ISSN: 2278-1862



Journal of Applicable Chemistry

2025, 14 (2): 471-547 (International Peer Reviewed Journal)





CNN - 66D

Iam (Intelligence Augmented /Assisted Method(s))

Transformers - -architectures & Fitz (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,	Dr. M.R.Appa Rao Campus,	Visakhapatnam 530 003,
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: Ashish Vaswani et al. published a paper entitled "Attention is All You Need" in the year 2017. It brought renaissance not only in sequence data processing, but also in computational paradigm with other data structures. The new approach won the favour of data scientists as a whole. This new model gained popularity as Transformer net (TransF Net) or Transformer neural network (TransF NN).

TransF NN contains two important modules, viz., attention layer and MLP-NN which help to carry out Natural Language processing (NLP).

Attention: In 2014, Bahdanau et al. proposed the idea of attention in the context of sequence-to-sequence models, used for neural machine translation (NMT). The attention mechanism targets at improving the performance of sequence-to-sequence models by allowing the model to focus on different parts of the input sequence during calculation of each output token. It overcomes limitation of encoding the entire input sequence into a single fixed-size vector which was earlier practice.

Luong et al. (2015) proposed different scoring mechanisms viz. Dot-Product (multiplicative attention), Generalized (Additive) attention and Concatenation-based attention. Another way of looking at it is global attention (where all encoder states considered) and local attention (focusing on only a subset of encoder states). The multiplicative attention is simple, efficient, and scalable for large datasets. It became a standard approach in later models and is it was a direct precursor to the self-attention mechanism in the Transformer model. The local attention is instrumental in mitigating the computational bottleneck for long input sequences, as it focuses on a smaller subset of tokens.

Attention: The attention in Transformer NN is calculated as

Attention(Q, K, V) = softmax
$$\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

where Q is Query matrix (current token's information), K: Key matrix (context of all tokens in the sequence), V: Value matrix (values to be passed along after attention weighting) and dk: Dimension of the key vectors (used for scaling).

Cross-Attention: It is essentially same except using the query from one sequence with the key and value from another sequence.

Mix-attention: It is computed as

MixAttention =
$$\alpha \cdot \text{SelfAttention} + (1 - \alpha) \cdot \text{CrossAttention}$$

This is also widely employed in sequence text modelling. ALiBi (Attention with Linear Biases) allows more efficient handling of long sequences in transformer models.

Multi-Head Attention: The multiple attention heads are concatenated and passed through a final linear transformation

MultiHead(Q, K, V) = Concat(head₁, ..., head_h)
$$W^{0}$$

and

head_i = Attention
$$(QW_i^Q, KW_i^K, VW_i^V)$$

Where W_i^Q, W_i^K, W_i^V are learned weight matrices for each attention head i.

FlashAttention: It is an optimized CUDA-based approach (on frameworks like PyTorch and TensorFlow) designed to efficiently compute the attention with GPUs for both training and inference of LLMs((GPT, BERT). It results in higher speed and increased scalability.

The evolution of architecture of TransF NN, attention mechanism, and hybridization with other approaches, during these few years, revolutionized computational modelling. This approach is in the front-line in dealing with multi-modal data (viz.Text, numerical time-series, sound (speech), image/video sequence, and tactile-sense-output) with local and global inter-dependencies

A few of Transformer net (TransF-N) architectures documented in this state-of-knowledge-methodsmodule for dataTOknowledge transformation are

la PCNet,

- (MM-HiFuse),
- △ Estimating energy expenditure based on video (E3V) using E3V-K5 dataset,
- *△* Performer with Graph Self-attention Mechanism,
- Grouped Attention and Cross-Layer Fusion Network (GACLFNet),
- learning Neural Ordinary Differential Equation (N-ODE) Transformer,
- *B*ioMechanically Accurate Neural Inverse KINematics solver (MANIKIN),
- Gradient Origin Embeddings (GOEmbed),
- lacktriangletic America Contraction International America Amer
- *△* Triplet convolutional twin transformer,
- lierarchical Multi-Task Learning (HirMTL),
- △ Memory-augmented Deformable Detection TRansformer (MD-DETR),
- Geolocalization with Adapters and Auto-Regressive Transformers (GAReT),
- le Vision Transformer (ViT),
- A RealViformer.,
- 👃 AuraLLM,
- *△* 3D Transformers,
- Multi-Relational Graph Contrastive Learning architecture (MRGCL) using a Multi-relational Graph Hierarchical Attention Networks (MGHAN),
- ⊖ TWiX, C-TWiX,
- △ Squeeze and Excitation based UNet TRansformers (SE-UNETR),
- Squeeze and Excitation based High-Quality Resolution Swin Transformer Network (SE-HQRSTNet)
- lacktriangleright BiMKANsDformer,
- A Hybrid Transformer with Multi-level Fusion for Multimodal Knowledge Graph Completion (MKGformer),
- △ Transformer Choice Net (TCNet) [Transformer Neural Network for Choice of Transformer],
- Ĝ Transformer-Based Framework
- A Physics-Informed Neural Networks (PhysIinfNNSFormer), and
- △ EfficientNet and Vision Transformer (ViT)-based Swin Ttransformer (SwinE-Net)

