

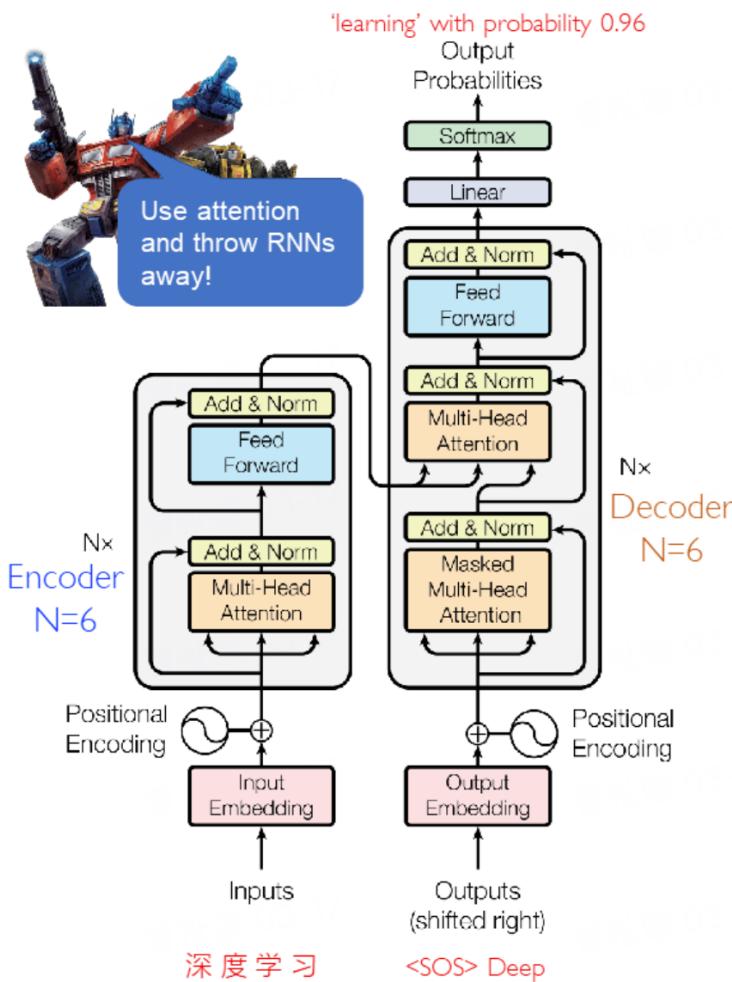
Solving Partial Differential Equation By Transformer

BY ZHIJUN ZENG

2024年3月17日

1 Recap Of Transformer

Transformer



- Main components:

- ① Scaled Dot-Product Attention
- ② (Masked) Multi-Head Attention
- ③ Position-wise FFN
- ④ Residual Connections
- ⑤ Layer Normalization
- ⑥ Positional encoding

- Encoder-decoder architecture:

- ① Encoder
- ② Decoder with Masking
- ③ Encoder-Decoder Attention

① Scaled Dot-Product Attention

- Recall relevance $e_{ij} = v_a^T \tanh(W_a s_{i-1} + U_a h_j)$: too complex!
- Given **query** q (to match others) and **key** k (to be matched):
 - Bilinear

$$a(q, k) = q^T W k \longrightarrow \text{Parameter-heavy}$$

- Dot Product

$$a(q, k) = q^T k \longrightarrow \text{Dimension-sensitive}$$

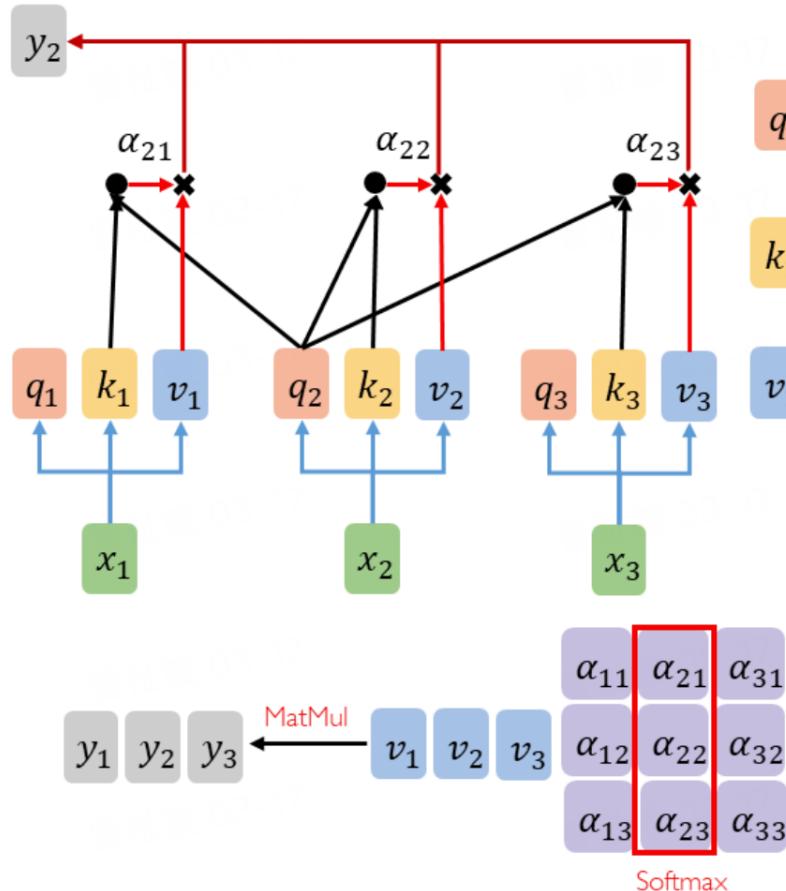
- Scaled Dot-Product

$$a(q, k) = q^T k / \sqrt{d_k} \longrightarrow \text{Dimension of } q \text{ and } k$$

- For independent q_i, k_i with mean **0** and variance **1**:
- Dot-product $q^T k = \sum_{i=1}^{d_k} q_i k_i$ has mean **0** and variance d_k .
 - » Will not cause large values going into the saturation of Softmax

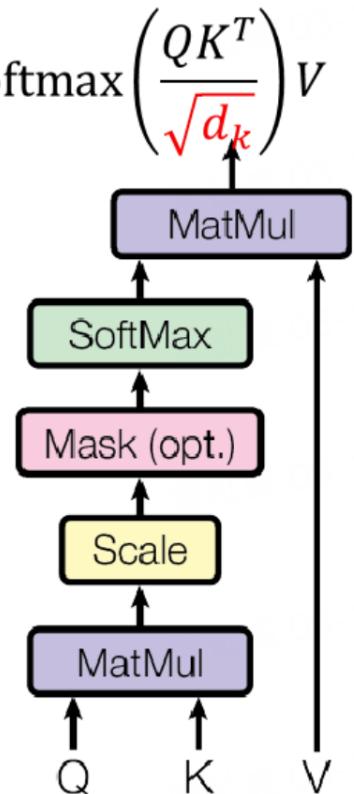
① Scaled Dot-Product Attention

$$y_1 = \sum_i \alpha_{1i} v_i \quad \alpha_{1i} = \frac{\exp(q_1 k_i)}{\sum_j \exp(q_1 k_j)}$$



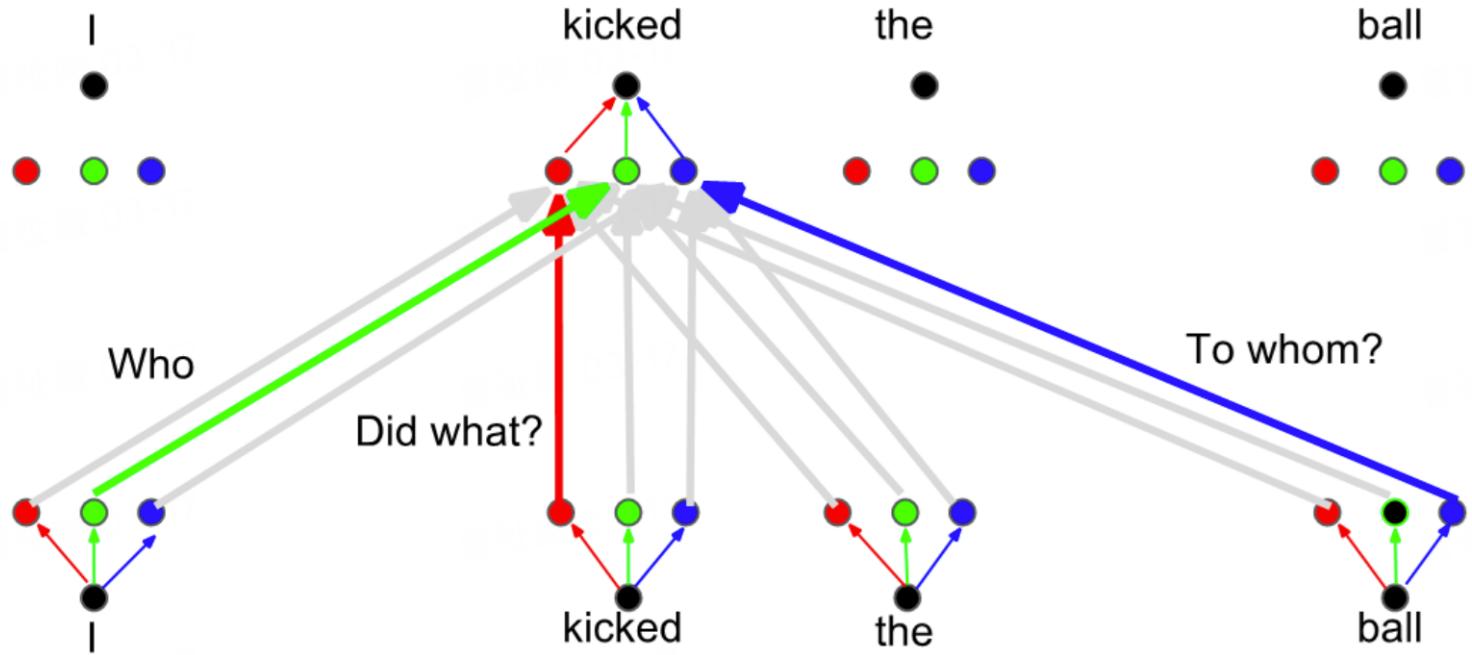
$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\begin{aligned} q_1 & q_2 & q_3 &= W^q & x_1 & x_2 & x_3 \\ Q & & & & & & \\ k_1 & k_2 & k_3 &= W^k & x_1 & x_2 & x_3 \\ K & & & & & & \\ v_1 & v_2 & v_3 &= W^v & x_1 & x_2 & x_3 \\ V & & & & & & \\ q_i &= W^q x_i \\ k_i &= W^k x_i \\ v_i &= W^v x_i \end{aligned}$$



Ignore $\sqrt{d_k}$ for simplicity

② Multi-Head Attention



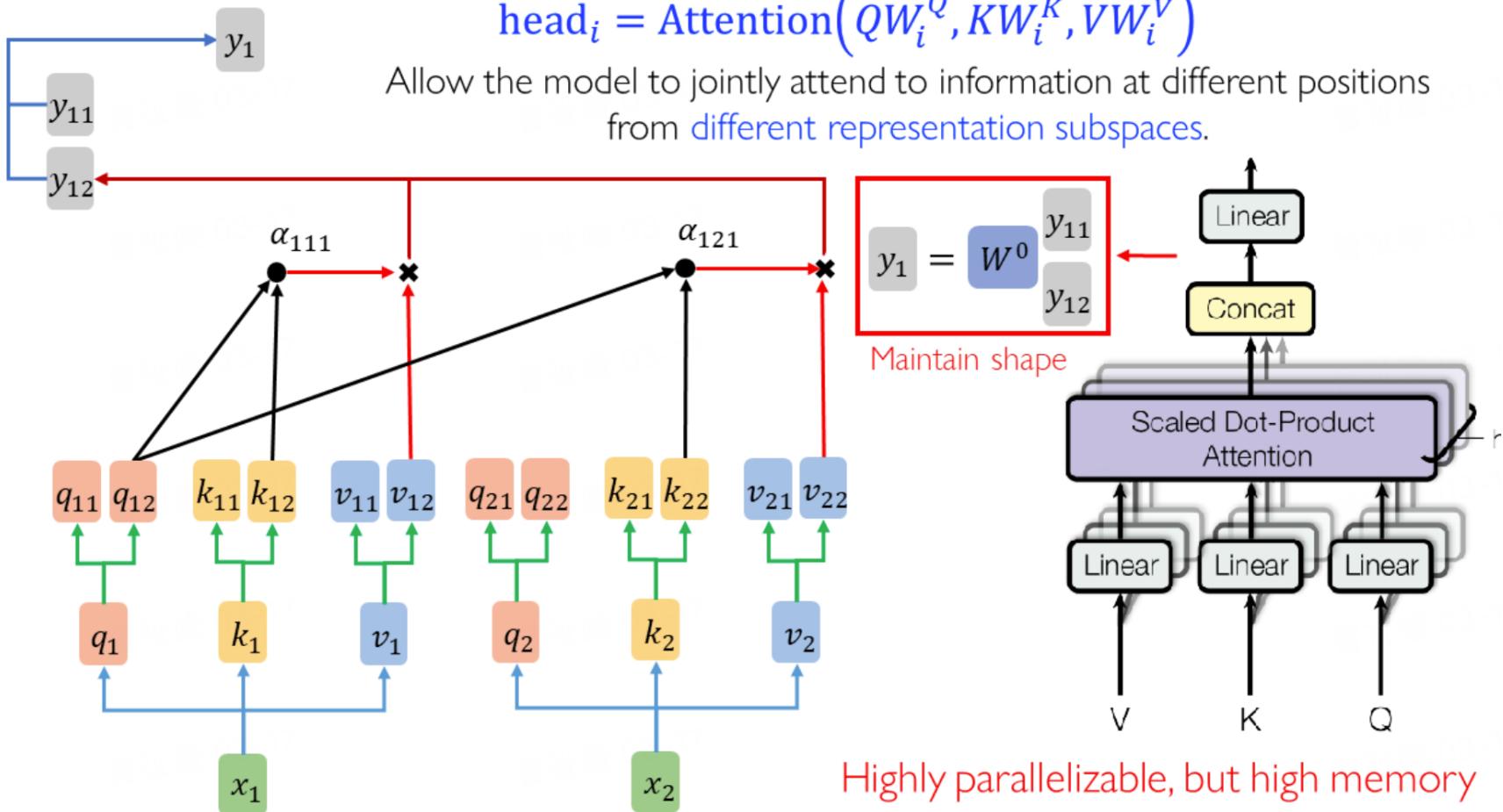
Different **heads** extract different information (like **channels** in CNNs)

② Multi-Head Attention

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

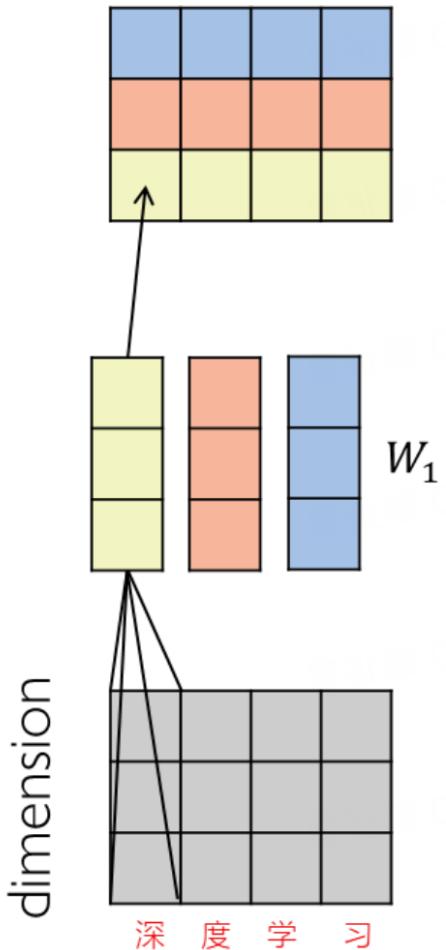
$$\text{head}_i = \text{Attention}\left(QW_i^Q, KW_i^K, VW_i^V\right)$$

Allow the model to jointly attend to information at different positions from different representation subspaces.

