Available online at www.joac.info

ISSN: 2278-1862



Journal of Applicable Chemistry

2025, 14 (2): 279-394 (International Peer Reviewed Journal)





CNN - 66B

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

Information Source	sciencedirect.com;	
S. Narasinga Rao M D	K. Somasekhara Rao, Ph 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,	Dr. M.R.Appa Rao Campus,	Visakhapatnam 530 003,
King George Hospital	Nuzvid-521 201, India	India
Visakhapatnam, A.F., 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: A Transformer net (TransF Net) or Transformer neural network (TransF NN) consists of an attention layer and MLP-NN to carry out Natural Language processing (NLP). The seminal paper of Ashish Vaswani et al. entitled "Attention is All You Need", in the year 2017 revolutionised sequence text data processing. The evolution of architecture of TransF NN, attention mechanism, and hybridization with other methods, during these few years, brought this paradigm to the fore-front of Data

Science dealing with Text, numerical time-series, sound (speech), images and video sequence with local and global inter-dependencies.

Data: Data are many types, viz., Boolean (AND (Conjunction), OR (Disjunction), NOT(Negation) ...), numerical (integer, floating point), sound (speech), text, image, video, and that from tactile senses (touch, smell, taste etc.). Real and complex (and quaternion, octonion, ...) belong to real and imaginary spaces. Depending upon bases of items, tensors (zero, first, second, third order or scalar, vector, matrix, tensor-of-third-order) are popular. Higher (>3) order tensors are also sometimes called as multi (4-, 5-, 6-) –way-data. Data are also referred as binary, octal, decimal and hexa-decimal with base 2,8,10 and 16 respectively. In images pixels for two-dimensional and voxels for 3-dimensional spaces are standard with black/white, RGB, and grey-colors. An image is also conceived as a culmination of patches.

Text data contain documents, pages, paragraphs, sentences, words and alphabetic characters. Each datum for example, a word in text, amino-acid in protein, real numerical value in stock market/forex field (time series), or patch in an image is called a token. The position of token in the sequence and numerical (embedded) value are the primary input for analysis. A few typical tasks in NLP are translation of text in one language (English) to another (German or French or ...), summarizing documents, generation of new stories in text format and so on.

The architectures of a few of Transformer neural nets (TransF-NN) documented in this state-of-knowledge-methods-module for Data2Knowledge transformation are

- A Multi-scale Acceleration Feature Fusion (MAFF) Transformer (TMAFF)
- Solution Video second-order transformer network (ViSoT)
- DHCT-GAN Transformer (DHCT: Dual-branch Hybrid CNN–Transformer Network, GAN: generative adversarial network)
- Discrete Cosine Transform Network (DCTNet)
- B Hierarchical Cross-Modal Similarity Search and Transfer (HCMSST)
- Enhanced Hybrid of CNN and Transformer Network (EHCTNet)
- Swin Transformer (SwinTrnF)
- Spatial-Temporal Transformer Network (ST-TNet)
- A Hybrid Attention-Dense Connected Transformer (HADT) Networks
- Efficient Conv-Transformer (ECTFormer)
- A Model architecture of MultiscAle Relational Transformer network (MART)
- A Multifractal Spectrum Transformer (MFSTransFormer)
- Residual sparse Transformer branch (RSTB) ;{[MSFN Mixed scale FFNN Top C sparse attention (TCSA)]}
- Parallel pyramidal Transformer [MSD (Multi-Scale Decoder), PPT, (Parallel Processing Transformer) FM. (Feature Modulation)]
- SCACD-Net [SCA (Spatial and Channel Attention) CD (Contextual Decoder)]
- Digital image correlation (DIC) transformer (DICTr)
- A TC (tea chrysanthemums) ViT
- SF-Hformer (Spatial-Frequency Hybrid Transformer)
- Multi-Patch De-raindrop Transformer for UAV images (MPDT) [contains frequency attention Transformer block (FATB) and Adaptive Feature Enhancement Module (AFEM)]

A	Swin Transformer-CNN Gait (STCG) Recognition
A	Histoformer [Dynamic-range Histogram Self-Attention (DHSA) module and the Dual-scale
•	Gated Feed-Forward]
	DTCN-EFFN-Transformer [DTCN: Deep time (Temporal) convolution network structure diagram; EFFN (Feed-Forward Networks based on Expansion factor) adds an expansion factor to the linear layer]
A	MDTNet (Multi-Domain Transformer Network)
A	Enhanced Hierarchical Vision Transformer (EHVT); local geometric transformer; global semantic transformer
A	Progressive alignment and interwoven composition network (PAIC-Net)
A	Denoising Adaptive Graph Transformer (DAGT)
A	Conformer (Convolution-augmented and confidence estimation network Transformer)
A	MEMAFormer
A	SAPPM
A	Multi-Scale Convolutional Prediction Network; Multi-Scale Convolutional PoolFormer block (MSCPNet)
A	MUltifuSion (MUST) Transformer
A	Enhanced Downsampling Attention Block (EDAB)
A	Frequency Prompting image restoration (FPro) Method
- e mae e nave e nave e nave e nave e nave e nave 11 i 1881 i 1887	
Keyw	ords: Artificial intelligence (AI); Capsule Neural Nets— MLP-
Attentio	on Mechanism-TransFormer Nets—Hybrid TransFormer Networks K(nowledge)Lab