Keywords: Artificial intelligence (AI); Convolution Neural Nets-- Capsule Neural Nets-- MLP-Attention Mechanism-TransFormer Nets--Hybrid TransFormer Networks--



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;]





Key point detection results

• Blue points: torso and limb keypoints

• Red points: supplementary extremity keypoints





Framework. E3SFormer

- Human skeleton sequence x is extracted using a pose estimator from the video
- Then fed into a backbone to obtain motion representation F.
- Then sent to an action recognition branch (upper) and an energy estimation regression branch (lower). The category-related joint-specific attention Ac from the action recognition branch is transferred to the energy estimation regression branch to boost its performance.
- The multi-modal data z are used for more personalized energy estimation



Table 1: Survey of deep/RL methods in solving TSPs.

Mothods	Learning type	Network structure	
Vinyals et al. (2015)	SL + AR	LSTM	
Bello et al. (2016)	RL + AR LSTM		
Khalil et al. (2017)	RL + AR	GNN	
Nazari et al. (2018)	RL + AR	LSTM	
Kool et al. (2018)	RL + AR	Transformer	
Deudon et al. (2018)	RL + AR	Transformer	
Ma et al. (2019)	RL + AR	GNN	
Bresson & Laurent. (2021)	RL + AR	Transformer	
Wu et al. (2022)	RL + AR	GNN	
Lei et al. (2022)	RL + AR	GAT	
Yang et al. (2023)	RL + AR	Transformer	
Zhu et al. (2023)	RL + AR	CosFormer	
Wang et al. (2023)	RL + AR	BERT	
lung et al. (2024)	RL + AR	CNN-Transformer	
Zhu et al. (2024)	RL + AR	GEIAM	
Zhao & Gu (2024)	RL + AR	GPN	
Nowak et al. (2016)	SL + NAR	GNN	
loshi et al. (2019)	SL + NAR	GCN	
Xiao et al. (2023)	RL + NAR	GNN	





Table 4: Performance comparisons of our method with benchmarking methods.

		TSP20		
Method	Туре	Length	Gap (%)	Time (ms)
Gurobi	Exact	24.25	0.00	1.77e + 01
Greedy algorithm (Dijkstra, 1959)	Heuristic	26.51	9.30	6.00e-02
Nearest insertion (Held & Karp, 1961)	Heuristic	29.00	19.58	1.80e-01
2-OPT (Lin & Kernighan, 1973)	Heuristic	26.16	7.87	8.45e + 00
Fastest insertion (Chvátal et al., 2010)	Heuristic	27.84	14.79	1.83e + 00
Christofides (Christofides, 2022)	Heuristic	26.51	9.32	3.50e + 00
A (Aarts & Korst, 1989)	Meta-heuristic	26.34	8.61	7.77e + 02
GA (Chung & Xu, 2012)	Meta-heuristic	25.16	3.73	1.77e + 04
Ant colony optimization (Zhao et al., 2016)	Meta-heuristic	26.33	8.56	4.55e + 02
Hill climbing (Yelmewad & Talawar, 2019)	Meta-heuristic	28.13	16.01	1.86e + 03
NAR4TSP(GS)	RL, NAR	25.75	6.19	3.61e-01
VAR4TSP(BS), B=100	RL, NAR	24.59	1.40	3.54e-01
NAR4TSP(BS), B=1000	RL, NAR	24.53	1.15	3.82e-01
Durs (GS)	TRL, NAR	25.20	3.90	3.51e-01
Durs (BS), B=100	TRL, NAR	24.50	1.05	3.65e-01
Durs (BS), B=1000	TRL, NAR	24.47	0.92	3.77e-01

• 'TRL' stands for tr ansfer r einforcement learning;