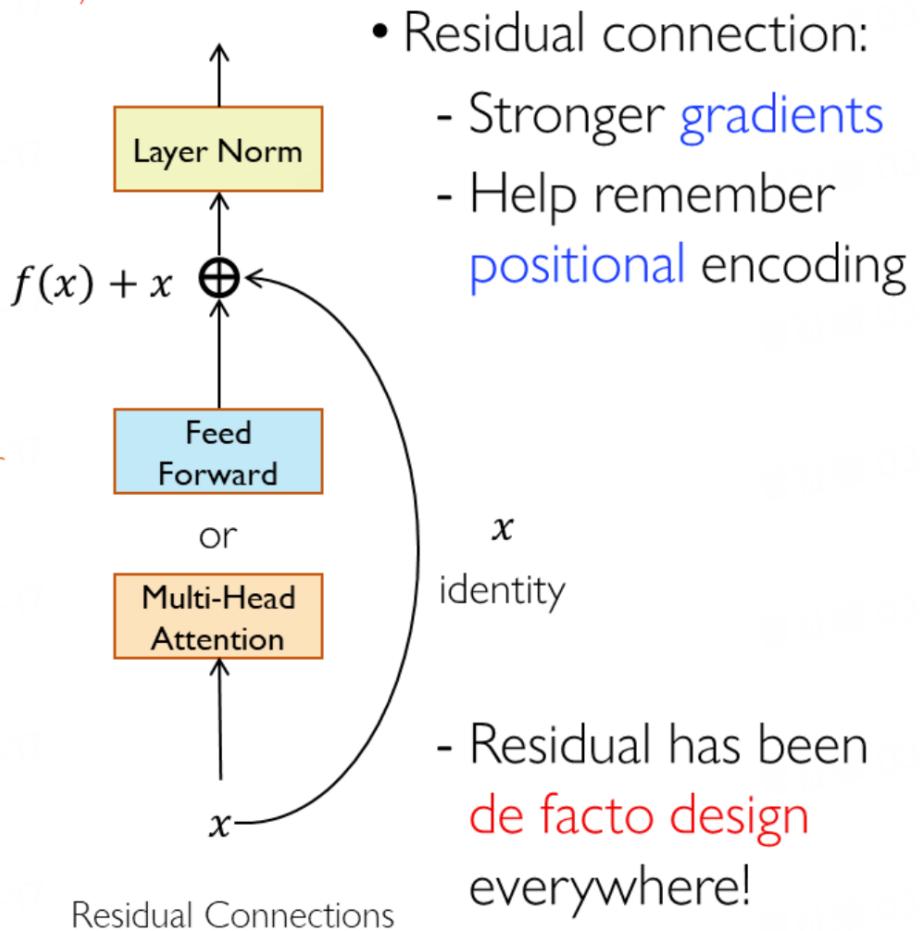
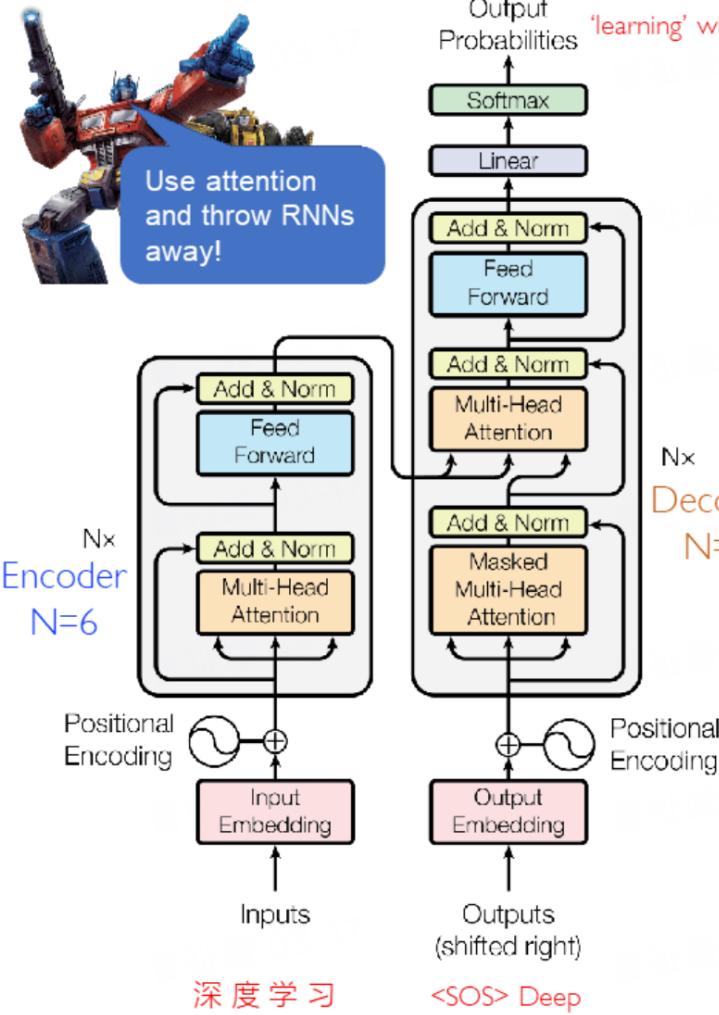


③ Position-wise FFN

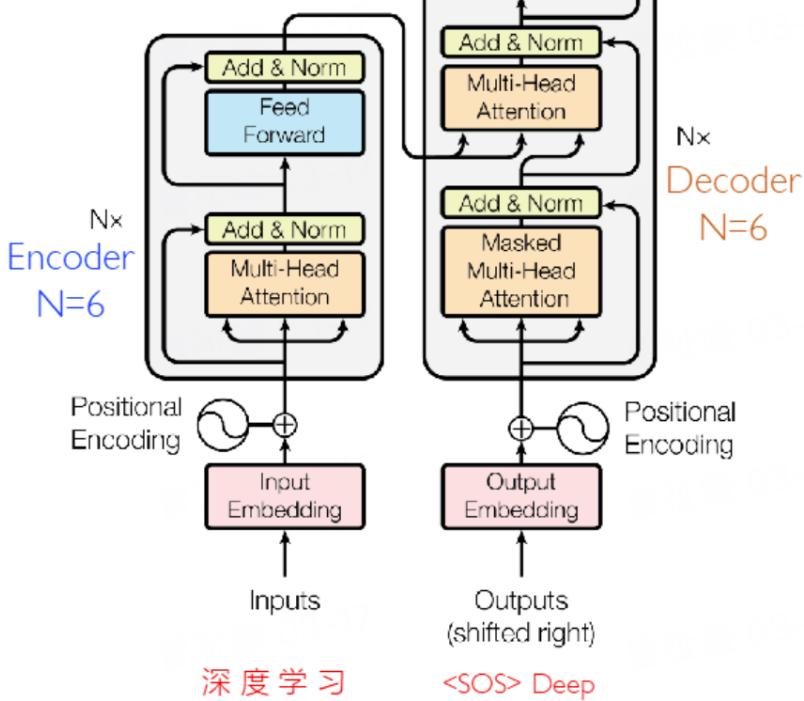
- Apply a two-layer position-wise MLP
 - Consisting of two linear transformations with a ReLU activation in between
$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$
 - Huge number of parameters.
- Equivalent to apply 1D convolution layers
 - Recall 1×1 convolution
 - Parameter are sharing among all positions
 - Applied to sequences with arbitrary length



④ Residual Connection



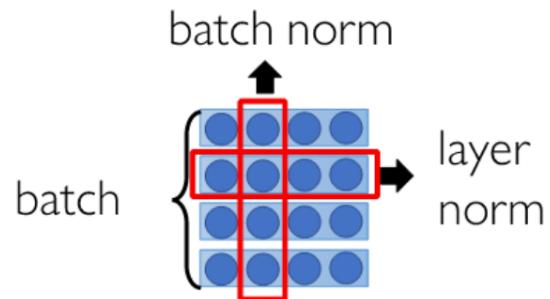
⑤ Layer Normalization



Output Probabilities
'learning' with probability 0.96

- Layer normalization:

- Center embeddings **around origin**, which helps attention layers
- No train-inference mismatch: Same computation throughout

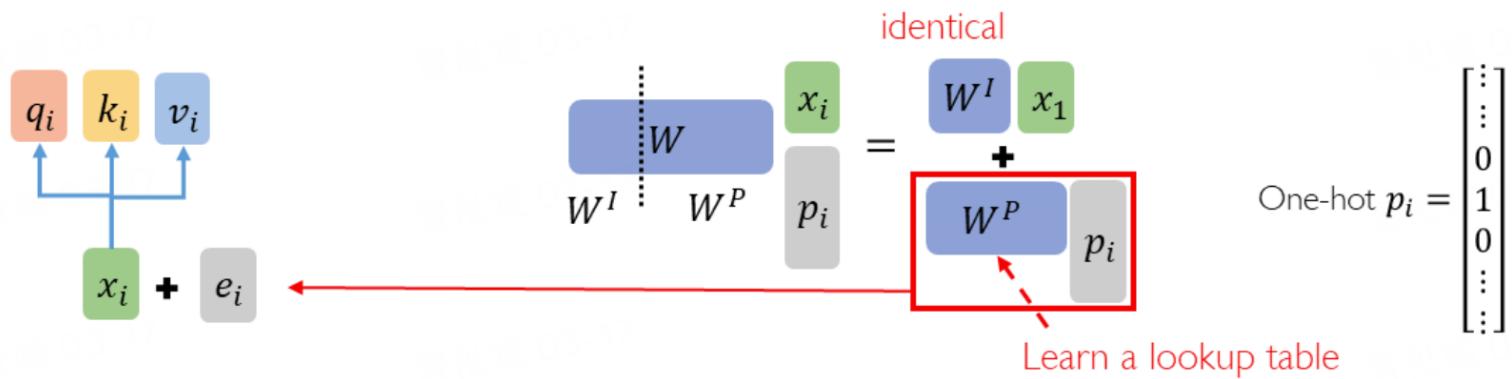


- DeepNet: Scaling Transformers to 1000 Layers. arXiv 2022. 😊

⑥ Positional Encoding

- No position information in self-attention 😞

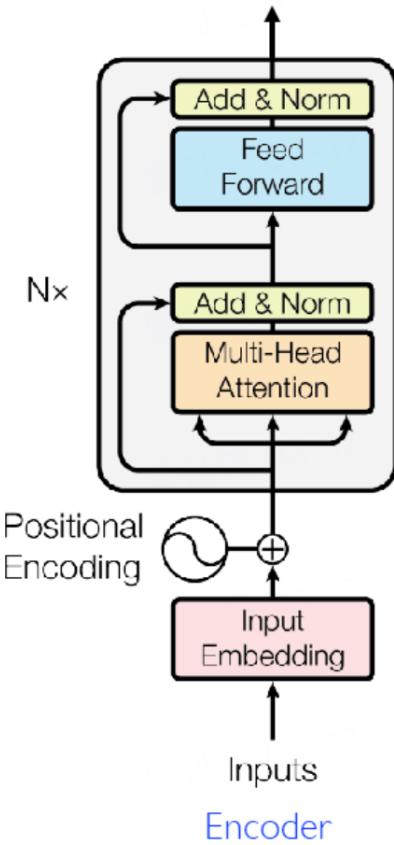
Inject some information about the relative or absolute position of the tokens in the sequence



- Each position i has a unique positional vector e_i yielded by fixed function

$$e_i(2j) = \sin\left(\frac{i}{10000^{\frac{2j}{d_{\text{model}}}}}\right), e_i(2j+1) = \cos\left(\frac{i}{10000^{\frac{2j}{d_{\text{model}}}}}\right)$$

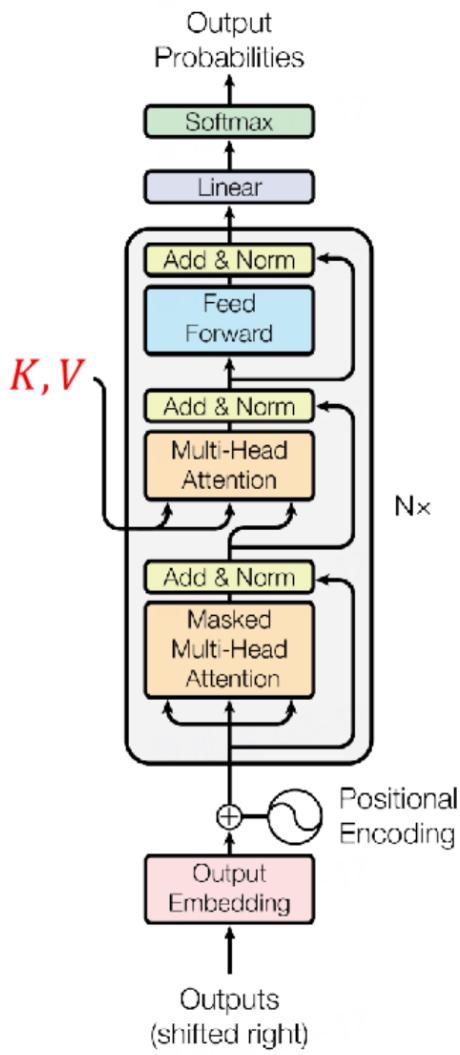
① Encoder



- The encoder stacks N encoder blocks
 - ($N = 6$ in original paper)
 - Multi-head self-attention
 - Position-wise FFN
 - Positional encoding at the bottoms of encoder
- Each block maintains shape
 - (#vectors, vector length)

- Weaknesses:
 - Complexity $O(n^2)$, n is sequence length
 - Difficult to train with huge parameters
 - Many heads are unimportant and can be pruned with little impact on performance

② Decoder



- The decoder stacks N decoder blocks
 - ($N = 6$ in original paper)
 - Multi-head self-attention
 - » K and V in the multi-head self-attention are from the encoder outputs
 - Masked multi-head self-attention
 - Position-wise FFN
 - Positional encoding at the bottoms of decoder
- Each block maintains shape
 - (#vectors, vector length)
- How does decoder work at inference and training?