CNN : [C [Computations; Computer; Chemistry, Cell, Cellestial, Cerebrum] NN [New News; News New; Neural Nets; Nature News; News of Nature;]] Fits : [Figure Image Table Script;]

K(nowledge)Lab rsr.chem1979







Woi	rkflow of SDI	D (Structural damage detection) based on TMAFF
	(a) (a) (b) (c)	tion and derivative. (a) ReLU; (b) ELU; (c) GELU
Nama	Table 4 C	Configurations of training process
IName	value	Description
Datah siza	250	The size of the date batch used in training
Batch size	256	The size of the data batch used in training
Batch size Initial learning rate	256 5×10 ⁻⁴	The size of the data batch used in training The initial learning rate for the optimization algorithm
Batch size Initial learning rate L2 rate	256 5×10 ⁻⁴ 1×10 ⁻³	The size of the data batch used in training The initial learning rate for the optimization algorithm The weight of L2 regularization
Batch size Initial learning rate L2 rate Random seed	256 5×10 ⁻⁴ 1×10 ⁻³ 42	The size of the data batch used in training The initial learning rate for the optimization algorithm The weight of L2 regularization The seed value used for random number generation
Batch size Initial learning rate L2 rate Random seed Patience	256 5×10 ⁻⁴ 1×10 ⁻³ 42 30	The size of the data batch used in training The initial learning rate for the optimization algorithm The weight of L2 regularization The seed value used for random number generation A parameter of early stopping, the acceptable duration if the performance does not increase



	Stand haad haad haad haad haad haad haad h	la ka
Transformer Net	2025-01	
Mari sari sari sari sari sari sari sari s		







(b) illustrates such operation from a 3D perspective of one video clip. White rectangle or cube refers to padding zeros.









Transformer Net	2025-03	
	2020 00	





n en	1 1	
Transformer Net	2025-03	



corresponds to the specific target category of prediction y_c . The dashed lines in the diagram represent corresponding relationships.





\checkmark The texts above and below the input images indicate the target classes of the visualizations





17/16
ar de

¥



		na na naina naina na naina na naina na naina na naina naina na naina na naina na naina na naina ata
Transformer Net	2025-05	





🟹 । नेवित निवित		, 1997 1997 1997 1997 1997 1997 1997 199
Transformer Net	2025-05	
	2023-03	



Transformer Net	2025-05	
	87 1887 1887 1887 1887 1887 1887 1887 1	9 48 48 48 48 48 48 48 4

Stage	Hyperparameter	Value
Freeze	Freeze epoch	50
	Freeze batch size	6
UnFreeze	UnFreeze epoch	300
	UnFreeze batch size	4
General	Initial learning rate	1e-3
	Optimizer	SGDM
Image	Size	368KB
	Resolution	512×512

Table 2 Experimental platform setup

Items	Parameters
CPU	Intel(R) Core(TM) i9-12900K
GPU	NVIDIA GeForce RTX 4080
Operating system	Windows 11
Acceleration environment	CUDA 12. 1
Development Platform	Pytorch 1.7.0





2025-06



Transformer Net	2025-07	
	2023-07	





(a): Bi-temporal input images; (b): Raw multi-scale feature images; (c): First-order feature images; (d): Semantic pixel maps; (e): Semantic difference map; (f): Second-order semantic difference map;

(g): Change detection heatmap





Network architectures for pretext tasks, face detection, and keypoint detection

- The encoder is connected to fully connected (#1) for classification and parallelly to the decoder for image inpainting and is trained on these pretext tasks.
- The pre-trained encoder is employed for face detection using fully connected (#2), and for facial keypoint detection, encoder is utilized alongside fully connected (#3) and Swin transformer block





















Figure 4: Swin base leaf cyclic shifting

Type of Layer	Output Vector	Configurations
Input Layer	(None, 224, 224, 3)	Image
VGG Block	(None 56 56 199)	Conv2D (3 x 3) (stride-
	(None, 56, 56, 128)	ing2D (2x2) (stride-2)
Inception Block		Conv2D (diff-diff kernel
	(None, 56, 56, 512)	size and stride), BN, ReLU, MaxPooling2D (2x2) (stride-1)
Stage 1		Conv2D (4 x 4) (stride-
	(None, 14, 14, 512)	tion, W-MSA, MLP, Proposedshifting-MSA
Stage 2		LayerNormalisation,
	(None, 7, 7, 512)	W-MSA, MLP, Proposedshifting-MSA
Classification	(None, 6)	LayerNormalisation, Avg-