- '2-OPT' stands for 2-OPT Local Search;
- 'GS' stands for greedy search;
- 'BS' represents beam search with a width *B*.
- 'Time' indicates the inference time of solving a TSP instance on average
- batch sizes:TSP20, TSP50, and TSP100 in searching solutions are 1024, 128, and 16,







AAA: 66D-Transformers-architectures & Fits





Proportion of Normal and Abnormal Traffic on the Catalytic Reforming Unit Process Platform







Fig. 1: We propose the **GOEmbed** (Gradient Origin **Embed**ding) mechanism that encodes source views (o^{ctxt}) and camera parameters (ϕ^{ctxt}) into arbitrary 3D Radiance-Field representations g(c,d) (sec. 3). We show how these general-purpose GOEmbeddings can be used in the context of 3D DFMs (Diffusion with Forward Models) (sec. 5) and for sparse-view 3D reconstruction (sec. 6).



Fig. 2: GOEmbed illustration. We demonstrate the mechanism here using the Triplane representation for g(c, d), but note that this can be applied to other representations as well. The GOEmbed mechanism (eq. 1) consists of two steps. First we render the origin ζ_0 from the context-poses ϕ^{ctxt} ; then we compute the gradient of the MSE between the renders and the source-views o^{ctxt} wrt. the origin ζ_0 which gives us the GOEmbed encoding ζ_{enc} .















Fig. 2: Schematic of the FPC-*i* module, where $i \in \{1, 2, 3\}$, highlighting its role in the hierarchical multi-task learning architecture. The FPC-*i* module is integral for promoting single-scale level interactions by propagating and enhancing features from higher-level semantic and saliency detection tasks into lower-level but larger resolution representations, thereby establishing a foundational layer for subsequent multi-scale and task-level feature fusion within the HirMTL framework.



Visual illustration of the TAF module effectiveness

+ Demonstrating its role in adaptively fusing features across multiple scales.

- ▲ Leftmost column: input images
- G Columns two to five: four scales of features
- △ Final two columns: taf's sophisticated multi-scale fusion results,
 - showcasing the heatmaps of fusion features for semantic segmentation (semseg) and salient object detection (sal),
- \checkmark The fusion ratios across scales are shown on the left side of each heatmap,
 - showcasing the module's adeptness in customizing feature fusion for different task requirements: focusing more on the boundaries of objects for semantic segmentation, while more on the centers of objects for salient object detection







Heat map visualizations for one class

- ✓ Reflecting the incremental impact of the FPC, TAF, and AICM modules on the semantic segmentation predictions of HirMTL on NYUD-v2
- ✓ The gradation from left to right corresponds to the successive addition of modules, vividly illustrating the refinement in predictions.



- \checkmark Given an input image x,
- ✓ query function Q(x, $\theta \nabla$, α) proposed to retrieve relevant memory units as a linear combination.
- ✓ The obtained information from the memory is utilized by the decoder across various decoding layers.
- ✓ The majority of the architecture remains frozen, encompassing the encoder and decoder; the trainable modules consist of memory units M, class embedding, bounding box embedding, and ranking function $g \phi$.







conditioned scenario. First, we sample an initial layout from a base distribution and a time t. Then, an intermediate sample g_t is calculated by linearly interpolating between the initial sample and the ground truth layout. Each intermediate element is embedded jointly with the given element condition \tilde{a}^k . Lastly, the Transformer architecture takes all the element embeddings to predict a vector field.







o (3) A classification module generating the final outcome scores













Fig. 2: Alignability-Verification based Metric Learning is proposed to is proposed to decide how well two video instances are alignable and produce an 'alignability score' for effective learning from a limited labeled set \mathbb{D}_l . Our approach employs a triplet loss (\mathcal{L}_{AT}) , considering videos from identical action classes as positive and those from different classes as negative. We selectively mine hard-negatives from the sampled minibatch based on alignment distance, presenting a challenging learning task for the model f_A . Additionally, we incorporate a matching loss \mathcal{L}_{score} to quantify the alignment between videos, serving as a verification task to determine whether a video pair belongs to the same class (i.e. alignable or target label = 1) or different classes (i.e. non-alignable or target label = 0). Further details are provided in Sec. 3.1.



Fig. 3: Collaborative Pseudo-labeling: The unlabeled instance $\mathbf{u}^{(i)}$ undergoes processing by both video encoders (f_E and f_A). For the Action Encoder f_E , its prediction (\mathbf{p}_E) is derived via its classification head. For the Alignability Encoder f_A , the embedding of $\mathbf{u}^{(i)}$ computes class-wise alignability scores against a gallery of labeled embeddings A. These scores are then used to generate a class-wise prediction \mathbf{p}_A using the non-parametric classifier ϕ_A . As these predictions stem from distinct supervisory signals— \mathbf{p}_E from video-level and \mathbf{p}_A from alignability-based supervision—they offer complementary insights, resulting in a refined collaborative pseudo-label.





