RNNs.

The frames proposed during this decadal period are Generative Pre-trained Transformer (GPT-x: x=1 to 4). BART (Bidirectional and Auto-Regressive Transformers), BERT (Bidirectional Encoder Representations from Transformers), T5 (Text-to-Text Transfer Transformer), PaLM (Pathways Language Model), CLIP (Contrastive Language-Image Pre-training), DALL·E (Multimodal models), BARD, and LLaMA (Large Language Model Meta AI)-2023.

The evolution of Tranformers led to LinFormer, Longformer and Performer with high end technical features and applications.

A few of Transformer neural nets (TransF-NN) or Transformer nets (TransF-N) architectures used in this state-of-knowledge-methods-module for dataToknowledge transformation are

- APT: Alarm Prediction Transformer,
- *E* FastPCI, Swin Transformer Transformer,
- la Cross-scale prototype learning transformer (Cplformer),
- lacktriangle Multiscale Network (MSNet),
- lacktrian (Automatic Fusion Networks (AutoFuse),
- *⊖* U-shape transformer,
- *△* cross-wise transformer module (CTM),
- △ Transformer with Sliding-Window Dissimilarity Cross-Attention (SWDCA),
- la Neural Networks To Vision Transformers (NN2ViT),
- △ Wavelet-domain Convolution (WeConv), Forest2Seq, RS-MOCO,
- *▲* RBMDC-Net,
- la RISurConv: Rotation Invariant Surface,
- 👃 S-JEPA,
- △ REDIR: Refocus-free Event-based De-occlusion Image Reconstruction,
- △ DOLFIN: Diffusion Layout Transformers Without Autoencoder, FasterSTS:
- 🚨 A Faster Spatio-Temporal Synchronous Graph Convolutional Networks,
- lacktrian Revenue A Robust Polar Transformation Network,
- line Comparison and Alignment Network,
- △ CVT-Occ: Cost Volume Temporal Fusion for 3D Occupancy,

- △ Transformer network integrated into a chatbot interface,
- △ Swin transformers with the U-Net architecture incorporating residual blocks (RBs) and
- *△* Attention mechanism and Trident Transformers (TT).











Swin Transformer-based coding structure

- ✓ Swin Transformer decomposes the entire encoder into several layers, each consisting of multiple small Transformer blocks
- ✓ It initially partitions image into several image blocks and encodes them into vectors
- ✓ These vectors are then passed through a linear coding layer to produce lowdimensional
- ✓ vectors, which are utilized in the subsequent hierarchical Transformer block
- Swin Transformer block comprises multiple grouped convolutional layers and several drift window self-attention modules



Reference View

Projection View Multi-view projection effect based on optical flow

Source View

Source





Overview of CrossDiff framework for generating human motion from textual Descriptions

✓ Framework incorporates both 3D and 2D motion data, using unified encoding and cross-decoding components to process mixed representations obtained from random projection.



Fig. 3: Overview of Mixture Sampling

- Original noise is sampled from a 2D gaussian distribution.
- From time-step T to α, crossdiff predicts the clean 2D motion x2d,0 and diffuses it back to x2d,t-1. In the remaining α steps, crossdiff denoises in the 3D domain and
 Finally actains the clean 2D motion
- ✓ Finally obtains the clean 3D motion

 \checkmark









SemiVDN

- ✓ Mean teacher scheme with a student model
- ✓ Teacher model
- Prior-guided Temporal Decoupling Experts
 - o to decompose the physical components that make up a snow video in a temporal spirit
- ✓ After that compute supervised losses for labeled data and unsupervised losses for unlabeled data
- ✓ Based on the decomposed component features (F' B and F ' S) in representation space,
- ✓ Develop a Distribution-driven Contrastive Regularization to highlight the snow-invariant information by replacing the snow-specific feature in ultra-positive samples replacing the background in negative samples



- A Model. Physics Transformer Block (PTB)
- A Module. Temporal Decoupling Experts module
 - To generate physics-specific components (i.e. S, A and T) for recovery.
- **A** Temporal Decomposition Router

To compute the temporal weights Qij from the temporal dimension, which are subsequently employed to compute a linear combination of all input temporal tokens and Qij.