Table 1: Layered description of the proposed model



Figure 5: Proposed leaf random shifting

Table 2: Soybean Real Field Dataset

Subject	Information
Type of data	Raw data
How data were acquired	Data curation performed with the help of local farmers, ex-
	perts, and pathologists using different devices. The Sony Cyber-
	shot DSC-300 camera, and smartphones, including the Samsung
	Galaxy M31 equipped with a 64-megapixel camera and an Ap-
	ple iPhone with a 12-megapixel camera.
Data format	The images are in JPEG format with a standard size of 224×224
	pixels.
Description for data col-	The images were captured around 11:30 a.m. and 3:30 p.m.
lection	during ideal illumination and mild temperatures for data cura-
	tion. captured primarily concentrate on the upper regions of
	soybean leaves affected by insect infestation, as well as leaves
	that exhibit normal, healthy conditions.
Data source location	Season 1 is collected from Vidisha district, whereas season 2 is
	from the guna district of Madhya Pradesh, India.


Table 4: Accuracy comparison on soybean dataset

Classification Models	Accuracy	Precision	Recall	F-1 score	Parameters in Milions
Ghost-convolution enlightened Transformer [32]	0.67	0.68	0.67	0.66	25.39
PlantXVIT [33]	0.78	0.80	0.78	0.77	6.40
Convolutional Swin Trans- former [27]	0.89	0.90	0.88	0.89	27.47
Inception convolutional trans- formers (ICVT) [34]	0.92	0.96	0.95	0.95	11.16
Former-Leaf [35]	0.74	0.76	0.74	0.74	60.00
Con-Vit [23]	0.73	0.75	0.73	0.75	38.00
RIC-Net [3]	0.88	0.89	0.87	0.88	6.71
MobileNet-V2 [36]	0.80	0.85	0.76	0.80	5.32
MFSwin Trans [28]	0.74	0.75	0.74	0.74	12.00
Proposed	0.94	0.93	0.94	0.92	5.20

Transformer Net

2025-12

Algorithm 1 :Proposed Model

Algorithm 1 .1 Toposed Model
1: Input: Dataset
2: Output: Optimal Classification
3: Set N_{epochs} as Number of Epochs
4: for epoch = $1, 2, \dots, N_{epochs}$ do
5: Set previous state \leftarrow image pixels position
6: for batch = $1, 2,, T$ do
7: Collect $x \in P(x)$ size(224, 224, 3) from dataset to train the model
8: current state \leftarrow previous state
9: window size $\leftarrow 7x7$
10: array b have (224,224) position \leftarrow 50176
11: $n_windows \leftarrow 1024$
12: for window $=1,2,,n$ do
13: $array_c \leftarrow elements of current window$
14: $array_b \leftarrow array_c$
15: $array_d \leftarrow random 49(7X7)$ elements from remaining array_b to
16: $array_b \leftarrow array_b - array_d$
17: $array_b \leftarrow array_b + array_c$
18: end for
19: end for
20: end for
21: Save model after training

Type of Layer	Output Vector	Configurations		
Input Layer	(None, 224, 224, 3)	Image		
VGG Block		Conv2D (3 x 3) (stride-		
	(None, 56, 56, 128)	1), BN, ReLU, MaxPool-		
		ing2D (2x2) (stride-2)		
Inception Block				
		Conv2D (diff-diff kernel		
	(None 56 56 512)	size and stride), BN,		
	(1016, 50, 50, 512)	ReLU, MaxPooling2D		
		(2x2) (stride-1)		
Stage 1				
0		Conv2D (4×4) (stride-		
	(None 14 14 519)	4), LayerNormalisa-		
	(1000, 14, 14, 512)	tion, W-MSA, MLP,		
	2	Proposedshifting-MSA		
Stage 2		/)		
		LayerNormalisation,		
	(None, 7, 7, 512)	W-MSA, MLP,		
	\sim	Proposedshifting-MSA		
Classification				
	(None, 6)	LayerNormalisation, Avg-		
		Pooling, Linear		

Table 1: Layered description of the proposed model







- ✓ The input tensor is split to query Q and key K in half along the channel dimensions, which makes prevent an increase in the size of the model parameters.
- ✓ Due to the single-head mechanism, it can effectively alleviate the bottleneck of the attention operation, specifically the reshape operation, resulting in a reduction in the model's execution time