(a) Apply EAA unit and EAF unit to the sliding-window-based method (e.g., EDVR)

(b) Apply the EAA unit and the EAF unit to the bidirectional recurrent-based method (e.g., BasicVSR








Overview of proposed method

(a) showcases our pipeline, which adopts an innovative strategy focused on learning degradation residual and employs the information-rich condition to guide the diffusion process.

(b) illustrates the utilization of our prompt pool, which empowers the network to autonomously select attributes needed to construct adaptive weather-prompts.

(c) depicts the general prompts directed by depth-anything constraint to supply scene information that aids in reconstructing residuals.

(d) shows the contrastive prompt loss, which exerts constraints on prompts driven by two distinct motivations, enhancing their representations



Selection Frequency of Sub-Prompts for Different Tasks

Selection frequency of sub-prompts

- \checkmark Some similar selection frequencies reflect the network's ability to adaptively
- exploit common attributes in some similarity between tasks (e.g. rain and raindrop).
- ✓ At the same time, the unique prompt frequencies highlight the flexibility to adapt to the specific characteristics of each weather condition.



weather conditions



- ✓ (A) We begin by optimizing our image transformer encoders
- \checkmark (B) with street-view frame and matching small aerial image pair.
- ✓ (C) Then, for adapting our image encoder to video inputs, we add our GeoAdapter GA module and only optimize the adapter parameters with video pairs as inputs, i.e., a street-view video V s and corresponding large aerial image IaL For training,
- ✓ we sample every kth frame from the street-view video and partition the large aerial image into non-overlapping patches
- ✓ (D) In GA, we apply temporal selfattention (TSA) computation only on the CLS tokens.
- \checkmark For TSA computation, we reuse the spatial self-attention weights.
- (E) During inference, we first perform a Sequence-to-Image inference procedure, where given a query street-view video,
 - The unified module $U = \{T, GA\}$ produces feature embeddings for both the V s and IaL
 - Then, using embeddings, we retrieve the t nearest neighbor large aerial images (here we show t =
 - i. and construct a small aerial image gallery G. (F)
 - ii. Finally, GA is removed, and feature embeddings for Ia sk and V s k are obtained.
 - These features are then passed to our TransRetriever TAR model to obtain final frame-byframe GPS predictions to construct a GPS trajectory.







Fig. 1: (a) Designing a RWVSR transformer is not trivial. A Swin-based transformer suited for standard VSR hallucinates more lines than a RealBasicVSR, a convolutional state-of-the-art. We propose RealViformer based on our investigation of attention under the RWVSR setting. RealViformer generates details with fewer artifacts than RealBasicVSR [3] and the Swin-based VSR model. (b) Schematic for spatial and channel attention. Spatial attention aggregates features based on pixel representations. Channel attention takes $H \times W$ feature map for matching across channels.



Fig. 3: (a) The recurrent baseline in Sec. 3.2 has a shallow mapping module \mathcal{F} , reconstruction module \mathcal{R} , upsampling module \mathcal{U} and warping function W. W aligns the hidden state h_{t-1} to feature at t based on optical flow $s_{(t-1)\to t}^{f}$. All residual blocks are convolutional. The concatenation between f_t and \hat{h}_{t-1} are replaced with the spatial or channel attention modules in (b) to compare the effect of attention. (b) The attention module first applies layer normalization to f_t and \hat{h}_{t-1} and then performs channel or spatial attention according to Sec. 3.1. The output feature O_t^A concatenated with f_t is processed by the module \mathcal{R} in (a).



Fig. 5: Improved Channel Attention Module (ICA), showing self-attention for simplicity. The 'squeeze' convolution compresses the number of input feature channels $X \in \mathbb{R}^{C \times H \times W}$ by ratio r. The features are then rescaled by weights predicted from the $\frac{C}{r} \times \frac{C}{r}$ attention map before being expanded by the 'excite' convolution back to the original number of input channels.