- **A** Temporal Decomposition Router
 - To convexly combine all the component tokens.
- la Output combined features
- △ X k and physics-specific features pj k are subsequently input into the
 - Prior-guided recovery module decoder
 - To generate the ultimate desnowed results



Fig. 6: Samples of the proposed real-world video dataset for video snow removal.



specific weights $\hat{\theta}'$ using the hyper-network. **Bottom:** video decoding. NeRV-Dec generates final NeRV weights θ' and reconstruct video \hat{x} .

















Fig. 1. Illustration of existing empirically-designed fusion strategies (a-f) and our data-driven f strategy (g) in unsupervised and semi-supervised settings. (a) Early fusion: I_f and I_m are concatenal input. (b) Middle fusion: I_f and I_m enter separate encoders with intermediate features fused. (c) fusion: I_f and I_m enter separate networks with resultant features fused. (d) Loss fusion: ψ and S_f/S mutually constrained by joint loss functions. (e) Feature fusion: multi-task networks are used with p feature shared. (f) Input fusion: S_f/S_m are fed as input for registration. (g) Data-driven fusion: f strategy is optimized during training. FG = Fusion Gate (FG) module; I_f = fixed image; I_m = m image; ψ = registration output; S_f/S_m = segmentation output for I_f/I_m (for semi-supervised registration)







Optimal Transport (CPOT)

- Without CPOT, all click prompts tend to converge to one point, resulting in homogeneous prompt-activated masks and inferior mask prediction.
- (b) With the proposed CPOT, click prompts are encouraged to focus on distinct visual regions.
- Consequently, our model with CPOT predicts a more accurate mask by integrating diverse prompt-activated masks.



Overview of proposed Click Prompt Learning with Optimal Transport (CPlot)

- ✓ Given input image, click disk maps, and previous mask, the Image Encoder extracts visual features F.
- ✓ The Click Encoder initializes click prompts Pc with click coordinates.
- (a) Prompt-Pixel Alignment aims to align click prompts Pc with the visual features F in the feature space.
- (b) Click Prompt Optimal Transport adopts optimal transport plan to generate optimized mask S* from vanilla prompt-activated mask S. A lightweight mask decoder is used to implicitly analyze optimized prompt-activated mask with visual features and make mask predictions





Method	#Param (M)	MACs (G)
STANet-BAM(ResNet18) [4]	12.2	49.2
STANet-PAM(ResNet18) [4]	12.2	50.2
DTCDSCN(SE-Res34) [7]	41.1	60.9
L-Unet [31]	8.5	
CDNet [32]	14.3	
MSCANet [33]	16.4	
BiT(ResNet18) [22]	3.0	35.0
SNUNet [14]	3.0	46.9
ChangeFormer(MiT-b1) [21]	13.9	26.4
IFN(VGG-16) [9]	36.0	316.5
FHD [34]	11.8	
ChangeStar(MiT-b1) [35]	18.4	33.7
Xu et al. [11]	61.4	
ChangerEx(ResNet18) [18]	11.4	23.9
ChangeStar(ResNet18) [35]	16.4	32.7
CDNeXt [36]	39.4	31.5
TransUNetCD [23]	95.5	
BAT [27]	6.9	40.3
SWDCA Network	5.4	25.0

bi-temporal image pairs with a resolution of 512 × 512 pixels. The **optimal** value is indicated in red font, whereas the **second-best** value is represented in blue font. # means the number of.









Fig. 6: Comparison of image restoration on *low+haze+rain* (top) and *low+haze+snow* (bottom) synthetic samples.



Fig. 7: Comparison of image restoration on low+haze+rain (top) and low+haze+snow (bottom) samples in real-world scenarios. * represents the utilization of original weights published in the author's code.















Fig. 3: Our neural network architecture comprises five RISurConv layers to extract rotation invariant features followed by a Transformer Encoder to enhance the learnt features before fully connected layers for object classification. We add a decoder with skip connections for segmentation task.