Models	Params (M)	FLOPs (G)	Size	Top-1 (%)	Throughput	Latency (ms)	Acc./ Latency
MobileNetv2-×1.4 [9]	6.4	0.6	224 ²	74.7	4589.1	23.3	3.2
EfficientNet-B0 [32]	5.3	0.4	224 ²	77.1	5845.3	22.3	3.5
ConvNeXt	4.1	0.6	224 ²	74.9	5586.4	20.7	3.6
MobileViT-S [17]	5.6	1.4	256 ²	78.4	2613.2	55.8	1.4
EdgeNeXt-S [13]	5.6	1.3	256 ²	79.4	3263.7	47.6	1.7
EdgeNeXt-Sb	5.6	1.0	224 ²	77.3	4207.0	40.6	1.9
EdgeViT-XS [24]	6.8	1.1	224 ²	77.5	4551.7	33.8	2.3
ECTFormer-×1.0	4.4	0.7	224 ²	77.1	6553.5	19.7	3.9
ECTFormer-×1.25	6.7	1.1	224 ²	79.0	4919.8	25.6	3.1
ECTFormer-×1.5	9.6	1.6	224 ²	80.3	4244.9	30.0	2.7

^a Refers to ConvNeXt-Tiny with reduced channel dimensions and depths for each stage, as described in Section 3.4.

^b Denotes the performance trained with the input image sizes of 224 × 224 without multiscale sampler [17] for a fair comparison.



Visualization of the activation maps using Grad-CAM between ECTFormer and the other efficient models

Table 2

Performance comparison on various fine-grained datasets.

Datasets	# Classes	Models	Params (M)	Top-1 (%)	Top-5 (%)	Acc./ Latency
Aircraft [46] 100		ConvNeXt [†]	3.8	86.5	97.4	4.2
		EfficientNet-B0 [32]	4.1	85.8	96.6	3.8
	100	EdgeNeXt-S [‡]	5.3	87.7	96.8	2.2
	100	EdgeViT-XS [24]	6.4	88.0	97.5	2.6
		ECTFormer-×1.0	4.1	87.6	97.4	4.4
		ECTFormer-×1.25	6.4	87.9	97.1	3.4
Flowers [47] 102		ConvNeXt [†]	3.8	97.6	99.6	4.7
		EfficientNet-B0 [32]	4.1	98.4	99.7	4.4
	100	EdgeNeXt-S [‡]	5.3	97.9	99.8	2.4
	102	EdgeViT-XS [24]	6.4	97.9	99.7	2.9
		ECTFormer-×1.0	4.1	98.3	99.8	5.0
		ECTFormer-×1.25	6.4	98.6	99.7	3.9
	ConvNeXt [†]	3.8	87.8	96.9	4.2	
		EfficientNet-B0 [32]	4.1	89.0	97.5	4.0
Food [48] 10	1.01	EdgeNeXt-S [‡]	5.3	89.4	97.5	2.2
	101	EdgeViT-XS [24]	6.4	90.1	97.8	2.7
		ECTFormer-×1.0	4.1	89.5	97.5	4.5
		ECTFormer-×1.25	6.4	89.5	97.5	3.5

NARTEREN KARANTEREN KARANTEREN KARANTEREN KARANTEREN KARANTEREN KARANTEREN KARANTEREN KARANTEREN KARANTEREN KAR

StanfordCars [49]		ConvNeXt [†]	3.9	92.2	99.1	4.5
	196	EfficientNet-B0 [32]	4.3	91.9	98.8	4.1
		EdgeNeXt-S [‡]	5.3	92.4	99.0	2.3
		EdgeViT-XS [24]	6.5	93.1	99.1	2.8
		ECTFormer-×1.0	4.1	92.4	99.0	4.7
		ECTFormer-×1.25	6.4	93.1	99.0	3.6
Dec (50)	37	ConvNeXt [†]	3.8	91.6	98.7	4.1
		EfficientNet-B0 [32]	4.1	91.8	98.7	3.8
		EdgeNeXt-S [‡]	5.3	93.1	99.0	2.1
Pet [50]		EdgeViT-XS [24]	6.4	92.2	99.1	2.5
		ECTFormer-×1.0	4.1	92.5	99.0	4.4
		ECTFormer-×1.25	6.3	93.2	99.2	3.4
		ConvNeXt [†]	3.9	84.5	93.9	3.6
	256	EfficientNet-B0 [32]	4.3	85.4	94.4	3.5
Calcash DEC (E1)		EdgeNeXt-S [‡]	5.4	86.7	94.8	1.9
Catteen-256 [51]		EdgeViT-XS [24]	6.5	84.7	94.1	2.3
		ECTFormer-×1.0	4.2	86.4	94.9	3.9
		ECTFormer-×1.25	6.4	87.5	95.1	3.1





Fig. 7: (a) Visualization of attention weights in PRT. (b) Visualization of attention weights in HRT (left), the group incidence matrix G (middle), and the corresponding groups (right). The sample results from the identical test data sample used in Figure 6. For the attention visualization, min-max normalization is applied to ensure a clearer visual understanding, and darker colors indicate higher attention weights.