2 Fourier/Galerkin Transformer

2.1 Attention Mechanism

Then encoder contains a stack of identical simple attention-based layer. For simplicity, we consider a single attention head that maps $\mathbf{y} \in \mathbb{R}^{n \times d}$ to another element in $\mathbb{R}^{n \times d}$ and defines trainable projection matrices/latent representation Q,K,V

$$W^Q, W^K, W^V \in \mathbb{R}^{d \times d}, Q := \mathbf{y}W^Q, K := \mathbf{y}W^K, V := \mathbf{y}W^V$$

Shortly, Galerkin/Fourier Attention uses a mesh-weighted normalization without softmax and allows a scaling to propagate through the encoder

$$\text{Attn}_{\text{sp}}: \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}, \tilde{\mathbf{y}} \leftarrow \mathbf{y} + \text{Attn}_+(\mathbf{y}), \mathbf{y} \rightarrow \tilde{\mathbf{y}} + g(\tilde{\mathbf{y}})$$

$$(Fourier) \quad \mathbf{z} = \text{Attn}_{\text{f}}(\mathbf{y}) := (\tilde{Q} \tilde{K}^\top) V / n$$

$$(Galerkin) \quad \mathbf{z} = \text{Attn}_{\text{g}}(\mathbf{y}) := Q (\tilde{K}^\top \tilde{V}) / n$$

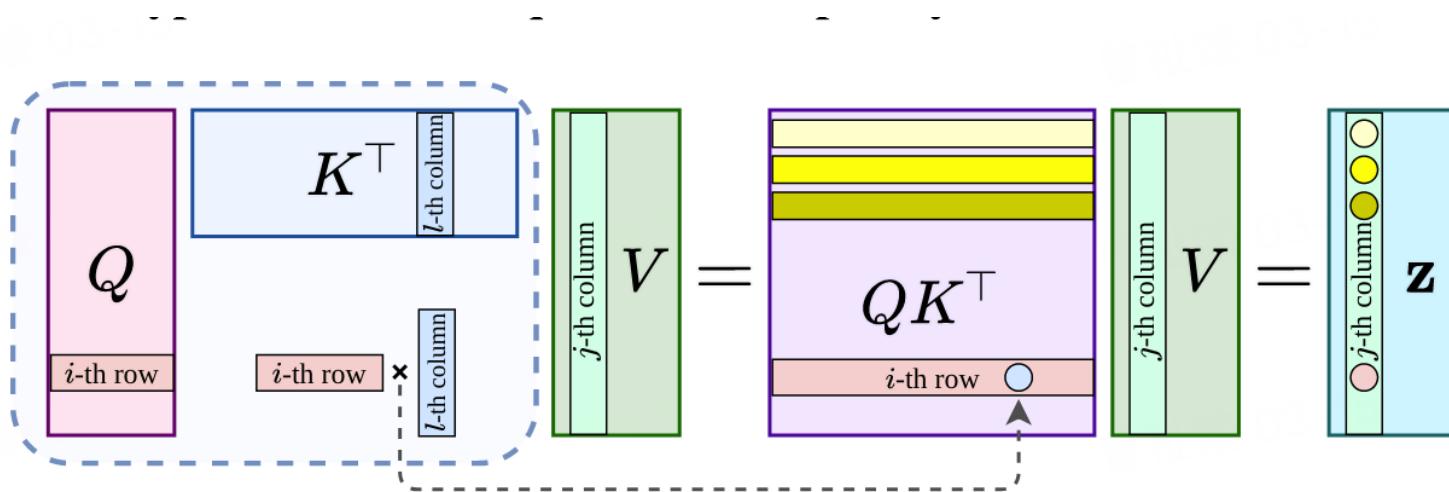
here \sim denotes layernorm and g is standard 2-layer FFN

$$\mathbf{y} \mapsto \text{Ln}(\mathbf{y} + \text{Attn}_+(\mathbf{y}) + g(\text{Ln}(\mathbf{y} + \text{Attn}_+(\mathbf{y})))), \text{ where } \text{Attn}_+(\mathbf{y}) := Q (K^\top V) / n$$

Consider an operator learning problem with domain $\Omega \subset \mathbb{R}$, $\{x_i\}^n$ denotes the set of grid points in the discretized Ω such that the weight $h = \frac{1}{n}$ is the mesh size.

Let $\zeta_q(\cdot), \phi_k(\cdot), \psi_v(\cdot): \Omega \rightarrow \mathbb{R}^{1 \times d}$ denotes the feature map of Q,K,V, i.e the i -th row of Q,K,V written as $\mathbf{q}_i = \zeta_q(x_i), \mathbf{k}_i = \phi_k(x_i), \mathbf{v}_i = \psi_v(x_i)$.

Using $V \in \mathbb{R}^{n \times d}$ with full column rank as example. its column contains a set of bases $\{v_j(\cdot)\}_{j=1}^d$ evaluated at grid point, similarly to $\{q_j(\cdot)\}_{j=1}^d, \{k_j(\cdot)\}_{j=1}^d, \{z_j(\cdot)\}_{j=1}^d$, and we denotes \mathbf{v}^j as the j -th column of V , then $(\mathbf{v}^j)_l = v_j(x_l)$

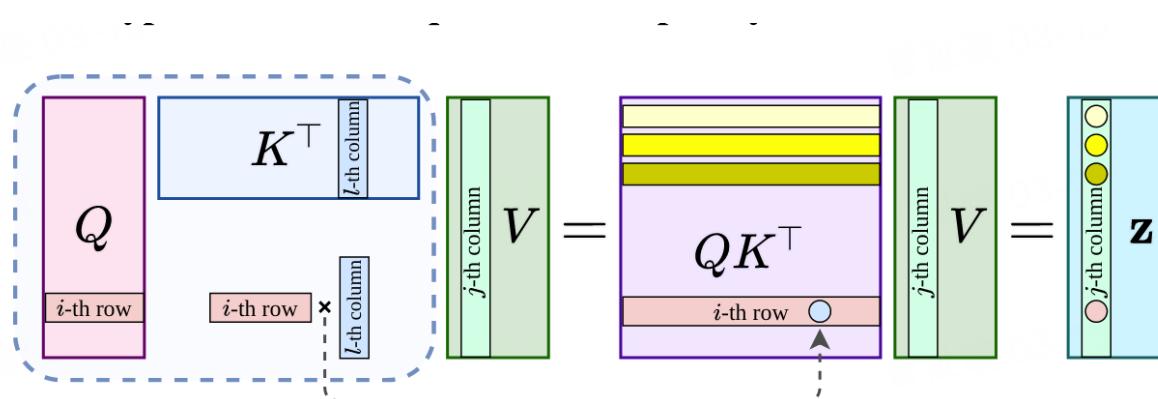


2.2 Fourier-type attention with quadratic complexity

In the Fourier-type attention, Q, K are assumed to be normalized for simplicity, the j -th column in the i -th row of \mathbf{z} is computed by

$$\begin{aligned} (\mathbf{z}_i)_j &= h(QK^\top)_i \cdot \mathbf{v}^j = h(\mathbf{q}_i \cdot \mathbf{k}_1, \dots, \mathbf{q}_i \cdot \mathbf{k}_l, \dots, \mathbf{q}_i \cdot \mathbf{k}_n)^\top \cdot \mathbf{v}^j \\ &= h \sum_{l=1}^n (\mathbf{q}_i \cdot \mathbf{k}_l) (\mathbf{v}^j)_l \approx \int_{\Omega} (\zeta_q(x_i) \cdot \phi_k(\xi)) v_j(\xi) d\xi \end{aligned}$$

where the h - weight facilitates the numerical quadrature interpretation of the inner product. Concatenating columns $1 \leq j \leq d$ yeilds the i - row \mathbf{z}_i of output $\mathbf{z} : \mathbf{z}_i \approx \int_{\Omega} (\zeta_q(x_i) \cdot \phi_k(\xi)) \psi_v(\xi) d\xi$, therefore the output at i -th row computes approximately an integral transform of a non-symmetric learnable kernel $\kappa(x, \xi) = \zeta_q(x)\phi_k(\xi)$ evaluated at x_i .



2.3 Galerkin-type attention with linear complexity

For Galerkin-type attention, K, V are assumed to be normalized, we consider the i -th entry in j -th column of \mathbf{z}

$$(\mathbf{z}^j)_i = h \mathbf{q}_i^\top \cdot (K^\top V)_{\cdot j},$$

which is the inner product of the i -th row of Q and the j -th column of $K^\top V$. Thus,

$$\mathbf{z}^j = h \begin{pmatrix} | & | & | & | \\ \mathbf{q}_1 & \mathbf{q}_2 & \cdots & \mathbf{q}_n \\ | & | & | & | \end{pmatrix}^\top (K^\top V) \cdot j = h \left((K^\top V)_{\cdot j} \begin{pmatrix} - & \mathbf{q}^1 & - \\ - & \vdots & - \end{pmatrix} \right)^\top$$

$(K^\top V) \cdot j$ contains the coefficients for the linear combination of the vector representations $\{\mathbf{q}^l\}_{l=1}^d$ of the bases stored in Q 's column space to form the output \mathbf{z}

$$\mathbf{z}^j = h \sum_{l=1}^d \mathbf{q}^l (K^\top V)_{lj}, \quad \text{where } (K^\top V)_{\cdot j} = (\mathbf{k}^1 \cdot \mathbf{v}^j, \mathbf{k}^2 \cdot \mathbf{v}^j, \dots, \mathbf{k}^d \cdot \mathbf{v}^j)^\top$$

$$z_j(x_i) := (\mathbf{z}^j)_i = h \sum_{l=1}^d (\mathbf{k}^l \cdot \mathbf{v}^j)(\mathbf{q}^l)_i \approx \sum_{l=1}^d \left(\int_{\Omega} v_j(\xi) k_l(\xi) d\xi \right) q_l(x_i)$$

Example 1: $n = 8192$				Encoders only: $n = 8192, d = 128, l = 10$				Computational complexity of the dot-product per layer	
	Mem	CUDA Mem	Speed	GFLOP	Mem	CUDA Mem	Speed	GFLOP	
ST	18.39	31.06	5.02	1393	18.53	31.34	4.12	1876	$O(n^2 c_e d)$
FT	10.05	22.92	6.10	1138	10.80	22.32	5.46	1610	$O(n^2 d)$
LT	2.55	2.31	12.70	606	2.73	2.66	10.98	773	$O(n(d^2 + c_e d))$
GT	2.36	1.93	27.15	275	2.53	2.33	19.20	412	$O(nd^2)$

Table 2: (a) Evaluation relative error ($\times 10^{-3}$) of Burgers' equation 5.1. (b) Evaluation relative error ($\times 10^{-2}$) of Darcy interface problem 5.2.

	(a)			(b)		
	$n = 512$	$n = 2048$	$n = 8192$		$n_f, n_c = 141, 43$	$n_f, n_c = 211, 61$
FNO1d [57]	15.8	14.6	13.9	FNO2d [57]	1.09	1.09
FNO1d 1cycle	4.373	4.126	4.151	FNO2d 1cycle	1.419	1.424
FT regular Ln	1.400	1.477	1.172	FT regular Ln	0.838	0.847
GT regular Ln	2.181	1.512	2.747	GT regular Ln	0.894	0.856
ST regular Ln	1.927	2.307	1.981	ST regular Ln	1.075	1.131
LT regular Ln	1.813	1.770	1.617	LT regular Ln	1.024	1.130
FT Ln on Q, K	1.135	1.123	1.071	FT Ln on Q, K	0.873	0.921
GT Ln on K, V	1.203	1.150	1.025	GT Ln on K, V	0.839	0.844
ST Ln on Q, K	1.271	1.266	1.330	ST Ln on Q, K	0.946	0.959
LT Ln on K, V	1.139	1.149	1.221	LT Ln on K, V	0.875	0.970

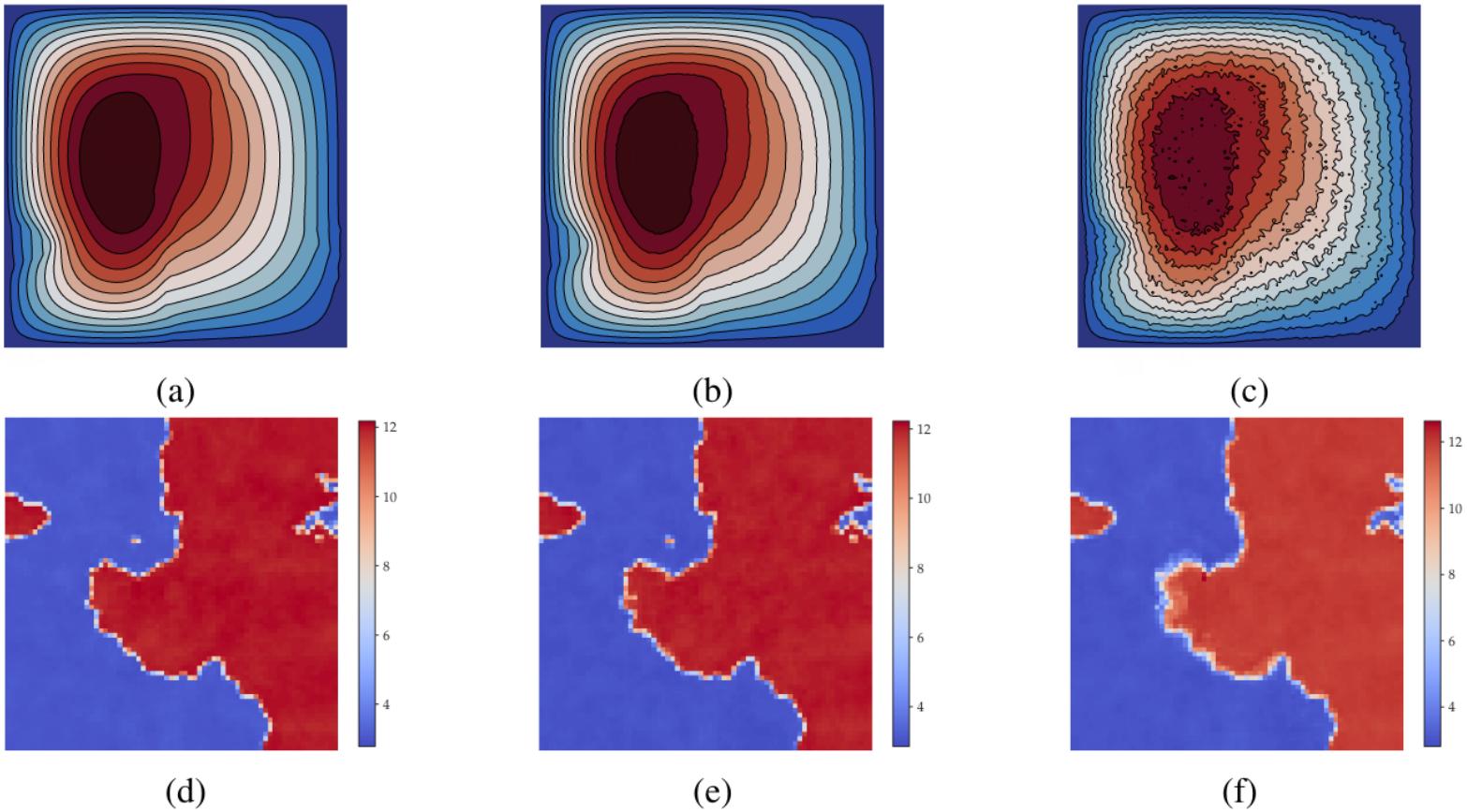


Figure 11: Galerkin transformer evaluation for the inverse interface coefficient identification problem using the same sample with Figure 10, model trained and evaluated under the same amount of noises: (a)–(c) the input $u_h(x)$ with noise level 0, 1%, and 10% on 211×211 grid; (d)–(f) the recovered coefficient through evaluation with noise level 0, 1%, and 10% on 71×71 grid with relative error being 0.0160, 0.0292, and 0.0885, respectively.