(c) Reconstruction Module

Fig. 6: The framework of RealViformer. (a) Overview of RealViformer, following a unidirectional recurrent framework. The outputs of the Forward module are propagated to the next time step and upsampled by module \mathcal{U} to get HR frames. (b) Explanation of the Forward module in (a), where W denotes the warping function. The reconstruction module \mathcal{R} takes current frame I_t^L and warped hidden state \hat{h}_{t-1} as inputs. (c) Reconstruction module \mathcal{R} . The shallow feature of I_t^L and \hat{h}_{t-1} are fused by CAF and then forwarded to Transformer blocks with U-shape connection [39]. Module GDFN follows Restormer [39]. Details of CAF and ICA modules are stated in Fig. 7 and Fig. 5.



Fig. 7: Details of Channel Attention Fusion (CAF) module. CAF gets the query from current frame feature f_t and {key, value} from hidden state \hat{h}_{t-1} . The attention output is concatenated with f_t to process for module output O_t .



Fig. 11: (a) Visual comparison between RealViformer and its ablations. Red circles highlight the improved details. (b) Radial Power Spectrum (RPS) of model predictions. Using ICA improves the power of high-frequency components (blue region).







We visualize Mamba-3D as an example.

- Given 3D input, we patchify it into L patches.
- During this process, we maintain the original 3D structure of the input.
- This sequence is then passed through K Mamba-ND blocks, each of which consists of a chain of 1D Mamba layers that process the sequence in alternating orderings.
- In 3D space, we use the order H+H-W+W-T+T-. In 2D space, the sequence would be H+H-W+W-. Finally, the sequence is reshaped back to its original 3D structure and passed to task-specific heads for downstream processing



Fig. 3: Variations of SSM Layer Design. Col 1 represents the standard 1D SSM layer. Col 2 represents Bi-SSM, which adds bidirectionality in a similar fashion as LSTM. Col 3 represents ND-SSM block, which extends Bi-SSM to more directions. Col 4 represents multi-head SSM block inspired by multi-head attention in Transformers.





(a) Different ways of arranging Mamba tion, there is only 1 sequence. Col 2: Factor**layers.** The first row visualizes alternating- izing the 3D sequence into D 2D sequences, directional design. The second row visualizes where D is the length of a single dimension. bidirectional design. The third row visualizes Col 3: Factorizing the 3D sequence into D^2 quad-directional design.

(b) Visualization of various Scan-Factorization policies. Col 1: No factoriza-1D sequences.

Visualization of block level design and factorization policies

Table 6: Ablation Study on Layer Designs. We report top-1 accuracy on the ImageNet-1K validation set. The Alt-Directional design is the top-performing one.

		IN1K↑	HMDB-51 \uparrow
Alt-Directional	Block Level	79.4	59.0
Multi-Head-SSM	Layer-Level	77.6	51.5
ND-SSM	Layer-Level	77.2	46.7
1D-SSM	-	76.4	34.9
Bi-SSM	Layer-Level	74.6	32.1

877 / 877 / 877 / 877 / 877 / 877 / 877 / 877 / 877 / 877 / 8

Transformer Net

2025-128

Reference	Methods	Time series	Metrics	Data	Best model
(Masum et al., 2018)	ARIMA and LSMT	Electric load, day	RMSE	Great Britain, Poland and Italy	LSTM
(de Oliveira and Cyrino Oliveira, 2018)	Bagging ARIMA and Exponential Smoothing	Mid-long term electric energy consumption, monthly	ASM, sMAPE, RMSE and TIC	Canada, France, Italy, Japan, Brazil, Mexico and Turkey	Remainder Sieve Bootstrap
(Khan and Osińska, 2021)	Fractional-order Grey Model and ARIMA	2019 statistical review of world energy, yearly	MAPE and MSE	Brazil, Russia, China, India, and the Republic of South Africa (BRICS)	ARIMA
(Banik et al., 2021)	Random Forest, Ensemble Learning, Boosting and XGBoost	Electricity load, hourly, weekly and monthly	R ² and RMSE	Tripura in India	Ensemble RF-XGBoost
(Farsi et al., 2021)	ARIMA, LSTM, and Parallel LSTM-CNN Network.	Hourly load consumption	MAPE, R ² and RMSE	Germany and Malaysia	Parallel LSTM-CNN Network
(Zhao et al., 2021)	GRU, LSTM, RNN, TCN, CNN+LSTM and Transformer	Day-ahead load forecasting	MAPE	Australia	Transformer network
(Chaturvedi et al., 2022)	SARIMA, LSTM RNN and Facebook Prophet	Monthly energy demand, monthly	MAPE and RMSE	India	Facebook Prophet
(Panagiotou and Dounis, 2022)	RNN, ANFIS, and LSTM	Energy consumption, hourly	MSE, R, R ² , MAPE and CI	Hospital Building's Energy Consumption	ANFIS LSTM
(Henzel et al., 2022)	Naive Methods, Linear Regression, LSTM, and the Facebook Prophet Method	Energy consumption, hourly	MSE and MAPE	Digital-Twin Model of the Building	Facebook Prophet