<u>1</u> 1 - 1867 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1887 - 1		7 100 1
Transformer Net	2025_16	
	2023-10	
MARIN MARIN MARIN MARIN MARINA NA ING		

































neighboring nodes



























Datasets		Rain200L		Rain200H	
Metrics		PSNR	SSIM	PSNR	SSIM
Duing based worth a da	DSC [1]	27.16	0.8663	14.73	0.3815
Prior-based methods	GMM [2]	28.66	0.8652	14.50	0.4164
CNN-based methods	PReNet [29]	37.80	0.9814	29.04	0.8991
	RCDNet [3]	39.17	0.9885	30.24	0.9048
	MPRNet [5]	39.47	0.9825	30.67	0.9110
	DualGCN [30]	40.73	0.9886	31.15	0.9125
	SPDNet [4]	40.50	0.9875	31.28	0.9207
	Uformer [31]	40.20	0.9860	30.80	0.9105
	Restormer [13]	40.99	0.9890	32.00	0.9329
Transformer-based methods	IDT [15]	40.74	0.9884	32.10	0.9344
	DRSformer [14]	41.23	0.9894	32.18	0.9330
	Ours	41.54	0.9898	32.52	0.9349







Fig. 2: An overview of the proposed Dual-Path Adversarial Lifting method. During inference in the target domain, the domain shift token P_l , the prediction network Φ_l , and the update network Ψ_l are updated given each mini-batch testing samples. The dual-path lifting transformer (Left). The details of the lifting block in each layer (Right).



to learn domain-specific knowledge for different domains across layers.




Fig. 1: Illustration of the *Think Before Placement* process. While other methods directly determine object placement based on image features, our approach adds a thinking process to this procedure. First, a large multi-modal model assesses which positions are reasonable, then it proceeds with the placement, achieving a more rational positioning of objects. I_{fg} , M_{fg} , I_{bg} and I_{comp} stands for foreground image, foreground mask, background image and composite image respectively.







✓ (a) Illustration of floorplan representation including the sparse corner-based (left) and our dense uniform sampling representation (right). Valid contour vertices (outlined) and corner vertices (filled) with supervision during training are colored according to the room on the left.
✓ (b) Illustration of the self-attention variants including the room-aware self-attention and vanilla self-attention. Our room-aware self-attention is a combination of intra-room and inter-room self-attention, which works among different vertices in a single room and among different rooms. And the vanilla self-attention performs on the flattened queries.

















Figure 5: Our proposed loss function for penalizing the misalignment between two lines. If the scribble line touches the top or the bottom patch boundary, then we use interpolation to calculate the loss as shown in the figure on the right.



Figure 7: First, candidate patches of various sizes $m = \{64, 128, \dots 512\}$ are sampled from the document. Line-parameter predictions are made for each patch by resizing them to the input shape expected by the model (s = 256)and then passing them through the Line-Parameter Generator (Sec. 3.2). An interline-gap value δ for each patch is calculated using the line predictions for the patch. These patch-level interline gap predictions are then scaled to the size of the original document ($\delta := \delta \cdot \frac{m}{s}$) and averaged to obtain an estimate of the average interline gap value for the whole document. This averaged value is multipled by a scaling factor ζ , which roughly represents the expected number of text lines seen in a patch to obtain the final context-adapted patch size (t). Patches of size t are then sampled from the original document, and are fed to the Line-Parameter Generator to get the final predictions.





















- ✓ We visualize the refined
- ✓ similarity on an airplane, with a red query point located at the wing.
- ✓ In the refined similarity generated by our method, high response weights are exclusively assigned to points belonging to the wing section

Table 2. Part segmentation results on ShapeNetPart dataset.