	Method	Format	Input Size	Params.
Traditional	VoxNet [21]	voxel	30^{3}	0.90M
	SubVolSup [24]	voxel	30^{3}	17.00M
	PointNet [23]	xyz	1024×3	3.50M
	PointCNN [17]	xyz	1024×3	0.60M
	PointNet++ [25]	xyz + nor	1024×6	1.40M
	DGCNN [33]	xyz	1024×3	1.84M
	RS-CNN [20]	xyz	1024×3	1.41M
	Pt Transformer [43]	xyz	1024×3	- 1
	Pt Transformer v2 [34]	xyz	1024×3	-
	Spherical CNN [7]	voxel	2×64^2	0.50M
	RIConv [41]	xyz	1024×3	0.70M
ţ.	SPHNet [22]	xyz	1024×3	2.90M
Rotation-invariant	SFCNN [27]	xyz	1024×3	- '
	ClusterNet [3]	xyz	1024×3	1.40M
	GCAConv [40]	xyz	1024×3	0.41M
	RIF [16]	xyz	1024×3	
	RI-GCN [15]	xyz + nor	1024×6	4.38M
	RIConv++[42]	xyz	1024×3	0.40M
	RIConv++[42]	xyz + nor	1024×6	0.40M
	Ours (w/o normal)	xyz	1024×3	14.0M
	Ours (w/ normal)	xyz + nor	1024×6	14.0M



Table 6: Ablation study on the Self-Attention module.

Model SA1 SA2 Transformer Encoder Acc.							
A	1	~	~	96.0			
в		~	\checkmark	95.6			
\mathbf{C}	~		\checkmark	95.2			
D	~	~		94.3			
E				92.8			





Fig. 1: Modeling comparison. (a) Pixel-wise modeling [22] utilizes a predefined Graph2Pixel algorithm to rasterize the lane graph into a segmentation map and a direction map on dense BEV pixels, and heuristic Pixel2Graph post-processing is needed to recover the lane graph from the predicted segmentation map V_{pixel} and direction map D_{pixel} (direction map is not drawn here for simplicity). (b) Piece-wise modeling [6] utilizes a predefined Graph2Piece algorithm to split the lane graph into a set of pieces and the connectivity matrix among pieces, and then it merges the predicted pieces $\mathcal{V}_{\text{piece}}$ to the graph with the Piece2Graph algorithm based on predicted connectivity E_{piece} . (c) The proposed path-wise modeling translates the lane graph into complete paths with a predefined Graph2Path algorithm to recover the graph. We perform path detection and adopt a Path2Graph algorithm to recover the lane graph.





Fig. 1: Comparison between the prediction targets of previous work and S-JEPA (ours). Instead of raw 3D coordinates, S-JEPA predicts the abstract representations of 3D skeletons, embedded by a transformer encoder, effectively learning more informative high-level depth and context features for the action recognition task.



Fig. 2: Overview of S-JEPA. First, diverse skeleton views are obtained by applying geometric transformations on the 3D skeletons. The view skeletons are passed through the view encoder, after which learnable mask tokens are inserted at the locations of masked joints to get the view features \mathbf{F}_v . The predictor takes \mathbf{F}_v as input and outputs the predicted representations \mathbf{R}_p of the missing joints at the locations of the mask tokens. The target representations \mathbf{R}_t are obtained by the target encoder, which takes unmasked 3D skeletons as input, and is updated through the Exponential Moving Average (EMA) of the view encoder weights after each iteration (sg denotes stop gradient). The centering and softmax operations aid in stabilizing the training loss. At fine-tuning and test times, only the target encoder weights are used.








Fig. 2: Example of motion trajectory extraction of the point cloud sequence.

 Table 1: Performance comparison of action recognition with different methods of MSRAction-3D dataset.

	Algorithm	Accuracy (%)
	MeteorNet [20]	88.50
	PSTNet [9]	91.20
Supervised Learning	PSTNet++ [10]	92.68
	Kinet [45]	93.27
	PPTr [39]	92.33
	P4Transformer [7]	90.94
	PST-Transformer [8]	93.73
	PSTNet + PointCPSC [30]	92.68
	PSTNet + CPR [29]	93.03
End-to-end Fine-tuning	g PSTNet + PointCMP [28]	93.27
	P4Transformer + MaST-Pre [27]	91.29
	PST-Transformer + MaST-Pre [27]	94.08
	P4Transformer + M2PSC (ours)	93.03
	PST-Transformer + M2PSC (ours)	94.84