(Shohan et al., 2022)	ANN, LSTM, Facebook Prophet and LSTM- Facebook Prophet	Hourly load demand, hourly	MAPE, R ² , SSE and RMSE	Florida	LSTM- Facebook Prophet
(Ribeiro et al., 2022)	ARIMA, SVR, Random Forest, XGBoost, RNN, LSTM, and GRU	Energy consumption, hourly	RMSE, MAPE and MAE	ESCO (Energy Service Company) Ireland	XGBoost
(Shin and Woo, 2022)	Random Forest, XGBoost and LSTM	Energy consumption, monthly	RMSE and MAPE	Korea	LSTM and Random Forest
(Sulandari et al., 2023)	ARIMA, Exponential Smoothing, Facebook Prophet, Neural Networks, and Proposed Ensemble Methods	Generation of electricity, hourly	MAPE and RMSE	US, Ontario, England, Wales and Australia	Proposed Ensemble Methods
(Pierre et al., 2023)	ARIMA, LSTM, GRU, ARIMA-LSTM, and ARIMA-GRU	Electricity demand, hourly	MAPE and RMSE	Benin Electricity Company (CEB)	Hybrid Approach
(Koukaras et al., 2024)	HGBR, LGBMR, ETR, RR, BRR, CBR	Energy consumption, hourly	R ² , RMSE, CVRMSE, NRMSE and MAE	ITI/CERTH Smart House, Thessaloniki (Greece)	HGBR, LGBMR



- Most previous works only perform the fusion globally or locally [73]
- Present DVLO designs a local-to-global fusion strategy that facilitates the interaction of global information while preserving local fine-grained information
- Furthermore, a bi-directional structure alignment is designed to maximize the inter-modality complementarity





Overview. method for the UDA-OD task.

! Architecture. teacher–student Alg.

- Teacher model generates pseudo-labels based on weakly augmented target domain images
- Student model is trained using both downsampled and strongly augmented inputs.
- CSPC is enforced by supervising the student model with inputs of different resolutions improving detection of objects at various scales
- ! Temporal ensemble is employed for
 - Robust pseudo-label selection
 - Combining classification confidence and box matching based on
 - Intersection over Union (IoU) to ensure high-quality pseudo-labels.

! ICFC module

- o Aligns object-level features across scales and augmentations,
 - Utilizing contrastive learning
- \checkmark To ensure intra-class attraction and inter-class repulsion,
 - Enhancing the consistency of object representations

Method	Backbone	Detector	Split	Bicycle	Bus	Car
Source	R50-FPN	Faster R-CNN	0.02	38.4	27.5	44.9
Oracle	R50-FPN	Faster R-CNN	0.02	51.2	52.7	74.1
AT [4]	V16	Faster R-CNN	0.02	51.3	64.9	63.6
CMT [5]	V16	Faster R-CNN	0.02	51.2	66.0	63.7
MOTOR [32]	R50	Faster R-CNN	0.02	35.6	38.6	44.0
SFA [33]	R50	Deform-Detr	0.02	44.0	46.2	62.6
MTTrans [34]	R50	Deform-Detr	0.02	46.5	45.9	65.2
DA-DETR [35]	R50	Deform-Detr	0.02	46.5	45.9	63.1
O ² Net [36]	R50	Deform-Detr	0.02	45.9	47.6	63.6
AQT [37]	R50	Deform-Detr	0.02	46.4	53.7	64.4
MTM [38]	R50	Deform-Detr	0.02	47.7	54.4	67.2
TDD [39]	R50	Faster R-CNN	0.02	49.1	51.1	64.1
MRT [40]	R50	Faster R-CNN	0.02	47.1	58.1	68.7
SA-DA-Faster [7]	R50-FPN	Faster R-CNN	0.02	45.4	50.3	62.1
AT [4]	R50-FPN	Faster R-CNN	0.02	53.3	52.1	66.0
CMT [5]	R50-FPN	Faster R-CNN	0.02	53.1	55.0	66.7
Ours	R50-FPN	Faster R-CNN	0.02	57.7	57.5	69.0
Source	R50-FPN	Faster R-CNN	All	50.8	46.7	62.4
Oracle	R50-FPN	Faster R-CNN	All	55.0	57.9	72.2
PDA [41]	V16	Faster R-CNN	All	35.9	44.1	54.4
ICR-CCR [42]	V16	Faster R-CNN	All	34.6	36.4	49.2
AT [4]	R50-FPN	Faster R-CNN	All	53.1	59.0	69.5
CMT [5]	R50-FPN	Faster R-CNN	All	53.6	58.6	69.9
Ours	R50-FPN	Faster R-CNN	All	53.0	63.0	71.9