3 Operator Transformer

3.1 Cross Attention

Besides the standard attention (Galerkin Type or Fourier Type), the Q, K, V are computed on the same grid of input function. (Consider sound speed to wavefield). When decoding, one need to evaluate the solution on given point $\{y_i\}_{i=1}^m$, one can use cross attention. Cross attention can be viewed as $\{k_l(\cdot), v_l(\cdot), q_l(\cdot)\}_{l=1}^d$, where k, v are defined on the input grid $D(x):\{x_i\}_{i=1}^n$, q is defined on the output grid $D(y):\{y\}_{i=1}^m$

$$z_s(y_j) = \sum_{l=1}^d \frac{\sum_{i=1}^n k_l(x_i) v_s(x_i)}{n} q_l(y_j)$$

we add FFN to improve the nonlinearity of NN.

3.2 Encoder

Input encoder inputs the sampled value $a(x_i)$ of $\{x_i\}_{i=1}^n$. After point-wise MLP, these values are map to latent space $f^{(0)}$, then input encoder use self-attention block to extract feature of input

$$\mathbf{f}^{(l')} = \text{LayerNorm}(\mathbf{f}^{(l)} + \text{Attn}(\mathbf{f}^{(l)})), \quad \mathbf{f}^{(l+1)} = \text{LayerNorm}(\mathbf{f}^{(l')} + \text{FFN}(\mathbf{f}^{(l')}))$$

Query encoder using a point-wise MLP with Random Fourier Projection to extract feature of query point Y

$$\gamma(Y) = [\cos(2\pi YB), \sin(2\pi YB)]$$

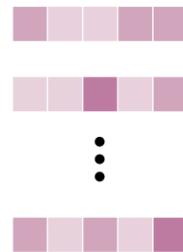
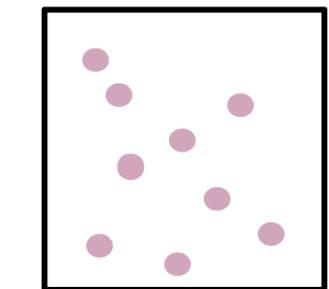
where $Y = [y_1, \dots, y_n]^T$ are coordinates and $B \in \mathbb{R}^{d_1 \times d_2}$ sampled from $\mathcal{N}(0, \sigma^2)$

Cross Attention module combines the encoding output $f^{(L)}$ with the query encoder output z^0 to obtain the final output

$$\mathbf{z}' = \mathbf{z}^{(0)} + \text{Cross-Attn}(\mathbf{z}^{(0)}, \mathbf{f}^{(L)}), \quad \mathbf{z} = \mathbf{z}' + \text{FFN}(\mathbf{z}')$$

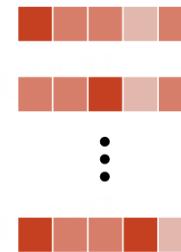
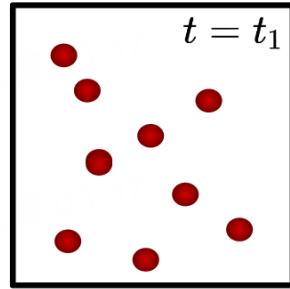
Also to encode the spatial information Operator Transformer use RoPE.

Steady state system

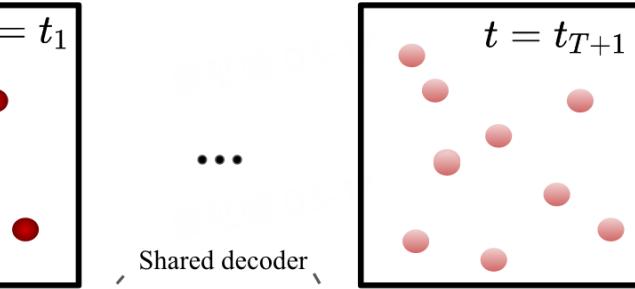


Latent encoding

Time-dependent system

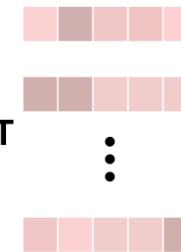
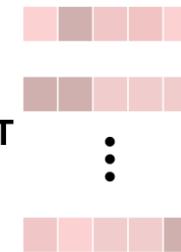


Shared propagator



...

Shared decoder



Data settings		Relative L_2 norm				
Case	ν, T	FNO-2D	FNO-3D	MWT	G.T.*	OFormer
NS1	$1 \times 10^{-3}, 50$	0.0128	0.0086	0.0062	0.0098	0.0104 ± 0.0005
NS2-part	$1 \times 10^{-4}, 30$	0.1559	0.1918	0.1518	0.1399	0.1755 ± 0.0059
NS2-full	$1 \times 10^{-4}, 30$	0.0834	0.0820	0.0667	0.0792	0.0645 ± 0.0011
NS3	$1 \times 10^{-5}, 20$	0.1556	0.1893	0.1541	0.1340	0.1705 ± 0.0007
NS-mix*	$[0.4, 1] \times 10^{-4}, 30$	0.1650	0.1654	0.1267	0.1462	0.1402 ± 0.0016
# of parameters (M)		0.41/2.37	6.56	179.18	1.56	1.85

Table 1: Benchmark on 2D Navier-Stokes equation with fixed 64×64 grid. For large dataset (NS-mix) we use a larger version of FNO-2D by increasing modes and width. We adopt Galerkin Transformer (G.T.) with a larger hidden dimension ($64 \rightarrow 96$) to match the size of other models.

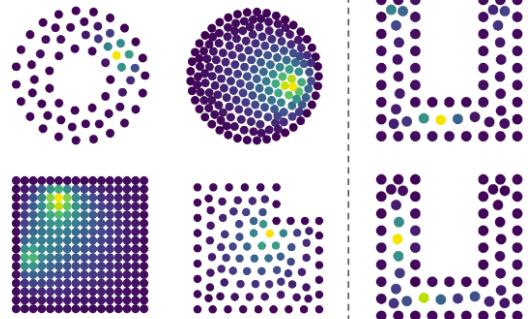
Res.	Relative L_2 norm ($\times 10^{-3}$)				
	FNO*	FNO+*	MWT	F.T.	OFormer
512	3.15	0.72	1.85	1.14	1.42 ± 0.03
2048	3.07	0.72	1.86	1.12	1.30 ± 0.07
8192	2.76	0.70	1.78	1.07	1.54 ± 0.10

Table 2: Benchmark on 1D Burgers' equation with different resolution. FNO uses ReLU (Nair & Hinton, 2010) whereas FNO+ uses GELU. Both FNO variants remove normalization compared to the original paper.

Res.	Relative L_2 norm ($\times 10^{-2}$)		
	FNO	G.T.	OFormer
141	1.09	0.84	1.26 ± 0.03
211	1.09	0.84	1.28 ± 0.02

Table 3: Benchmark on 2D Darcy flow with different resolution. MWT's result is not reported since its implementation requires resolution to be power of 2 and thus it is trained on a different data.

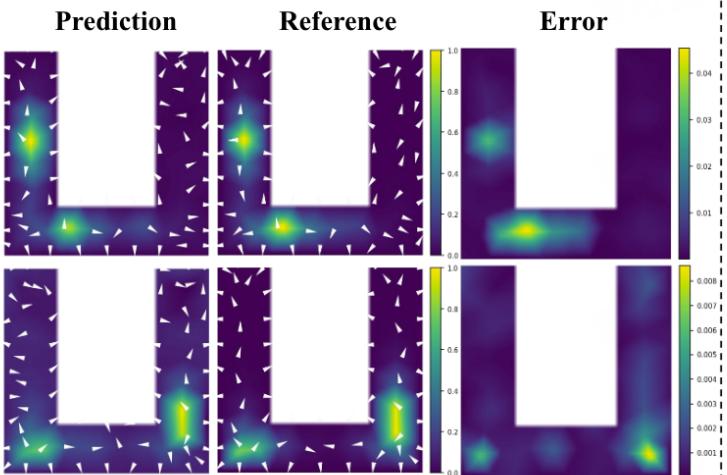
Training samples



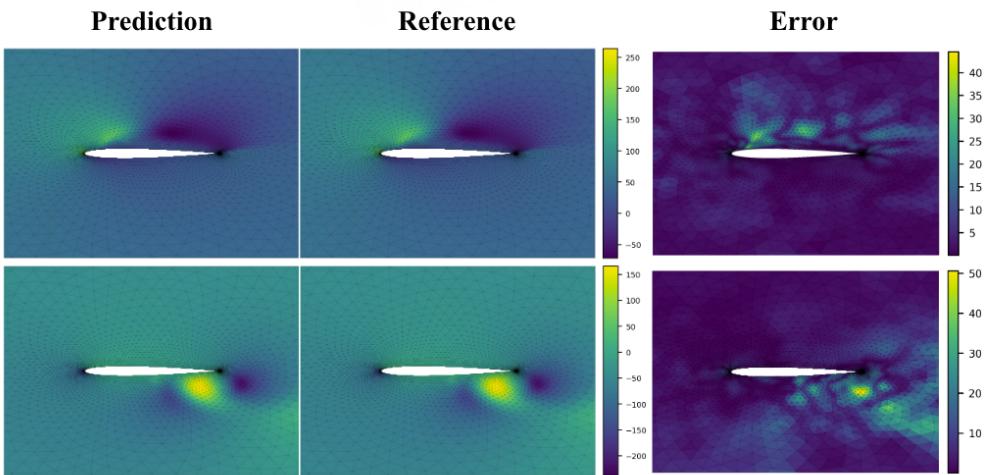
Testing samples

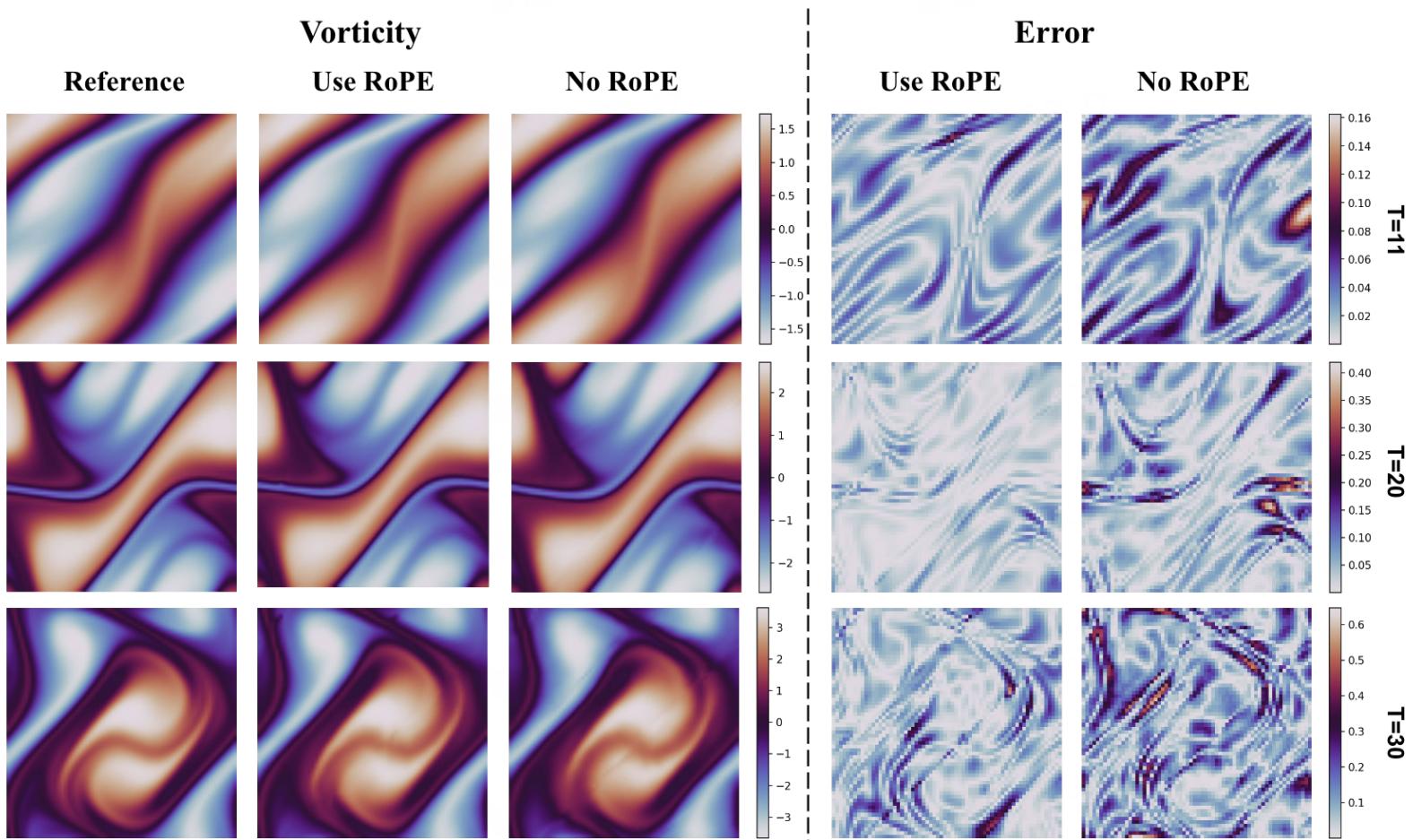
Mesh aug.	Model	Mean squared error ($\times 10^{-3}$)			
		Electric		Magnetic	
		Potential	Field	Potential	Field
No	OFormer	0.174 \pm 0.083	1.690 \pm 0.201	0.138 \pm 0.021	1.470 \pm 0.172
	GNN	0.327	2.723	0.161	2.640
Yes	OFormer	0.019 \pm 0.004	0.328 \pm 0.025	0.016 \pm 0.002	0.195 \pm 0.009
	GNN	0.006	1.821	0.006	1.338

(a) Electric field



(b) Airfoil





4 HIERARCHICAL TRANSFORMER

4.1 Multi-scale PDE

Multiscale PDEs, in a narrower sense, refer to PDEs with rapidly varying coefficients

Multiscale elliptic PDE:

$$\begin{aligned} -\nabla \cdot (a(x) \nabla u(x)) &= f(x) \quad x \in D \\ u(x) &= 0 \quad x \in \partial D \end{aligned}$$

where $0 < a_{\min} \leq a(x) \leq a_{\max}$, the solution map is $\mathcal{S}: L^\infty(D; \mathbb{R}_+) \rightarrow H_0^1(D; \mathbb{R}_+)$ such that $u = \mathcal{S}(a)$. The setup $a(x) \in L^\infty$ allows rough coefficients with fast oscillation $a(x) = a(x/\varepsilon)$ and high contrast ratio $\frac{a_{\max}}{a_{\min}} \gg 1$.

Navier-Stokes Equation becomes turbulent due to the simultaneous interaction of a wide range of temporal and spatial scales of motion.

Numerical solution of Helmholtz equation exhibits severe difficulties in the high wave number regime due to the interaction of high-frequency waves and numerical mesh.