- ✓ "Sum" column displays dual tracer images (combined FDG and Amyloid [FBP and FMM]),
- ✓ "Ref FDG" represents the reference FDG
- ✓ "Gen FDG" denotes the generated FDG, and
- ✓ "FDG_Bias" signifies the difference map between reference and generated FDG
- ✓ "Ref Amy" represents reference Amyloid,
- ✓ "Gen Amy" refers to generated Amyloid, and
- ✓ "Amy_Bias" indicates difference map between reference and generated Amyloid
- The image range spans from 0 to 3 SUVR, whereas the difference map range is between -0.2 and +0.2 SUVR
- Subject-related metrics, including amyloid status, gender, MMSE, and age, are summarized atop each panel













Overview of few-shot RSIs recognition framework

- Start with dense sampling to generate discrete tokens of RSIs, to which 2D positional codes are assigned successively.
- These tokens are then clustered into distinct parts with guidance of von-Mises-Fisher (vMF) loss function, and
- The parts are combined into global remote sensing scenes through a Set-Transformer with the constraint of Predictive Info-NCE loss.
- Finally, the feature representations of global remote sensing scenes are
- classified with a prototype-based classification head.





Fig. 4. Hierarchical contextual prediction for compositional representation. It consists of bottom-up composition and top-down prediction. The optimization process is constrained by the Predictive Info-NCE Loss function.



- Each point represents a cluster centroid (i.e., part).
- Figure (a), each centroid is clustered closely, resulting in posterior representations of tokens that cannot be distinguished.
- Figure (b) depicts the scenario where each cluster centroid is separated. Here, posterior representations of the tokens can be treated as distinct parts

TABLE II Overall accuracy(%) on AID dataset

Method	5-way 1-shot	5-way 5-shot
MAML [5] MetaSGD [27]	$43.20 \pm 0.77\%$ $45.01 \pm 0.98\%$	$60.37 \pm 0.75\%$ $62.58 \pm 0.80\%$
LLSR [39]	45.18	61.76
RS-MetaNet [28]	33.87 $58.51 \pm 0.84\%$	50.40 $73.76 \pm 0.69\%$
DLA-MatchNet [41] our method	$61.99 \pm 0.94\%$ $65.48 \pm 0.68\%$	$\begin{array}{r} 75.03 \pm 0.67\% \\ \textbf{79.91} \pm \textbf{0.41}\% \end{array}$



- ✓ Multiply them in the assignment matrix from cosine similarities to reduce outliers.
- ✓ Finally, 9D object poses of novel instances are retrieved by Umeyama algorithm [27] with RANSAC [7] from the dense correspondences

Figure 10. t-SNE visualization of 3D semantic features from partial 3D features (red) and full 3D features (blue) inside the attention region before and after feature fusion.

- △ Identification methods: AR, MA and ARIMA [20];
- Statistical methods: Lasso Regression (LASSO) [23], Support Vector Regression (SVR) [41], Random Forest (RF) [21], and eXtreme Gradient Boosting (XGB) [22];
- Deep methods: Long Short-Term Memory (LSTM) [26], Gated Recurrent Unit (GRU) [25], Transformer [28], Informer [42], AttentionMixer [7] and iTransformer [43].
- ✓ GRU decoder for transformer-based baselines performed to produce quality variable prediction