Methods	Ins. mIoU	Cat. mIoU	air.	bag	cap	car	cha.	ear.	gui.	kni.	lam.	lap.	mot.	mug	pis.	roc.	ska.	tab.
PointAttn. [5]	85.9	84.1	83.3	86.1	85.7	80.3	90.5	82.7	91.5	88.1	85.5	95.9	77.9	95.1	84.0	64.3	77.6	82.8
PointASNL[26]	86.1	83.4	84.1	84.7	87.9	79.7	92.2	73.7	91.0	87.2	84.2	95.8	74.4	95.2	81.0	63.0	76.3	83.2
PointMLP [12]	86.1	84.6	83.5	83.4	87.5	80.5	90.3	78.2	92.2	88.1	82.6	96.2	77.5	95.8	85.4	64.6	83.3	84.3
APES [25]	85.8	83.9	85.3	85.8	88.1	81.2	90.6	74.0	90.4	88.7	85.1	95.8	76.1	94.2	83.1	61.1	79.3	84.2
PointTran. [30]	86.5	83.7	85.8	85.3	86.8	77.2	90.5	82.0	90.8	87.5	85.2	96.3	75.4	93.5	83.8	59.7	77.5	82.5
PointGT [9]	85.8	84.2	84.3	84.5	88.3	80.9	91.4	78.1	92.1	88.5	85.3	95.9	77.1	95.1	84.7	63.3	75.6	81.4
PCT [7]	86.4	83.1	85.0	82.4	89.0	81.2	91.9	71.5	91.3	88.1	86.3	95.8	64.6	95.8	83.6	62.2	77.6	83.7
LGGCM [3]	86.7	84.8	85.1	85.9	90.3	80.8	91.6	75.4	92.7	88.1	86.5	96.1	77.0	94.2	84.5	63.6	80.2	84.3
Ours	87.5	85.6	86.1	87.2	88.1	79.4	92.4	82.3	92.0	88.4	86.9	96.7	78.7	95.6	85.8	63.8	79.3	85.3









- The PointNet-based local 3D encoder and 2D auto-encoder extracts local 3D features and local 2D features, respectively.
- Then, the 2D features are projected into 3D space to fuse with 3D features forming the super point features F.
- Based on the super points, keypoint embeddings E and 3D patches are extracted as input to the novel adaptive graph transformer to estimate accurate 3D keypoint coordinates by leveraging the dynamic kinematic correspondences and local details.
- Notably, during the training stage, the local 3D patches are shifted by random noises in order to enforce the model providing robust estimations

24.12.1	Mean keypoir	-	
Method	ICVL	NYU	Input
Ren-9x6x6 [17]	7.31	12.69	D
DeepPrior++ [32]	8.1	12.24	D
Pose-Ren [4]	6.79	11.81	D
DenseReg [44]	7.3	10.2	D
CrossInfoNet [10]	6.73	10.08	D
JGR-P20 [11]	6.02	8.29	D
SSRN [37]	6.01	7.37	D
PHG [36]	5.97	7.39	D
HandPointNet [13]	6.94	10.54	Р
Hand-Transformer [21]	6.47	9.80	P
Point-to-Point [16]	6.3	9.10	P
V2V [31]	6.28	8.42	V
HandFolding [8]	5.95	8.58	P
HandR2N2 [6]	5.70	7.27	P
IPNet [35]	5.76	7.17	D+P
HandDAGT (Ours)	5.66	7.12	D+P

Table 1: Comparison of the proposed method with previous state-of-the-art methods on the single-hand ICVL and NYU datasets. Input indicates the input type of 2D depth image (D), 3D voxels (V), or 3D point cloud (P).





Table 1.Comparison of results of different methods on randomised TSP dataset

Mashad		TSP20			TSP50			TSP100	
Method	Len.	Gap	Time	Len.	Gap	Time	Len.	Gap	Time
Concorde ^[2]	3.83	0.00%	2.3m	5.69	0.00%	14m	7.76	0.00%	1.1h
LKH-3 ^[3]	3.83	0.00%	21.2m	5.69	0.00%	27.3m	7.76	0.00%	55m
Bello et-al.*[11]	3.89	1.42%	-	5.95	4.46%	-	8.30	6.90%	-
Deudon et-al.*[12]	3.84	0.26%	-	5.81	2.07%	-	8.85	13.97%	-
Kool et-al. (greedy)*[15]	3.85	0.34%	Os	5.80	1.76%	2s	8.12	4.53%	6s
Bresson et-al. (greedy) ^[24]	3.89	1.57%	Os	5.75	1.05%	14s	8.01	3.22%	19s
Jung et-al. (greedy) ^[18]	3.84	0.25%	Os	5.75	0.98%	6s	8.00	3.00%	12s
Ours (greedy)	3.84	0.22%	Os	5.74	0.95%	1s	7.95	2.45%	6s
Kool et-al. (B=1280)*[15]	3.84	0.08%	5m	5.73	0.52%	24m	7.94	2.26%	1h
Bresson et-al. (B=2500)	3.85	0.34%	14m	5.75	0.97%	44.8m	7.86	1.26%	1.5h
Jung et-al. (B=2500) ^[18]	3.83	0.00%	1.4m	5.72	0.46%	26.2m	7.86	1.22%	1.83h
Ours (B=2500)	3.83	0.00%	1m	5.70	0.14%	20.4m	7.80	0.76%	1h



The AUC value of ground truth confidence is measured as 'Optimal'. The result with the lowestAUC value in each experiment is highlighted.