4.2 Hierachichal Attention Model

4.2.1 Hierarchical Discretization

Let $\mathcal{I}^{(r)}$ be the finest level index set, such that each index $i = (i_1, \dots, i_r) \in \mathcal{I}^{(r)}$ denotes the finest level spatial objects such as image pixels.

For any $i = (i_1, \dots, i_r) \in \mathcal{I}^{(r)}$ and $1 \leq m \leq r$, $i^{(m)} = (i_1, \dots, i_m)$ represents i 's m -th level parent node which is the aggregate of finer level objects. $\mathcal{I}^{(m)} := \{i^{(m)} : i \in \mathcal{I}^{(r)}\}$ is the m -th level index set.

Each index $i \in \mathcal{I}^{(m)}$ corresponds to a token, which can be characterized e.g. by its position $\mathbf{x}_i^{(m)}$ and a feature vector $\mathbf{f}_i^{(m)}$ with number of channels $C^{(m)}$.

4.2.2 Reduce Operation

The reduce operation defines the map from finer-level features to coarser-level features.

The reduce map can be abstractly defined as $\mathbf{f}_i^{(m),t} = \mathcal{R}^{(m)}(\{\mathbf{f}_j^{(m+1),t}\}_{j \in i^{(m,m+1)}})$, which maps the $(m+1)$ -th level features with indices in $i^{(m,m+1)}$ to the m -th level feature with index $i^{(m)}$. For simplicity, we use a simple linear layer for the operator $\mathcal{R}^{(m)}$ in our current implementation, namely, $\mathbf{f}_i^{(m),t} = R_0^{(m)}\mathbf{f}_{(i,0)}^{(m+1),t} + R_1^{(m)}\mathbf{f}_{(i,1)}^{(m+1),t} + R_2^{(m)}\mathbf{f}_{(i,2)}^{(m+1),t} + R_3^{(m)}\mathbf{f}_{(i,3)}^{(m+1),t}$, where $R_0^{(m)}, R_1^{(m)}, R_2^{(m)}, R_3^{(m)} \in \mathbb{R}^{C^{(m-1)} \times C^{(m)}}$ are parametrized by linear layers.

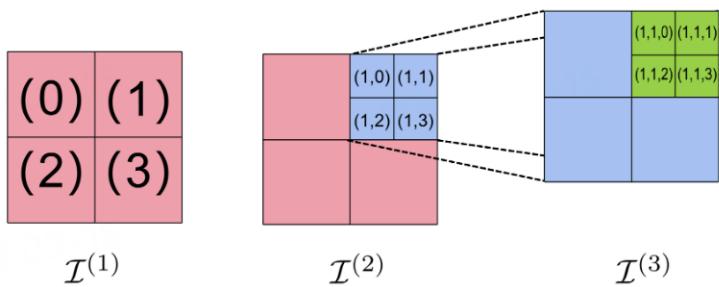


Figure 3.1: Hierarchical discretization and index tree. In this example, the 2D unit square is hierarchically discretized into three levels which are indexed by $\mathcal{I}^{(1)}$, $\mathcal{I}^{(2)}$ and $\mathcal{I}^{(3)}$, respectively. For example, we denote by $(1)^{(1,2)} = \{(1, 0), (1, 1), (1, 2), (1, 3)\}$ the set of the second level child nodes of the node (1).

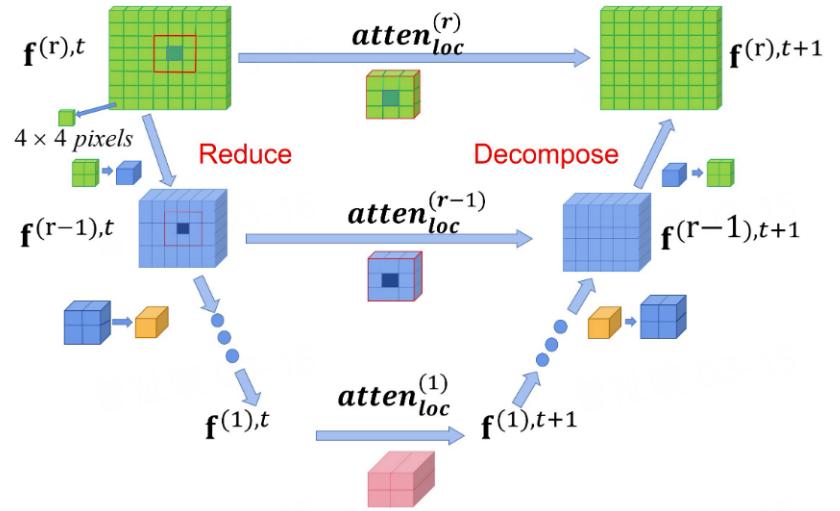


Figure 3.2: One V-cycle of the feature update.

4.2.3 Multilevel Local Aggregation

The finest level features $f_i^{(r),t} \in \mathbb{R}^{C^{(r)}}$ at the evolution step t , with index $i \in \mathcal{I}^{(r)}$, can be updated by the following token aggregation formula through vanilla attention,

$$\text{atten: } f_i^{(r),t+1} = \sum_{j=1}^{N^{(r)}} \mathcal{G}(f_i^{(r),t}, f_j^{(r),t}) v(f_j^{(r),t}), \quad \text{for } i \in \mathcal{T}^{(r)}, \quad (1)$$

with $N^{(r)} := |\mathcal{I}^{(r)}|$. For simplicity, we ignore the softmax normalizing factor and only consider single-head attention, and we assume that the interaction kernel \mathcal{G} is of the form $\mathcal{G}(f_i, f_j) = \exp(W^Q f_i \cdot W^K f_j) = \exp(q_i \cdot k_j)$, where $q_i := W^Q f_i$, $k_i := W^K f_i$, $v_i := v(f_i) = W^V f_i$, and $W^Q, W^K, W^V \in \mathbb{R}^{C^{(r)} \times C^{(r)}}$ are learnable matrices.

Instead of original attention, we proposed a local aggregation scheme. The local aggregation at the r -th level writes,

$$\text{atten}_{\text{loc}}^{(r)}: f_i^{(r),t+1} = \sum_{j \in \mathcal{N}^{(r)}(i)} \exp(q_i^{(r),t} \cdot k_j^{(r),t}) v_j^{(r),t}, \quad \text{for } i \in \mathcal{I}^{(r)}, \quad (2)$$

where $\mathcal{N}^{(r)}(i)$ is the set of the r -th level neighbors of $i \in \mathcal{I}^{(r)}$, $q_i^{(r),t} := \tilde{W}^Q f_i^{(r),t}$, $k_i^{(r),t} := \tilde{W}^K f_i^{(r),t}$, $v_i^{(r),t} := \tilde{W}^V f_i^{(r),t}$ with learnable matrices $\tilde{W}^Q, \tilde{W}^K, \tilde{W}^V \in \mathbb{R}^{C^{(r)} \times C^{(r)}}$. Then for $m = r - 1, \dots, 1$, we calculate the local attention in each level with nested $q_i^{(m),t}$, $k_i^{(m),t}$, $v_i^{(m),t}$

For the m -th level:

$$\text{atten}_{\text{loc}}^{(m)} : f_i^{(m),t+1} = \sum_{j \in \mathcal{N}^{(m)}(i)} \exp(q_i^{(m),t} \cdot k_j^{(m),t}) v_j^{(m),t}, \quad \text{for } i \in \mathcal{I}^{(m)}, \quad (3)$$

where $q_i^{(m),t} = \mathcal{R}^{(m)}(\{q_j^{(m+1),t}\}_{j \in i^{(m,m+1)}})$, $k_i^{(m),t} = \mathcal{R}^{(m)}(\{k_j^{(m+1),t}\}_{j \in i^{(m,m+1)}})$, $v_i^{(m),t} = \mathcal{R}^{(m)}(\{v_j^{(m+1),t}\}_{j \in i^{(m,m+1)}})$, and $\mathcal{N}^{(m)}(i)$ is the set of the m -th level neighbors of $i \in \mathcal{I}^{(m)}$.

4.2.4 Decompose Operation

To mix the multilevel features $f_i^{(m),t+1}$ for $m = 1, \dots, r$, we propose a decompose operation that reverses the reduce operation from level 1 to level $r - 1$.

The decompose operator

$$D^{(m)}: f_i^{(m),t+1} \mapsto \left\{ f_j^{(m+1),t+\frac{1}{2}} \right\}_{j \in i^{(m,m+1)}},$$

maps the m -th level feature with index i to $(m + 1)$ -th level features associated to its child set $i^{(m,m+1)}$. The $(m + 1)$ -th level features $f_i^{(m+1),t+\frac{1}{2}}$ are further aggregated to $f_i^{(m+1),t+1}$, such that $f_i^{(m+1),t+1} = \sum f_j^{(m+1),t+\frac{1}{2}}$ for $j \in i^{(m,m+1)}$. In the current implementation, we use a simple linear layer such that $f_{i,s}^{(m+1),t+\frac{1}{2}} = D^{(m),s} T f_i^{(m),t+1}$, for $s = 0, 1, 2, 3$, with parameter matrices $D^{(m),s} \in \mathbb{R}^{C^{(m)} \times C^{(m+1)}}$.

4.3 Algorithm

Algorithm 1 One V-cycle of Hierarchical Attention

Input: $\mathcal{I}^{(r)}, \mathbf{f}_i^{(r),t}$ for $i \in \mathcal{I}^{(r)}$.

STEP 0: Get the $\mathbf{q}_i^{(r),t}, \mathbf{k}_i^{(r),t}, \mathbf{v}_i^{(r),t}$ for $i \in \mathcal{I}^{(r)}$.

STEP 1: **For** $m = r - 1, \dots, 1$, **Do** the reduce operations $\mathbf{q}_i^{(m),t} = \mathcal{R}^{(m)}(\{\mathbf{q}_j^{(m+1),t}\}_{j \in i^{(m,m+1)}})$ and also for $\mathbf{k}_i^{(m),t}$ and $\mathbf{v}_i^{(m),t}$, for any $i \in \mathcal{I}^{(m)}$.

STEP 2: **For** $m = r, \dots, 1$, **Do** the local aggregation by equation 3.2 and equation 3.3 to get $\mathbf{f}_i^{(m),t+1}, m = 1, \dots, r$ for any $i \in \mathcal{I}^{(m)}$.

STEP 3: **For** $m = 1, \dots, r - 1$, **Do** the decompose operations $\{\mathbf{f}_j^{(m+1),t+\frac{1}{2}}\}_{j \in i^{(m,m+1)}} = \mathcal{D}^{(m)}(\mathbf{f}_i^{(m),t+1})$, for any $i \in \mathcal{I}^{(m)}$; then $\mathbf{f}_i^{(m+1),t+1} = \mathbf{f}_i^{(m+1),t+\frac{1}{2}}$, for any $i \in \mathcal{I}^{(m+1)}$

Output: $\mathbf{f}_i^{(r),t+1}$ for any $i \in \mathcal{I}^{(r)}$.

The decoder is chosen as spectral convolution layer.

Empirical L^2 loss function is defined as

$$L^2(\{(a_j, u_j)\}_{j=1}^N; \theta) := \frac{1}{N} \sum_{i=1}^N \frac{\|u_j - \mathcal{N}(a_j; \theta)\|_2^2}{\|u_j\|_2^2},$$

where $\|\cdot\|_2$ is the canonical l^2 vector norm. For any $\xi \in \mathbb{Z}_n^2 := \{\xi \in \mathbb{Z}^2 \mid -n/2 + 1 \leq \xi_j \leq n/2, j = 1, 2\}$, the normalized discrete Fourier transform (DFT) coefficients of f writes

$$F(f)(\xi) := \frac{1}{\sqrt{n}} \sum_{x \in G^2} f(x) e^{-2\pi i x \cdot \xi}.$$

The empirical H^1 loss function is given by,

$$L^H(\{(a_j, u_j)\}_{j=1}^N; \theta) := \frac{1}{N} \sum_{i=1}^N \frac{\|u_j - \mathcal{N}(a_j; \theta)\|_h}{\|u_j\|_h},$$

where $\|u\|_h := \sqrt{\sum_{\xi \in \mathbb{Z}_n^2} |\xi|^2 |F(u)(\xi)|^2}$

Model	Runtime (s)	Darcy smooth		Darcy rough		Multiscale trigonometric	
		L^2	H^1	L^2	H^1	L^2	H^1
FNO2D	7.278	0.620	3.883	1.646 ± 0.021	11.955 ± 0.088	1.794	12.605
FNO2D-48	8.062	0.619	2.620	1.220 ± 0.018	5.138 ± 0.093	1.565	11.093
FNO2D-96	10.969	0.575	2.437	1.216 ± 0.024	5.140 ± 0.281	1.518	10.106
MWT	19.715	—	—	1.138 ± 0.010	4.107 ± 0.008	1.021	7.245
GT	38.219	0.945	3.365	1.790 ± 0.012	6.269 ± 0.418	1.052	8.207
SWIN	41.417	—	—	1.622 ± 0.047	6.796 ± 0.359	1.489	13.385
HT-NET	33.375	0.291	0.815	0.571 ± 0.001	1.371 ± 0.001	0.603	2.633

— MWT (Gupta et al., 2021) only supports resolution with powers of two.

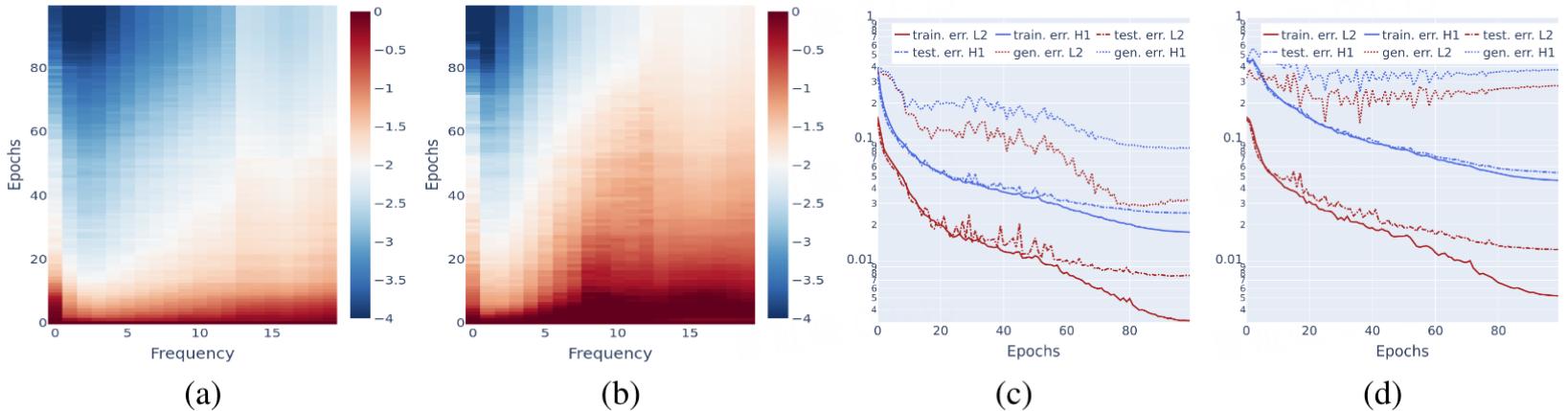


Figure 4.2: In (a) HT-Net trained with H^1 loss, and (b) HT-Net trained with L^2 loss, we show the evolution of errors with x-axis for frequency, y-axis for training epochs, and colorbar for the magnitude of L^2 error on each frequency in \log_{10} scale, the error for each frequency is normalized frequency-wise by the error at epoch 0. The loss curves with training, testing, and generalization errors are shown in (c) for HT-Net trained with H^1 loss, and in (d) for HT-Net trained with L^2 loss.