Methods	H=1			
	MAE	R2		
Dataset: Hegan	g Station			
AR	0.135 ± 0.000	0.064 ± 0.000		
MA	0.133 ± 0.000	0.110 ± 0.000		
ARIMA	0.135 ± 0.000	0.063 ± 0.000		
LASSO	0.045 ± 0.000	$0.282_{\pm 0.000}$		
SVR	$0.023_{\pm 0.000}$	$0.890_{\pm 0.000}$		
XGBoost	0.016 ± 0.000	0.943 ± 0.000		
LSTM	$0.021_{\pm 0.001}$	$0.846_{\pm 0.011}$		
GRU	$0.023_{\pm 0.001}$	0.832 ± 0.008		
Transformer	0.015 ± 0.002	$0.927_{\pm 0.022}$		
Informer	$0.042_{\pm 0.017}$	0.562 ± 0.269		
Transformer	$0.014_{\pm 0.001}$	$0.919_{\pm 0.026}$		
AttentionMixer	$0.011_{\pm 0.003}$	0.962 ± 0.012		
DeepFilter	0.012 ± 0.001	0.963 ± 0.006		
Dataset: Jinan	Station	_		
AR	0.106 ± 0.000	0.744 ± 0.000		
MA	0.112 ± 0.000	0.710 ± 0.000		
ARIMA	0.101 ± 0.000	0.759 ± 0.000		
LASSO	0.029 ± 0.000	0.742 ± 0.000		
SVR	0.061 ± 0.000	0.640 ± 0.000		
XGBoost	0.031 ± 0.000	0.903 ± 0.000		
LSTM	0.015 ± 0.000	0.948 ± 0.001		
GRU	0.016 ± 0.001	0.939 ± 0.005		
Transformer	0.016 ± 0.004	0.955 ± 0.023		
Informer	0.018 ± 0.002	$0.917_{\pm 0.013}$		
Transformer	0.012 ± 0.000	0.971 ± 0.005		
AttentionMixer	0.011 ± 0.004	0.981 ± 0.012		
DeepFilter	$0.010_{\pm 0.002}$	0.986±0.003		
mparative Study C)n The Hegang A	nd Jinan Datasets		

Fig. 4: Glimpse selection step-by-step: AdaGlimpse explores 224×224 images from ImageNet with 32×32 glimpses of variable scale, zooming in on objects of interest and stopping the process after reaching 75% predicted probability. The rows correspond to: A) glimpse locations, B) pixels visible to the model (interpolated from glimpses for preview), C) predicted label, D) prediction probability.

Model and training conf	liguration
-------------------------	------------

✓ Model and training configuration details specified

Parameter name	Value
Encoder	
ViT type	base
native patch size	16
transformer embed dim	768
transformer blocks	12
attention heads	12
mlp ratio	4
Decoder	
transformer embed dim	512
transformer blocks	8
attention heads	16
mlp ratio	4
RL agent	
hidden dim	256
action distribution	TanhNorma
sac target entropy value	-3
sac initial alpha value	1

training epochs	100
backbone pre-training epochs	600
backbone lr	10^{-5}
rl agent lr	$5 * 10^{-4}$
backbone weight decay rate	10^{-4}
rl agent weight decay rate	10^{-2}
lr scheduler type	one cycle
lr warmup epochs	10
minimum lr	10^{-8}
initial random action batches	10000
initial frozen backbone epochs	10
rl loss function	L2
rl batch size	256
backbone batch size	128
replay buffer size	10000

in Snangnal i ech dataset							
Mathad	Part	A	Part B				
Method	MAE	MSE	MAE	MAE			
MCNN[10]	110.2	173.2	26.4	41.3			
CSRNet[14]	68.2	115.0 107.8	10.6	16.0 14.0			
DSPNet[37]	68.2		8.9				
C-CNN[37]	88.1	141.7	14.9	22.1			
DNCL[37]	73.5	112.3	18.7	26.0			
ED-CNN[40]	69.8	114.7	10.2	14.9			
ICC[42]	76.9	130.1	8.4	15.2			
TransCrowd-token[43]	69.0	116.5	10.6	19.7			
TransCrowd-GAP[43]	66.2	105.1	9.3	16.1			
AutoScale_loc[44]	65.8	112.1	8.6	6 13.9			
MACC+SM[44]	67.7	113.0	9.8	12.9			
MSC-FFN[32]	65.8	105.9	7.6	11.8			
CCD Net[46]	70.0	118.3	-	-			
DA ² Net[33]	74.1	128.4	7.9 1				
T ² CNN[47]	85.3	137.4	18.8	29.2			
CSFNet(ours)	66.1	103.2	7.5	11.8			

Table 2: Comparison of CSFNet with other advanced counting methods in ShanghaiTech dataset

Figure 1: Demonstrating rotation-invariant face detection in complex, real-life scenarios with inevitable in-plane face rotations captured during gatherings, sports, and artistic performances.