Fig. 3: The training pipeline of our DDIR model. It consists of double branches, which are the deformation branch and the SR branch. Each branch is composed of an encoder and an MLP, taking the LR image and query coordinates as the inputs. The appearance embedding l_a is computed as the spatial average pooling of the 2D feature map from the encoder E_{ϕ}^{sr} of the SR branch, which is fed into the decoding function $f_{\theta'}^d$ of the deformation branch by concatenation. The RGB output of the deformation branch is supervised by the deformation field. Then, the predicted deformation field feeds into the decoding function f_{θ}^{sr} of the SR branch outputs the target high-resolution RGB values at the query coordinates. Combining the appearance embedding and the deformation field, our DDIR model learns the dual-level deformable implicit representation to address the deformations at the image and pixel levels simultaneously.



Fig. 3: The training pipeline of our DDIR model. It consists of double branches, which are the deformation branch and the SR branch. Each branch is composed of an encoder and an MLP, taking the LR image and query coordinates as the inputs. The appearance embedding l_a is computed as the spatial average pooling of the 2D feature map from the encoder E_{ϕ}^{sr} of the SR branch, which is fed into the decoding function $f_{\theta'}^d$ of the deformation branch by concatenation. The RGB output of the deformation branch is supervised by the deformation field. Then, the predicted deformation field feeds into the decoding function f_{θ}^{sr} of the SR branch outputs the target high-resolution RGB values at the query coordinates. Combining the appearance embedding and the deformation field, our DDIR model learns the dual-level deformable implicit representation to address the deformations at the image and pixel levels simultaneously.



Name	DatasetID	# instances	# features	# classes
robert	41165	10000	7200	10
riccardo	41161	20000	4296	2
guillermo	41159	20000	4296	2
dilbert	41163	10000	2000	5
christine	41142	5418	1636	2
cnae-9	1468	1080	856	9
fabert	41164	8237	800	7
Fashion-MNIST	40996	70000	784	10
KDDCup09_appetency	1111	50000	230	2
mfeat-factors	12	2000	216	10
volkert	41166	58310	180	10
APSFailure	41138	76000	170	2
jasmine	41143	2984	144	2
nomao	1486	34465	118	2
albert	41147	425240	78	2
dionis	41167	416188	60	355
jannis	41168	83733	54	4
covertype	1596	581012	54	7
MiniBooNE	41150	130064	50	2
connect-4	40668	67557	42	3
kr-vs-kp	3	3196	36	2
higgs	23512	98050	28	2
helena	41169	65196	27	100
kc1	1067	2109	21	2
numerai28.6	23517	96320	21	2
credit-g	31	1000	20	2
sylvine	41146	5124	20	2
segment	40984	2310	16	7
vehicle	54	846	18	4
bank-marketing	1461	45211	16	2
Australian	40981	690	14	2
adult	1590	48842	14	2
Amazon_employee_access	4135	32769	9	2
shuttle	40685	58000	9	7
airlines	1169	539383	7	2
car	40975	1728	6	4
jungle_chess_2pcs_raw_endgame_complete	41027	44819	6	3
phoneme	1489	5404	5	2
blood-transfusion-service-center	1464	748	4	2















• The output features of ica-gcn are passed through an average global pooling operation











Converting them to binary spikes during inference 0















Blue colour: paths travelled in the air by the robot









		1 100 100 100 100 100 100 100 100 100 100 100 100 100 100 100 100 100 100
Transformer Net	2025-150	
l. Marina sur		



FIGURE 1. Four different subjects depicting various clinical statuses (normal, MCI, DEM) for FDG, FBP, and FMM. The "Sum" column displays dual tracer images (the 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 the 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.