5 Transformer For BVP

Mathematical Formulation of EIT: Identify the unknown inclusion D buried in Ω or equivalently $p = \mathcal{I}^D(x)$

Forward Model of EIT:

$$\nabla \cdot (\sigma \nabla u) = 0 \text{ in } \Omega, \quad \sigma = \sigma_1 \text{ in } D, \sigma = \sigma_2 \text{ in } D^c$$

Current: $g = \sigma \nabla u \cdot \mathbf{n}|_{\partial\Omega}$ Neumann B.C

Voltages: $f = u|_{\partial\Omega}$ Dirichlet B.C

NtD Mapping

$$\Lambda_\sigma: g = \sigma \nabla u \cdot \mathbf{n}|_{\partial\Omega} \rightarrow f = u|_{\partial\Omega}$$

Forward Model $\mathcal{F}: \mathbb{P} \rightarrow \mathbb{M}$

$$\mathbf{m} = \mathcal{F}(\mathbf{p})$$

where $\mathbf{m} \in \mathbb{M}$ is the measurement of $\partial\Omega$, $\mathbf{p} \in \mathbb{P}$ is the coefficient to be recovered.

Inverse problem $\mathcal{F}^{(-1)}: \mathbb{M} \rightarrow \mathbb{P}$, $\mathbf{p} = \mathcal{F}^{(-1)}(\mathbf{m})$, essentially

$$\Lambda_\sigma \rightarrow \sigma$$

5.1 Finite Number of Data Pair

One measurement is a single data pair (g, f) , and all the measurements form:

$$\mathbb{M} := H^{-1/2}(\partial\Omega) \times \Lambda_\sigma H^{-1/2}(\partial\Omega) = \{(g, f) : g \in H^{-1/2}(\partial\Omega), f = \Lambda_\sigma g\}$$

Infinite data pair: $\mathbb{M} = \{(g, f) : g \in H^{-1/2}(\partial\Omega), f = \Lambda_\sigma g\}$,

Finite data pair : $\mathbb{M}_L = \text{span}\{\{(g_l, f_l)\}_{l=1}^L : g_l \in H^{-1/2}(\partial\Omega), f_l = \Lambda_\sigma g_l\}$

Operator

$$\mathcal{F}_L^{-1} : \mathbb{M}_L \rightarrow \mathbb{P}$$

may not be well-defined.

5.2 Direct Sampling Method

The index function:

$$\mathcal{I}^D(x) \approx \mathcal{I}_1^D(x) := \frac{\mathbf{d}(x) \cdot \nabla \phi(x)}{|f - \Lambda_{\sigma_0} g|_{L^2(\partial\Omega)} |\eta_x|_Y}.$$

$$f - \Lambda_{\sigma_0} g \rightarrow \phi \rightarrow \mathbf{d} \rightarrow \eta_x$$

1. $f - \Lambda_{\sigma_0} g = (\Lambda_\sigma - \Lambda_{\sigma_0})g$: the difference of NtD mapping (known data)
2. ϕ is the harmonic extension of $f - \Lambda_{\sigma_0} g$ with noise ξ

$$-\Delta \phi = 0 \quad \text{in } \Omega, \quad \mathbf{n} \cdot \nabla \phi = (f - \Lambda_{\sigma_0} g) + \xi \quad \text{on } \partial\Omega, \quad \int_{\partial\Omega} \phi \, ds = 0$$

3. $\mathbf{d}(x)$ probing direction : empirical choice $\mathbf{d}(x) = \frac{\nabla \phi}{\|\nabla \phi(x)\|}$
4. $-\Delta \eta_x = -\mathbf{d}(x) \cdot \nabla \delta_x$ in Ω , $\mathbf{n} \cdot \nabla \eta_x = 0$ on $\partial\Omega$, $\int_{\partial\Omega} \eta_x \, ds = 0$ where $\{\eta_x\}$ restricted on boundary is equipped with $|\cdot|_Y$ semi-norm

5.3 Integral Form of DSM

$$\mathcal{I}_1^D(x) := \frac{\mathbf{d}(x) \cdot \nabla \phi(x)}{|f - \Lambda_{\sigma_0} g|_{L^2(\partial\Omega)} |\eta_x|_Y}$$

Heuristics: the global information of ϕ should be used to locate a point x

$$\hat{\mathcal{I}}_1^D(x) := \frac{\int_{\Omega} \mathbf{d}(x) \cdot \mathcal{K}(x, y) \nabla \phi(y) dy}{\|f - \Lambda_{\sigma_0} g\|_{L^2(\partial\Omega)} |\eta_x|_Y},$$

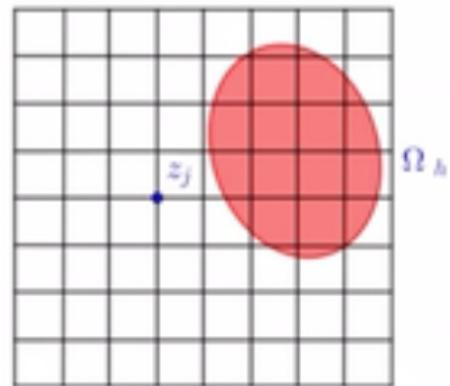
$$d(x) := \int_{\Omega} \mathcal{Q}(x, y) \nabla \phi(y) dy,$$

$$|\eta_x|_Y^2 := (\mathcal{V}\eta_x, \eta_x)_{L^2(\partial\Omega)}.$$

Question: How to choose $\mathcal{K}(x, y)$, $\mathcal{Q}(x, y)$, $\mathcal{V}(x, y)$?

5.4 Matrix Form: Spatial Discretization

- Ω_h : a mesh of Ω with size h
- $\{z_j\}_{j=1}^M$ grid points
- $u \approx u_h \in \mathbb{R}^M$
- ϕ_n : interpolation $\partial_{x_n} \phi$ at $\{z_j\}$, $n = 1, 2$



$$\mathcal{I}_1^D(x) \!:=\! \frac{d(x)\cdot\nabla\,\phi(x)}{|f-\Lambda_{\sigma_0}g|_{L^2(\partial\Omega)}|\eta_x|_Y},$$

$$\hat{\mathcal{I}}_1^D(x) \!:=\! \frac{\int_\Omega \boldsymbol{d}(x)\cdot \mathcal{K}(x,y)\,\nabla\,\phi(y){\rm d}y}{|f-\Lambda_{\sigma_0}g|_{L^2(\partial\Omega)}|\eta_x|_Y},$$

$$d(x) := \int_\Omega \mathcal{Q}(x,y)\,\nabla\,\phi(y){\rm d}y,$$

$$|\eta_x|_Y^2 \!:=\! (\mathcal{V}\,\eta_x,\,\eta_x)_{L^2(\partial\Omega)},$$

$$\int_\Omega \mathcal{K}(z_i,y)\,\partial_{x_n}\phi(y){\rm d}y \approx \sum_j \omega_j \mathcal{K}(z_i,z_j)\,\partial_{x_n}\phi(z_j) =: \boldsymbol{k}_i^\top\,\boldsymbol{\phi}_n.$$

$$d_n(z_i) \approx \sum_j \omega_j \,\mathcal{Q}(z_i,z_j)\,\partial_{x_n}\phi(z_j) =: \boldsymbol{q}_i^\top\,\boldsymbol{\phi}_n.$$

$$|\eta_{z_i}|_Y^2 \approx \bar{\boldsymbol{\eta}}_i^\top\,V\bar{\boldsymbol{\eta}}_i \approx \sum_n c_{n,i}\,\boldsymbol{\phi}_n^\top\,\boldsymbol{q}\,\boldsymbol{q}_i^\top\,\boldsymbol{\phi}_n = \sum_n \boldsymbol{\phi}_n^\top\,\boldsymbol{v}_{n,i}\,\boldsymbol{v}_{n,i}^\top\,\boldsymbol{\phi}_n$$

$$\hat{\mathcal{I}}_1^D(z_i) \approx C_{f,g} \bigg(\sum_n \boldsymbol{\phi}_n^\top\,\boldsymbol{v}_{n,i}\,\boldsymbol{v}_{n,i}^\top\,\boldsymbol{\phi}_n \bigg)^{-1/2} \sum_n \boldsymbol{\phi}_n^\top\,\boldsymbol{q}_i\,\boldsymbol{k}_i^\top\,\boldsymbol{\phi}_n,$$

$$\text{where }C_{f,g}=\|f-\Lambda_{\sigma_0}g\|_{L^2(\partial\Omega)}^{-1}\text{ is a constant.}$$

$$\begin{aligned}\mathcal{I}_1^D(x) &:= \frac{d(x) \cdot \nabla \phi(x)}{|f - \Lambda_{\sigma_0} g|_{L^2(\partial\Omega)} |\eta_x|_Y}, \\ \hat{\mathcal{I}}_1^D(x) &:= \frac{\int_{\Omega} \mathbf{d}(x) \cdot \mathcal{K}(x, y) \nabla \phi(y) dy}{|f - \Lambda_{\sigma_0} g|_{L^2(\partial\Omega)} |\eta_x|_Y}, \\ d(x) &:= \int_{\Omega} \mathcal{Q}(x, y) \nabla \phi(y) dy, \\ |\eta_x|_Y^2 &:= (\mathcal{V}\eta_x, \eta_x)_{L^2(\partial\Omega)},\end{aligned}$$

$$[\hat{\mathcal{I}}_1^D(z_i)]_{i=1}^M \approx C_{f,g} \left(\phi_h W^Q * \phi_h W^K \right) / \left(\phi_h W^V * \phi_h W^V \right)^{1/2}$$

How to choose $\mathcal{K}, \mathcal{Q}, \mathcal{V}$: Learn matrices W^Q, W^K, W^V using Transformer!

5.5 Generate Data

- Spatial discretization: Ω_h : a mesh of Ω , $\{z_j\}_{j=1}^M$: \mathcal{M} grid points, $u \approx u_h \in \mathbb{R}^M$
- Sampling of Coefficient: σ . Generate $N \gg 1$ samples of D with different shape and location
 - p_h^k : a discretization of \mathcal{I}^D on Ω_h
 - m_h^k : a data function associated with $\mathcal{F}_L(p_h^k)$
- Sampling of NtD mapping: For each sample of D , generate L input-output data pair $\{(g_l, f_l)_{l=1}^L\}$ to sample NtD mapping Λ_σ which are used to generated L channels in an input image.

5.6 Operator Learning Problems for EIT

Find a parametrization mapping \mathcal{G}_θ to approximate \mathcal{F}^{-1} by minimizing :

$$\mathcal{J}(\theta) = \frac{1}{N} \sum_{k=1}^N \| \mathcal{G}_\theta(\mathbf{m}_h^k) - \mathbf{p}_h^k \|_{\mathbb{P}}^2$$

The parameters N, h, L will affect the approximation

1. h determines the resolution
2. N affect the representativity of training dataset
3. L decides how much of a finite portion of the infinite spectral information can be accessed

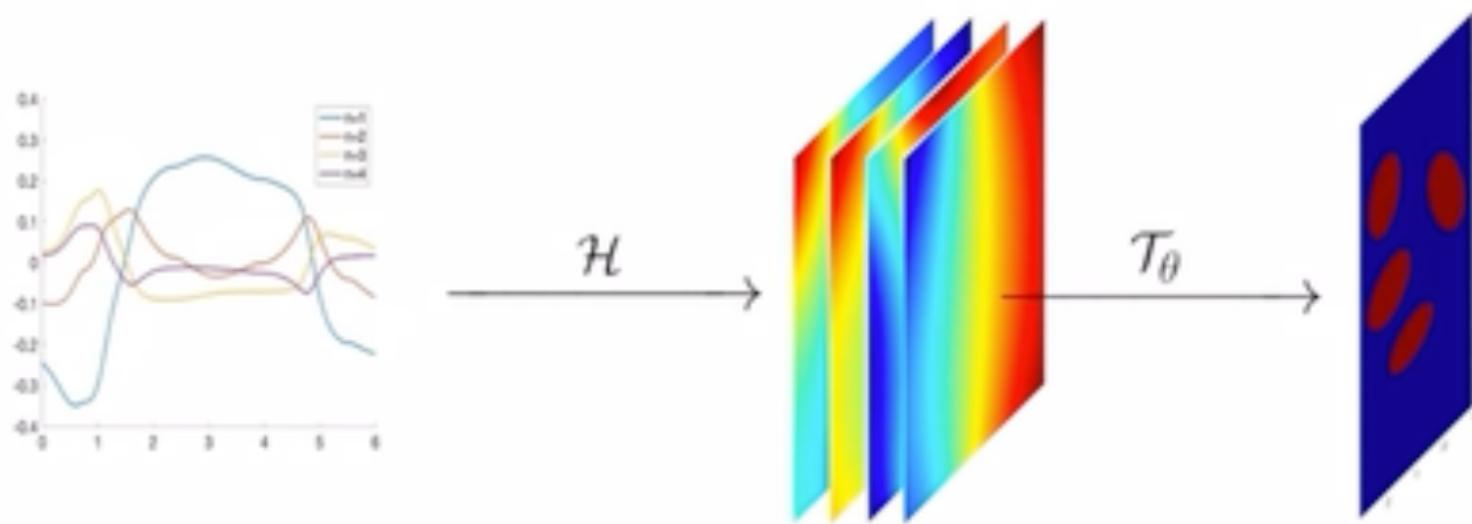
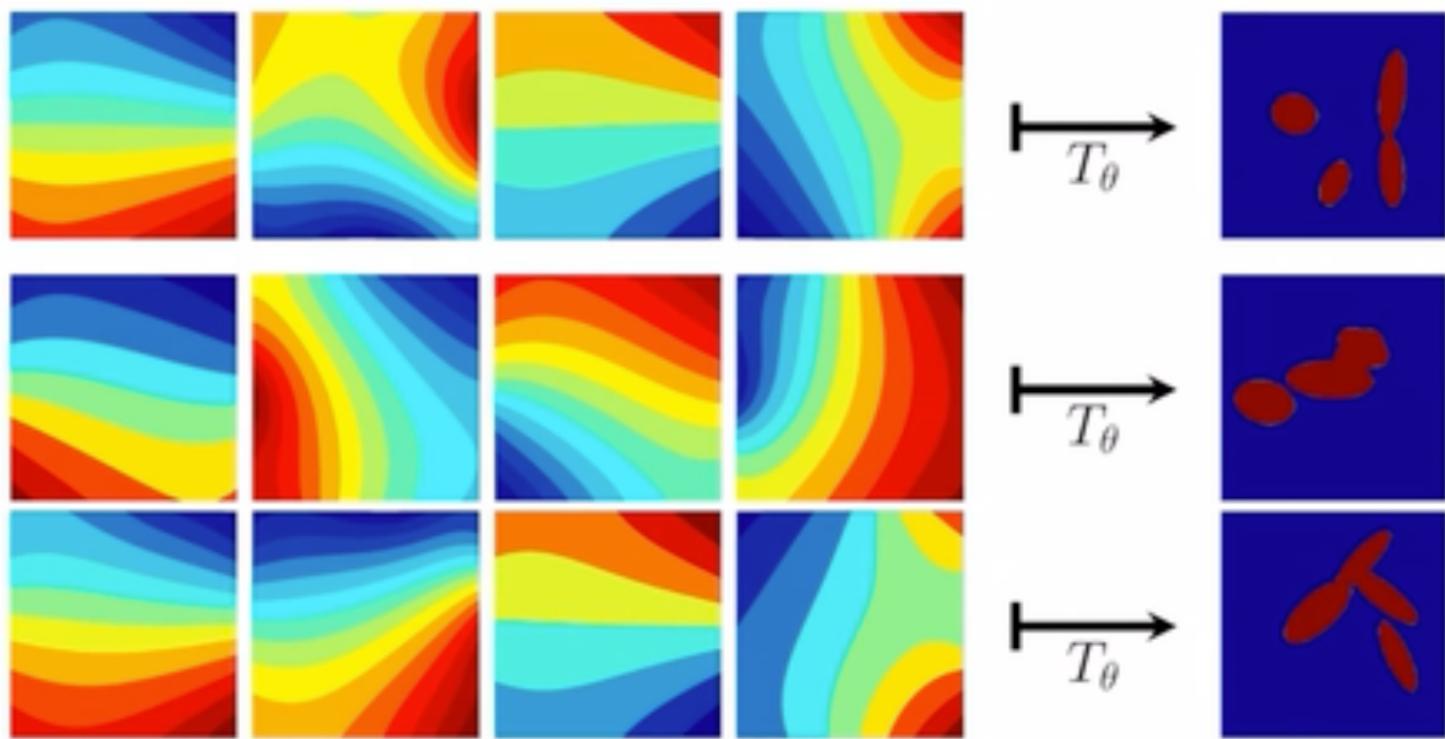
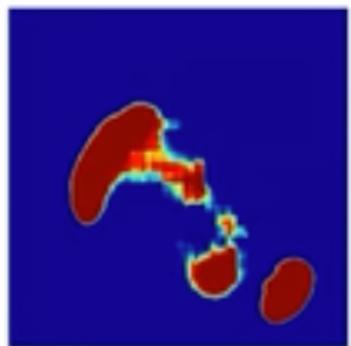