RP-Net: A Robust Polar Transformation Network for Rotation-Invariant Face Detection

- ✓ During the test, we only retain baseline model and MGFML module after Stage 4 to achieve the cross-modality retrieval
- DMANet: Dual-Modality Alignment Network for Visible-Infrared Person Re-Identification

MoE S	Stacked .	Att-Shared	MoE-Shared	Params	FLOPs]	-AUROC	P-AUROO
				4.6M	2.16G	96.1	96.3
		\checkmark		2.8M	2.16G	95.7	96.1
\checkmark				$10.9 \mathrm{M}$	2.17G	97.4	96.9
\checkmark		\checkmark		9.1M	2.17G	97.2	96.8
\checkmark			\checkmark	4.6M	2.17G	97.1	96.8
\checkmark		\checkmark	\checkmark	2.8M	2.17G	97.1	96.8
\checkmark	\checkmark			19.3M	2.18G	98.1	97.2
~	\checkmark	\checkmark		17.5M	2.18G	97.9	97.1
\checkmark	\checkmark		\checkmark	6.7M	2.18G	97.7	97.0
~	\checkmark	\checkmark	\checkmark	4.9M	2.18G	97.7	97.0

Fig. 1: Comparison of Temporal Fusion Methods. Illustrated are four key approaches: (1) Temporal Self-Attention [20], leveraging attention mechanisms for temporal integration; (2) Warp and Concat [8,35,37], combining features across frames and fusing them through convolution; (3) Cost Volume Construction in image space [25], constructing cost volume from image input of different frames and leveraging plane-sweep volumes for depth map generation; and (4) Our Proposed Method, which involves constructing a temporal cost volume in 3D space to enhance feature refinement. In the figure, (\triangle) and \otimes represent coordinate alignment and element-wise product, accordingly.

Fig. 2: Overall Architecture of CVT-Occ. The image backbone extracts multiscale features from multi-view images, which are transformed into 3D volume features denoted as $\mathbf{V} \in \mathbb{R}^{H \times W \times Z \times C}$. The Cost Volume Temporal Module samples points along the line of sight within the current volume and projects them onto K - 1 historical frames, resulting in $K \times N$ 3D volume features. These features are concatenated to construct cost volume features $\mathbf{F} \in \mathbb{R}^{H \times W \times Z \times (K \times N) \times C}$. Convolution layers are then applied to generate weights $\mathbf{W} \in \mathbb{R}^{H \times W \times Z}$, refining the depth of 3D voxel. Finally, an occupancy decoder produces 3D semantic occupancy predictions. In the figure, (a), (c), and \otimes represent coordinate alignment, concatenation, and element-wise product, respectively.

(a) The single disperser CASSI imaging process. HSI data cube is captured by a monochromatic sensor.

(b) GC-GAP projection.

- (c) Latent encoder.
- (d) Cincellified Dee

(d) Simplified Denoiser.

(e) The measurement y and masks A pass through an N-stage DUN, where

each stage is composed of a GC-GAP projection and a denoiser.

The denoiser follows a U-shape structure and consists of five Trident Transformers (TT), where each

TT is assisted with prior knowledge zGT generated from the diffusion model

Season	Model	MAE	MSE	RMSE	R^2	PICP	PINAW	WS
	IVMD-BO-QRB	21.12	1079.12	32.85	0.73	56.18%	15.51%	47.47%
	IVMD-BO-QRTT	21.22	1022.08	31.97	0.74	31.46%	46.88%	16.71%
	IVMD-BO-QRTLT	21.67	1089.00	33.00	0.74	51.69%	31.15%	35.59%
	IVMD-BO-QRTBT	21.79	1112.22	33.35	0.74	83.15%	17.99%	68.19%
Spring	IVMD-MOBO-QRB	20.69	1090.98	33.03	0.74	92.14%	12.94%	80.22%
	IVMD-MOBO-QRTT	21.53	1011.24	31.80	0.74	96.63%	18.21%	79.03%
	IVMD-MOBO-QRTLT	19.72	996.03	31.56	0.76	91.91%	14.56%	78.53%
	IVMD-MOBO-QRTBT	19.98	1060.80	32.57	0.76	96.40%	15.11%	81.83%