Model	1D-Wave			
WIGGET	Loss	rMAE	rRMSE	
PINNs	1.93e-2	0.326	0.335	
PINNsFormer	1.38e-2	0.270	0.283	
PINNs + NTK	6.34e-3	0.140	0.149	
PINNsFormer + NTK	4.21e-3	0.054	0.058	



Model	Navier-Stokes			
	Loss	rMAE	rRMSE	
PINNs	6.72e-5	13.08	9.08	
QRes	2.24e-4	6.41	4.45	
FLS	9.54e-6	3.98	2.77	
PINNsFormer	6.66e-6	0.384	0.280	

✓ **PINNsFormer outperforms all baselines on all metrics.**

Results for solving convection and 1D-reaction equations using Transformer architecture with different activation functions

Activation	Convection			1D-Reaction		
Activation	Loss	rMAE	rRMSE	Loss	rMAE	rRMSE
ReLU	0.5256	1.001	1.001	0.2083	0.994	0.996
Sigmoid	0.1618	1.112	1.223	0.1998	0.991	0.993
Sin	0.3159	1.074	1.141	4.9e-6	0.017	0.032
ReLU+LN	0.7818	1.001	1.002	0.2028	0.992	0.993
Sigmoid+LN	0.0549	0.941	0.967	0.2063	0.992	0.990
Sin+LN	0.3219	1.083	1.156	4.7e-6	0.016	0.033
Wavelet	3.7e-5	0.023	0.027	3.0e-6	0.015	0.030
Wavelet+LN	NaN	NaN	NaN	3.9e-6	0.018	0.037

• PINNsFormer (withWavelet activation) consistently outperforms all other activation functions in terms of training loss, rMAE, and rRMSE





nsformer Net		2025-155			
Component	$\mu_{\mathbf{x}_{\mathrm{out}}}$	$\sigma^{2}_{\mathbf{x}_{\mathrm{out}}}$	$\sigma^{2}_{\mathbf{g}_{\mathrm{in}}}$	$\mathbf{r}^{\mathbf{l}}_{\mathbf{x}_{out}}$	$r^{l}_{\mathbf{g}_{in}}$
Embeddings	0	$\sum \sigma_{w_{\mathrm{embd}}}^2$	-	$\frac{\pi^2}{18*\log(V)^2} + \frac{2}{9}$	-
Linear $(d_{\rm in} \rightarrow d_{\rm out})$	0	$d_{\rm in}\sigma_w^2(\sigma_{x_{\rm in}}^2+\mu_{x_{\rm in}}^2)$	$d_{\rm out}\sigma_w^2\sigma_{g_{\rm out}}^2$	$\frac{r_{x_{\rm in}}^l + \mu_{x_{\rm in}}^2/\sigma_{x_{\rm in}}^2}{1 + \mu_{x_{\rm in}}^2/\sigma_{x_{\rm in}}^2}$	$r^l_{g_{ m out}}$
ReLU	$rac{\sigma_{x_{ m in}}}{\sqrt{(2\pi)}}$	$\frac{(\pi-1)}{(2\pi)}\sigma_{x_{\rm in}}^2$	$\frac{1}{2}\sigma_{g_{\rm out}}^2$	$0.7r_{x_{\rm in}}^l + 0.3r_{x_{\rm in}}^{l^{-2}}$	$(\frac{1}{2} + \frac{\sin^{-1}\left(r_{x_{\mathrm{in}}}^{l}\right)}{\pi})\mathbf{r}_{\mathrm{gout}}^{\mathrm{l}}$
LayerNorm (d)	0	1	$rac{\sigma_{g_{ m out}}^2}{\sigma_{x_{ m in}}^2}$	$r^l_{x_{\mathrm{in}}}$	$r^l_{g_{ m out}}$
Dropout (p)	$\mu_{x_{in}}$	$\frac{\sigma_{x_{\mathrm{in}}}^2 + p \mu_{x_{\mathrm{in}}}^2}{1-p}$	$\frac{1}{1-p}\sigma_{g_{\rm out}}^2$	$\frac{r_{x_{\rm in}}^l(1-p)}{1+p\mu_{x_{\rm in}}^2/\sigma_{x_{\rm in}}^2}$	$(1-p)r_{g_{\mathrm{out}}}^l$
SHA-without V	0	$r_{x_{ ext{in}}}^{l}\sigma_{x_{ ext{in}}}^{2}$	$r_{g_{ m out}}^{l}\sigma_{g_{ m out}}^{2}$	1	1
Softmax	$\frac{1}{L}$	$\frac{e^{(1-r_{x_{\rm in}}^d)\sigma_{x_{\rm in}}^2}-1}{L^2}$	$\frac{e^{(1-r_{x_{\mathrm{in}}}^d)\sigma_{x_{\mathrm{in}}}^2}}{L^2}\sigma_{g_{\mathrm{out}}}^2$		-
Signal propag	ation for	r forward and bac	ckward passes t	hrough componen	ts of a transformer