Methods	Flying	Things 25	Driving 25		KITTI	2012 26	KITTI 2015 27		
	PSM	GA-Net	PSM	GA-Net	PSM	GA-Net	PSM	GA-Net	
CCNN-S 9	9.783	3.598	13.150	5.047	0.079	0.906	0.015	0.039	
ConfNet-S 2	10.060	3.719	13.502	4.898	0.040	1.078	0.012	0.062	
LGC-Net-S 2	10.272	4.185	12.822	5.313	0.041	0.818	0.017	0.039	
LAF-Net-S	8.372	3.350	12.260	4.358	0.032	0.575	0.014	0.035	
ConFormer-S	7.726	2.963	10.784	3.168	0.029	0.339	0.011	0.023	
Optimal	4.502	1.602	4.040	1.399	0.002	0.086	0.0002	0.0004	







Method	Backbone	Parameters (M)	mIou			
			ADE20K	Cityscapes		
SETR-PUP (80 k) [25]	ViT-B	98	46.3	-		
	ViT-L	318	48.6	82.15		
Swin transformer [32]	Swin-T	60	46.1	-		
	Swin-S	81	49.3	-		
PVT [50]	PVT v2-B0	7.6	37.2	-		
	PVT v2-B1	18	42.6	-		
Segmenter [60]	Seg-B	86	48.1	80.5		
	Seg-B/Mask	86	50.1	80.6		
HRFormer (150k) [67]	HRFormer-S	14	44.0	80.0		
	HRFormer-B	50	46.3	81.4		
MaskFormer [68]	Swin-T	42	46.7	78.5		
	Swin-S	63	49.8	-		
Mask2Former [69]	Swin-T	47	47.7	82.1		
	Swin-B	216	-	83.3		
Segformer [26]	MiT-B0	3.8	37.4	76.2		
	MiT-B1	13.7	42.2	78.5		

Ablation experiments of transformer variants and MEMAFormer

Method	Backbone	Parameters (M)	GFLOPs	FPS	mlou
FCN [20]	MobileNetV2	9.8	39.6	64.4	19.7
PSPNet [53]	MobileNetV2	13.7	52.9	57.7	29.7
DeepLabV3+ [61]	MobileNetV2	15.4	69.4	43.1	34.0
Segformer [26]	MiT-B0	3.8	8.4	50.5	37.4
FeedFormer [62]	MiT-B0	4.6	7.8	79.5	39.2
U-MixFormer [66]	MiT-B0	6.1	9.1	55.4	41.2
MEMAFormer	MiT-B0	6.4	14.3	49.4	42.2





The workflow of proposed method to classify security and non-security requirements that rely on LLMs fine-tuning and Few-Shot learning with transfer learning models leveraging pre-trained models using data augmentation

Table 5. Results of the fine-tuned models in terms of accuracy, precision, recall, and F1-score on the four datasets.

Model	Accuracy	Precision	Recall	F1-Score
	Normal	dataset		
BERT-based-uncased	0.89	0.89	0.76	0.82
DistilBERT-base-uncased	0.89	0.89	0.86	0.82
RoBERTa-base	0.9	0.86	0.86	0.86
xlnet-base-uncased	0.9	0.89	0.81	0.85
C	ontextual word a	ugmented dataset	b)	
BERT-based-uncased	0.94	0.9	0.9	0.9
DistilBERT-base-uncased	0.89	0.85	0.81	0.83
RoBERTa-base	0.84	0.92	0.57	0.71
xlnet-base-uncased	0.92	0.94	0.81	0.87
B	ack translation a	ugmented dataset		
BERT-based-uncased	0.92	0.9	0.96	0.88
DistilBERT-base-uncased	0.92	0.9	0.96	0.88
RoBERTa-base	0.84	0.87	0.62	0.71
xlnet-base-uncased	0.9	0.83	0.9	0.86







Step 5: Apply 1 × 1 convolution and Layer Normalization

• Apply 1 × 1 convolution to reduce dimensions of concatenated features:

 $\mathbf{F}_{\text{final}} \leftarrow \text{Conv}_{1 \times 1}(\mathbf{F}_{\text{concat}})$

Apply layer normalization:

 $\mathbf{F}_{enh} \leftarrow LayerNorm(\mathbf{F}_{final})$

Step 6: Global Average Pooling and Classification

· Perform global average pooling on the normalized features:

$$\mathbf{F}_{pool} \leftarrow \mathbf{GAP}(\mathbf{F}_{enh})$$

· Flatten the pooled features and pass through fully connected layers:

$$\mathbf{z}_1 \leftarrow \text{ReLU}(\mathbf{W}_1 \mathbf{F}_{\text{pool}}), \quad \hat{y} \leftarrow \mathbf{W}_2 \mathbf{z}_1$$

Output: Return predicted class probabilities \hat{y} .