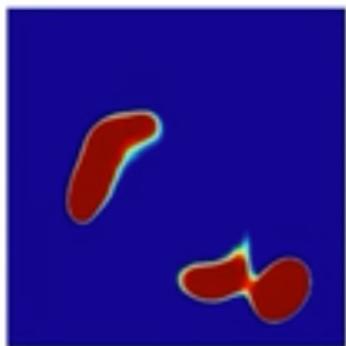
Image to Image



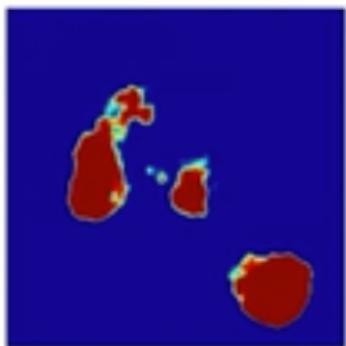
Learn a *single* parametrized operator T_θ that maps harmonic extension features to reconstruct the inclusions.



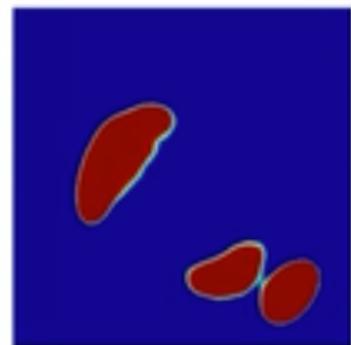
U-Net (7.7m, $L = 1$)



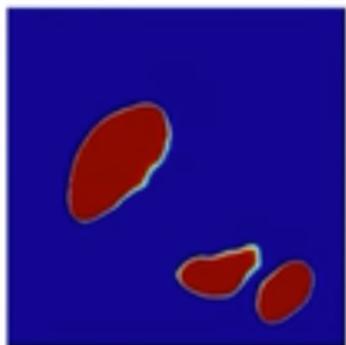
U-Net big (31m, $L = 3$)



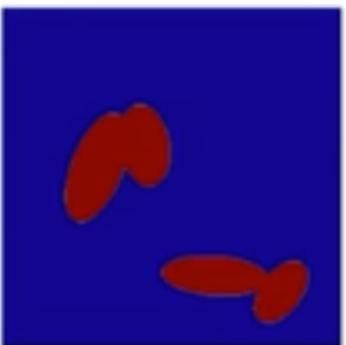
FNO (33m, $L = 1$)



UIT (11.4m, $L = 1$)



UIT (11.4m, $L = 3$)



Ground truth

6 PDEFormer

Consider 1D time-dependent PDEs on $(t, x) \in [0, 1] \times [0, 1]$ with periodic boundary condition

$$\mathcal{F}(u, c_1, \dots) = 0, u(0, x) = g(x)$$

where $c_i \in \mathbb{R}$ are real-valued coefficients, $g(x)$ is the initial condition. Assume that \mathcal{F} has symbolic expression which involves differential and algebraic operations.

Goal: Construct a surrogate of the solution $(\mathcal{F}, g, c_1, \dots) \rightarrow u$ operator.

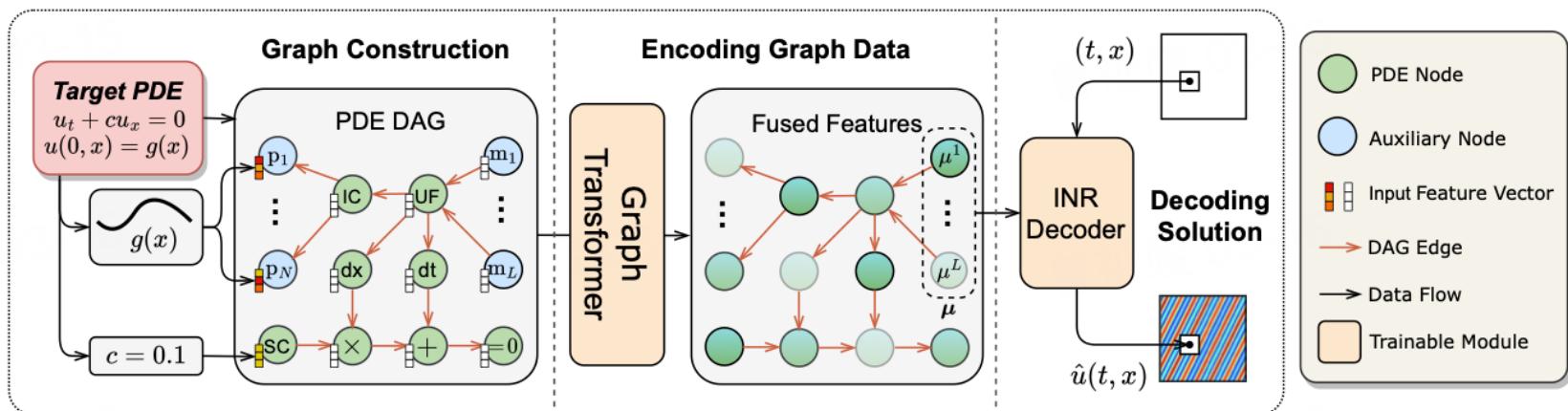


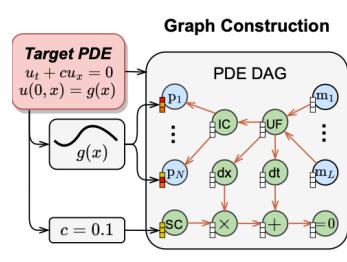
Figure 1: PDEformer architecture, taking $\mathcal{F}(u, c) = u_t + cu_x$ as the example.

6.1 Graph Construction

For \mathcal{F} , the symbolic information specifying the PDE, we use Computational Graph to represent.

Node: feature vector \mathbb{R}^{d_f}

1. Unknown field variable (UF): a node that represent unknown field u , we introduce $m_1 \dots, m_L$ additional nodes that connected to UF, which is used to decode the prediction.
2. scalar coefficient (SC): c repeated d_f times
3. Initial condition (IC): $g(x)$ given at an equi-space with n_x grid point and divide into $N = n_x / d_f$ patches, yielding N vector $g_1, \dots, g_N \in \mathbb{R}^{d_f}$, connect these node with IC node
4. Operation:node dx and dt for example.



6.2 Encoding Graph Data

We use graph-Transformer to encode symbolic and numeric information in the graph data into a latent code $\mu = [\mu^1, \dots, \mu^L] \in \mathbb{R}^{L \times d_e}$.

For $l = 1, \dots, L$ we let $\mu^l \in \mathbb{R}^{d_e}$ be the embedding vector assigned to the node with type m_l in the output.

6.2.1 GNN

Let $G = (V, E)$ be a graph, feature vector of node v_i be x_i , GNNs aims to learn representation of nodes and graphs. Modern GNNs follow a learning schema that iteratively updates the representation of a node by aggregating representations of its first or higher-order neighbors.

$h_i^{(l)}$ denotes the representation of v_i at the l -th layer and $h_i^{(0)} = x_i$. The AGGREGATE-COMBINE step as

$$a_i^{(l)} = \text{AGGREGATE}^{(l)}(\{h_j^{(l-1)} : j \in \mathcal{N}(v_i)\}), \quad h_i^{(l)} = \text{COMBINE}^{(l)}(h_i^{(l-1)}, a_i^{(l)}),$$

The AGGREGATE function is used to gather the information from neighbors. Common aggregation functions include MEAN, MAX, SUM, which are used in different architectures of GNNs. The goal of COMBINE function is to fuse the information from neighbors into the node representation.

for graph representation tasks, a READOUT function is designed to aggregate node features $h_i^{(L)}$ of the final iteration into h_G

$$h_G = \text{READOUT}(\{h_i^{\{L\}} | v_i \in G\})$$

6.2.2 Structural Encodings in Graphomer

Node centrality measures how important the a node is in the graph. We use the degree centrality as an additional signal to NN.

Assigns each node two real- valued embedding vectors according to its indegree and outdegree.

$$h_i^{(0)} = x_i + z_{\deg^-(v_i)}^- + z_{\deg^+(v_i)}^+$$

where $z^-, z^+ \in \mathbb{R}^d$ are learnable embedding vectors specified by indegree $\deg^-(v_i)$ and out degree $\deg^+(v_i)$

6.2.3 Spatial Encoding

For graphs, nodes are not arranged as a sequence. To encode structural information of graph G , we consider a function $\phi(v_i, v_j) : V \times V \rightarrow \mathbb{R}$ which measures the spatial relation between v_i, v_j if they connect. In GraphTransformer they use the shortest path as ϕ .

Modified Attention

$$A_{ij} = \frac{(h_i W_Q) (h_j W_K)^T}{\sqrt{d}} + b_{\phi(v_i, v_j)},$$

where $b_{\phi(v_i, v_j)}$ is a learnable scalar indexed by $\phi(v_i, v_j)$ and shared across all layers.

6.2.4 Edge Encoding in Attention

To encode edge features into attention layers, we

$$A_{ij} = \frac{(h_i W_Q) (h_j W_K)^T}{\sqrt{d}} + b_{\phi(v_i, v_j)} + c_{ij}, \text{ where } c_{ij} = \frac{1}{N} \sum_{n=1}^N x_{e_n} (w_n^E)^T$$

where x_{e_n} is the feature of the n -th edge e_n in $SP_{ij} = (e_1, \dots, e_N)$, $w_n^E \in \mathbb{R}^{d_E}$ is the n -th weight embedding and d_E is the dimensionality of edge feature.

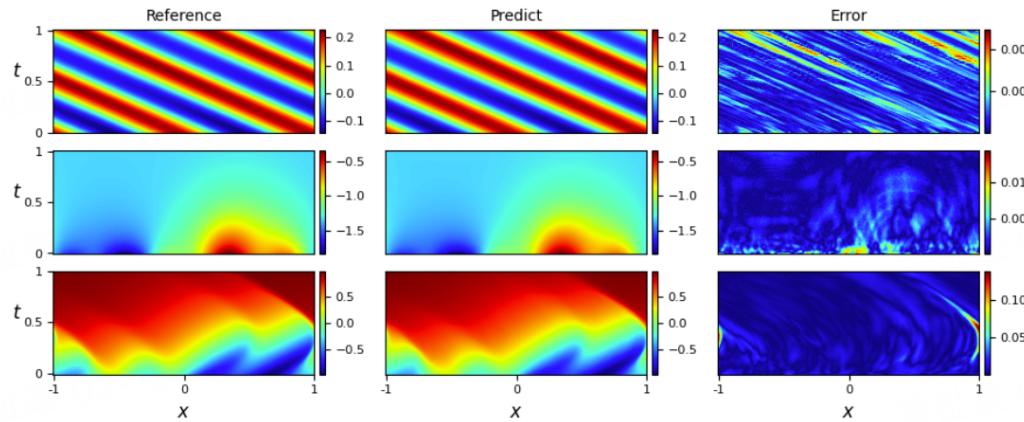
$$\begin{aligned} h'(l) &= \text{MHA}(\text{LN}(h^{(l-1)})) + h^{(l-1)} \\ h^{(l)} &= \text{FFN}(\text{LN}(h'(l))) + h'(l) \end{aligned}$$

6.3 Decoder

To obtain a mesh-free prediction, we employ implicit neural representation that takes (t, x) as input and produces $\hat{u}(t, x)$ according to μ .

Table 1: Test relative L^2 error on PDEBench. Format the first and second best outcomes in bold and underline respectively.

Model	Burgers			Adection		Reaction-Diffusion
	$\nu = 0.1$	$\nu = 0.01$	$\nu = 0.001$	$\beta = 0.1$	$\beta = 1$	$\nu = 1, \rho = 1$
U-Net (Ronneberger et al., 2013)	0.1627	0.2253	0.2431	0.0873	0.2655	0.0126
DeepONet (Lu et al., 2021)	0.0699	0.1791	0.2010	0.0186	0.0187	<u>0.0015</u>
FNO (Li et al., 2021)	0.0155	0.0445	<u>0.0700</u>	<u>0.0089</u>	<u>0.0097</u>	0.0018
PDEformer (Ours)	<u>0.0103</u>	<u>0.0309</u>	0.0921 (OoD)	0.0119	0.4000 (OoD)	0.7399 (OoD)
PDEformer-FT (Ours)	0.0046	0.0146	0.0295	0.0043	0.0075	0.0009



7 GNOT:General Neural Operator Transformer

Consider PDEs in the domain $\Omega \subset \mathbb{R}^d$ and the function space \mathcal{H} over Ω , including boundary shapes and source functions.

Goal: learning an operator \mathcal{G} from the input function space \mathcal{A} to the solution space \mathcal{H} , i.e., $\mathcal{G}: \mathcal{A} \rightarrow \mathcal{H}$.

Input function space: \mathcal{A} could contain multiple different types, like boundary shapes, source functions distributed over Ω , and vector parameters of the systems.