Hyperparameter	Configurations					
Input Size	224					
	Adagrad					
	AdamW					
Optimizers	SGD					
	Adam					
	RMSprop					
Batch Size	8, 16, 32					
Learning Rate	0.01, 0.001, 0.0001					
Early Stopping	Patience of 10 epochs					
Loss Function	Categorical Cross Entropy					
Number of Epochs	200					

Performance comparison of different pre-trained models, truncated MobileNetV2, and the proposed MSCPNet model.

Model	Accuracy	Precision	Recall	F1-Score	MCC	Inference Time (s)
DenseNet121	95.53%	94.37%	95.15%	94.69%	0.9396	0.0289
ResNet50	95.21%	93.71%	95.41%	94.37%	0.9356	0.0120
ShuffleNetV2	94.09%	92.20%	93.98%	92.82%	0.9209	0.0107
SqueezeNet	93.61%	92.13%	92.25%	92.12%	0.9135	0.0091
MobileNetV2	96.65%	96.26%	95.85%	96.04%	0.9542	0.0110
Truncated MobileNetV2	95.85%	95.02%	95.19%	95.09%	0.9434	0.0083
Proposed MSCPNet	97.44%	96.76%	97.37%	97.04%	0.9653	0.0111

Performance Comparison of different backbone models with the Proposed MSCPNet block
Backbone Model	Accuracy	Precision	Recall	F1-Score	MCC	Inference Time (s)	Total FLOPs
DenseNet121	95.85%	94.82%	95.76%	95.24%	0.9436	0.0302	3,066,045,440
ResNet50	95.53%	94.16%	94.42%	94.29%	0.9390	0.0180	4,915,031,552
ShuffleNetV2	95.53%	93.96%	95.31%	94.51%	0.9396	0.0150	318,950,600
SqueezeNet	94.41%	92.84%	93.78%	93.25%	0.9241	0.0136	1,385,824,800
MobileNetV2	97.28%	96.37%	97.03%	96.67%	0.9631	0.0181	517,634,624
Proposed MSCPNet	97.44%	96.76%	97.37%	97.04%	0.9653	0.0111	315,258,752

TABLE 7: Performance comparison with different optimizers.

Optimizer	Accuracy	Precision	Recall	F1-Score	MCC
Adagrad	96.33%	95.62%	94.96%	95.26%	0.9498
Adam	97.44%	96.76%	97.37%	97.04%	0.9653
AdamW	96.33%	95.44%	95.75%	95.58%	0.9500
RMSprop	96.49%	95.84%	95.90%	95.87%	0.9521
SGD	97.12%	96.83%	96.16%	96.45%	0.9609



(e) Gray Leaf Spot Heatmap (Predicted Class: Gray Leaf Spot)



(f) Gray Leaf Spot Grad-CAM visualizations for different disease classes









Fig. 2. Structure of Multi-scale Iterative Tamper Detection Net.

- A The network primarily consists of the global-local feature synchronization phase and the tampered area refinement phase.
- G During the global-local feature synchronization stage, the adoption of a parallel network architecture ensures that the network can simultaneously capture feature information at different scales.
- Generation Attention Block (EDAB) is introduced to utilize the image's multi-level features, providing sufficient spatial information for each layer in the network and maintaining close connections between different scales, thereby effectively addressing the issue of scale variations.
- Accurate prediction of tampering masks is achieved through feature fusion, ensuring the comprehensive utilization of features at different scales in tampering detection.
- Additionally, boundary information is processed through the Edge Enhancement Module (EEM) edge supervision branch, which serves as an auxiliary to the main feature extraction process









Table	1:	Quantitative comparison	on
SPAD	[62]	for rain streak removal.	

Method	SPAD [62] PSNR ↑ SSIM ↑			
DDN [13]	36.16	0.9463		
PReNet [50]	40.16	0.9816		
RCDNet [61]	43.36	0.9831		
MPRNet [75]	43.64	0.9844		
SPAIR [46]	44.10	0.9872		
Uformer-S [64]	46.13	0.9913		
SCD-Former [15]	46.89	0.9941		
IDT [67]	47.34	0.9929		
Restormer [73]	47.98	0.9921		
DRSformer [5]	48.53	0.9924		
FPro (Ours)	48.99	0.9936		

Table 2: Quantitative comparison onAGAN-Data [47] for raindrop removal.

	AGAN-Data 47			
Method	$ PSNR\uparrow$	SSIM ↑		
Eigen's [12]	21.31	0.757		
Pix2pix [20]	28.02	0.855		
Uformer-S [64]	29.42	0.906		
WeatherDiff ₁₂₈ [42]	29.66	0.923		
TransWeather [55]	30.17	0.916		
DuRN [33]	31.24	0.926		
RaindropAttn [48]	31.37	0.918		
AttentiveGAN [47]	31.59	0.917		
IDT [67]	31.63	0.936		
Restormer [73]	31.68	0.934		
FPro (Ours)	31.96	0.937		