More formally, \mathcal{A} could be represented as $\mathcal{A} = \mathcal{H} \times \cdots \times \mathcal{H} \times \mathbb{R}^p$. For $\forall a = (a^1(\cdot), \dots, a^m(\cdot), \theta) \in \mathcal{A}$, $a^j(\cdot) \in \mathcal{H}$ represents boundary shapes and source functions, and $\theta \in \mathbb{R}^p$ represents parameters of the system, and $\mathcal{G}(a) = u \in \mathcal{H}$ is the solution function over Ω .

Dataset: $\mathcal{D} = \{(a_k, u_k)\}_{1 \leq k \leq D}$, where $u_k = \mathcal{G}(a_k)$. We discretize the input functions and the solution function on irregular discretized meshes over the domain Ω . For an input function a_k , we discretize it on the mesh $\{x_i^j \in \Omega\}_{1 \leq i \leq N_j}^{1 \leq j \leq m}$ and the discretized a_k^j is $\{(x_i^j, a_k^{i,j})\}_{1 \leq i \leq N_j}^{1 \leq j \leq m}$, where $a_k^{i,j} = a_k^j(x_i^j)$. In this way, we use $\mathcal{A}_k = \{(x_i^j, a_k^{i,j})\}_{1 \leq i \leq N_j}^{1 \leq j \leq m} \cup \theta_k$ to represent the input functions a_k .

Solution function: u_k is discretized on mesh $\{y_i \in \Omega\}_{1 \leq i \leq N'}$ and become $\{(y_i, u_k^i)\}_{1 \leq i \leq N'}$, here $u_k^i = u_k(y_i)$.

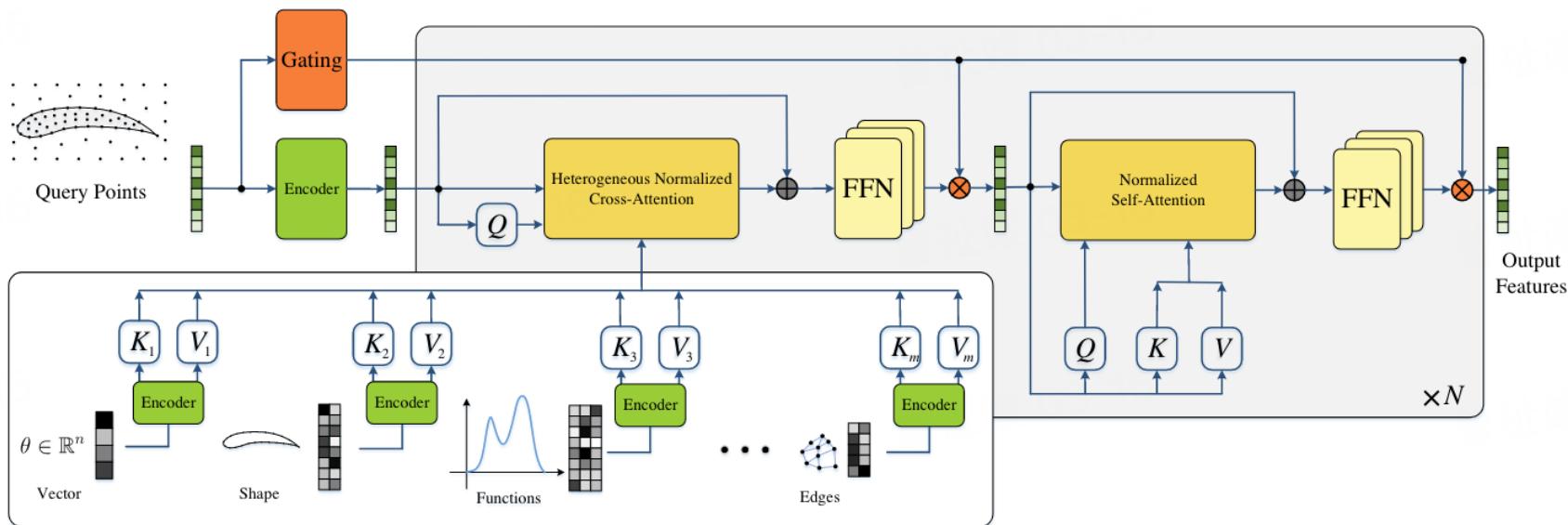
Operator \mathcal{G} : Use a parameterized neural network $\tilde{\mathcal{G}}_w$, which receives the input $\mathcal{A}_k (k=1, \dots, D)$ and outputs $\tilde{\mathcal{G}}_w(\mathcal{A}_k) = \{\tilde{u}_k^i\}_{1 \leq i \leq N'}$ to approximate u_k . Our goal is to minimize the mean squared error(MSE) loss between the prediction and data as

$$\min_{w \in W} \frac{1}{D} \sum_{k=1}^D \frac{1}{N'} \|\tilde{\mathcal{G}}_w(\mathcal{A}_k) - \{u_k^i\}_{1 \leq i \leq N'}\|_2^2, \quad (4)$$

where w is a set of the network parameters and W is the parameter space.

7.1 Architecture Overview

1. General input encoder: embed different input functions with MLPs
2. Attention block (self/cross) to modelling the latent operator
3. Geometric gating mechanism to tackle multiscale properties



7.2 General Input Encoding

Input of model :

Query Point: $\{x_i^q\}_{1 \leq i \leq N_q}$ map into query embedding $X \in \mathbb{R}^{N_q \cdot n_e}$

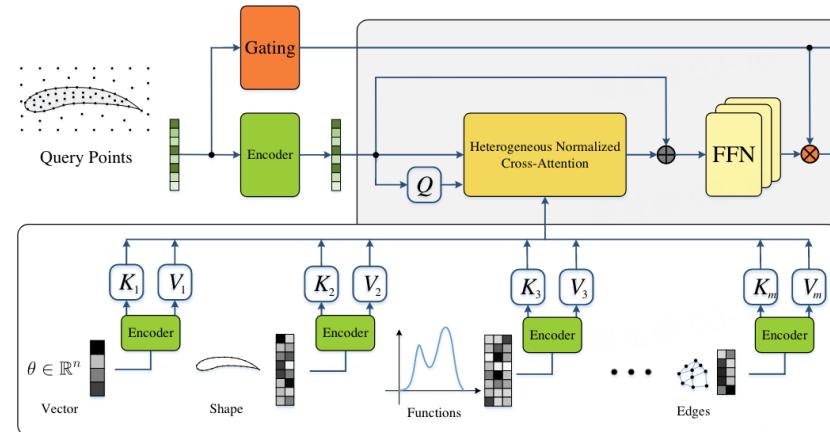
Parameter vector $\theta \in \mathbb{R}^p$ map into $Y = f_\omega(\theta) \in \mathbb{R}^{1 \times n_e}$

Boundary shape $\{x_i\}_{1 \leq i \leq N}$, map into $Y = (f_w(x_i))_{1 \leq i \leq N} \in \mathbb{R}^{N \cdot n_e}$

Domain distributed function $\{(x_i, a_i)\}_{1 \leq i \leq N}$ map into $Y = (f_w(x_i, a_i))_{1 \leq i \leq N} \in \mathbb{R}^{Nn_e}$

Mesh points and information $\{(x_i, z_i)\}_{1 \leq i \leq N}$ map into $Y = (f_w(x_i, z_i))_{1 \leq i \leq N}$

Edge information $\{(x_i^{\text{src}}, x_i^{\text{dst}}, e_i)\}_{1 \leq i \leq N}$ map into $Y = (f_w(x_i^{\text{src}}, x_i^{\text{dst}}, e_i))_{1 \leq i \leq N}$



7.3 Heterogeneous Normalized Attention Block

Here we have features of query points X and conditional embeddings $\{Y_l\}_{l=1}^L$. We need to design self attention for X and cross attention for X and Y_l .

Suppose we have three sequences called queries $\{q_i\}_{1 \leq i \leq N}$, keys $\{k_i\}_{1 \leq i \leq M}$ and values $\{v_i\}_{1 \leq i \leq M}$. The attention is computed as follows,

$$z_t = \sum_i \frac{\exp(q_t \cdot k_i / \tau)}{\sum_j \exp(q_t \cdot k_j / \tau)} v_i, \quad (5)$$

where τ is a hyperparameter.

For self-attention models, q, k, v are obtained by applying a linear transformation to input sequence $X = (x_i)_{1 \leq i \leq N}$, i.e., $q_i = W_q x_i$, $k_i = W_k x_i$, $v_i = W_v x_i$. For cross attention models, q comes from the query sequence X while keys and values come from another sequence $Y = (y_i)_{1 \leq i \leq M}$, i.e., $q_i = W_q x_i$, $k_i = W_k y_i$, $v_i = W_v y_i$.

However, the computational cost of the attention is $O(N^2 n_e)$ for self attention and $O(NM n_e)$ for cross attention where n_e is the dimension of embedding.

We first normalize these sequences respectively,

$$\tilde{\mathbf{q}}_i = \text{Softmax}(\mathbf{q}_i) = \left(\frac{e^{q_{ij}}}{\sum_j e^{q_{ij}}} \right)_{j=1, \dots n_e}, \quad (6)$$

$$\tilde{\mathbf{k}}_i = \text{Softmax}(\mathbf{k}_i) = \left(\frac{e^{k_{ij}}}{\sum_j e^{k_{ij}}} \right)_{j=1, \dots n_e}. \quad (7)$$

Then we compute the attention output without softmax using the following equation,

$$\mathbf{z}_t = \sum_i \frac{\tilde{\mathbf{q}}_t \cdot \tilde{\mathbf{k}}_i}{\sum_j \tilde{\mathbf{q}}_t \cdot \tilde{\mathbf{k}}_j} \cdot \mathbf{v}_i. \quad (8)$$

We denote $\alpha_t = (\sum_j \tilde{\mathbf{q}}_t \cdot \tilde{\mathbf{k}}_j)^{-1}$ and the efficient attention could be represented by,

$$\mathbf{z}_t = \sum_i \alpha_t (\tilde{\mathbf{q}}_t \cdot \tilde{\mathbf{k}}_i) \cdot \mathbf{v}_i = \alpha_t \tilde{\mathbf{q}}_t \cdot \left(\sum_i \tilde{\mathbf{k}}_i \otimes \mathbf{v}_i \right). \quad (9)$$

We could compute $\sum_i \tilde{\mathbf{k}}_i \otimes \mathbf{v}_i$ first with a cost $O(M n_e^2)$ and then compute its multiplication with \mathbf{q} with a cost $O(N n_e^2)$. The total cost is $O((M + N) n_e^2)$ which is linear with respect to the sequence length.

To fuse the query embeddings with multiple conditional embeddings, we design cross attention with above mechanism:

Specifically, suppose we have L conditional embeddings $\{Y_l \in \mathbb{R}^{N_l \times n_e}\}_{1 \leq l \leq L}$ encoding the input functions and extra information. We first compute the queries $Q = (\mathbf{q}_i) = XW_q$, keys $K_l = (\mathbf{k}_i^l) = YW_k$ and values $V_l = (\mathbf{v}_i^l) = YW_v$, and then normalize every \mathbf{q}_i and \mathbf{k}_i to be $\tilde{\mathbf{q}}_i$ and $\tilde{\mathbf{k}}_i$. Then we compute the cross-attention as follows,

$$\mathbf{z}_t = \tilde{\mathbf{q}}_t + \frac{1}{L} \sum_{l=1}^L \sum_{i_l=1}^{N_l} \alpha_t^l (\tilde{\mathbf{q}}_t \cdot \tilde{\mathbf{k}}_{i_l}) \mathbf{v}_{i_l}, \quad (10)$$

$$= \tilde{\mathbf{q}}_t + \frac{1}{L} \sum_{l=1}^L \alpha_t^l \tilde{\mathbf{q}}_t \cdot \left(\sum_{i_l=1}^{N_l} \tilde{\mathbf{k}}_{i_l} \otimes \mathbf{v}_{i_l} \right). \quad (11)$$

where $\alpha_t^l = \frac{1}{\sum_{j=1}^{N_l} \tilde{\mathbf{q}}_t \cdot \tilde{\mathbf{k}}_j}$ is the normalization coefficient.

After applying such a cross-attention layer, we impose the self-attention layer for query features, i.e,

$$\mathbf{z}'_t = \sum_i \alpha_t (\tilde{\mathbf{q}}_t \cdot \tilde{\mathbf{k}}_i) \cdot \mathbf{v}_i, \quad (12)$$

where all of \mathbf{q} , \mathbf{k} and \mathbf{v} are computed with the embedding \mathbf{z}_t as

$$\mathbf{q}_t = W_q \hat{\mathbf{z}}_t, \mathbf{k}_t = W_k \hat{\mathbf{z}}_t, \mathbf{v}_t = W_v \hat{\mathbf{z}}_t. \quad (13)$$

We tile multiple layers and multiple heads similar to other transformer models. The embedding \mathbf{z}_t and \mathbf{z}'_t are divided into H heads as $\mathbf{z}_t = \text{Concat}(\mathbf{z}_t^i)_{i=1}^H$ and $\mathbf{z}'_t = \text{Concat}(\mathbf{z}'_t^i)_{i=1}^H$

7.4 Geometric Gating Mechanism

To handle multi-scale problems, we introduce our geometric gating mechanism based on mixture-of-experts (MoE): a geometric gating network that inputs the coordinates of the query points and outputs unnormalized scores $G_i(x)$ for averaging these expert networks.

In each layer of our model, we use K subnetworks for the MLP denoted by $E_i(\cdot)$. The update of \mathbf{z}_t and \mathbf{z}'_t in the feedforward layer is replaced by the following equation when we have multiple expert networks as

$$\mathbf{z}_t \leftarrow \mathbf{z}_t + \sum_{i=1}^K p_i(x_t) \cdot E_i(\mathbf{z}_t). \quad (14)$$

The weights for averaging the expert networks are computed as

$$p_i(x_t) = \frac{\exp(G_i(x_t))}{\sum_{i=1}^K \exp(G_i(x_t))}, \quad (15)$$

where the gating network $G(\cdot): \mathbb{R}^d \rightarrow \mathbb{R}^K$ takes the geometric coordinates of query points x_t as inputs. The normalized outputs $p_i(x_t)$ are the weights for averaging these experts.

7.5 Experiment

“A”, “B”, and “C” represent the problem has irregular mesh, has multiple input functions, and is multi-scale

Dataset	Type	MIONet	FNO(-interp)	GK-Transformer	Geo-FNO	OFormer	Ours
	Challenge	Subset					
Darcy2d	-	-	5.45e-2	1.09e-2	8.40e-3	1.09e-2	1.24e-2
NS2d	-	part	—	1.56e-1	1.40e-1	1.56e-1	1.71e-1
	-	full	—	8.20e-2	7.92e-2	8.20e-2	6.46e-2
Elasticity	A	-	9.65e-2	5.08e-2	2.01e-2	2.20e-2	1.83e-2
NS2d-c	A, C	u	2.74e-2	6.56e-2	1.52e-2	1.41e-2	2.33e-2
		v	5.51e-2	1.15e-1	3.15e-2	2.98e-2	4.83e-2
		p	2.74e-2	1.11e-2	1.59e-2	1.62e-2	2.43e-2
NACA	A, C	-	1.32e-1	4.21e-2	1.61e-2	1.38e-2	1.83e-2
Inductor2d	A, C	A_z	3.10e-2	—	2.56e-1	—	2.23e-2
		B_x	3.49e-2	—	3.06e-2	—	2.83e-2
		B_y	6.73e-2	—	4.45e-2	—	4.28e-2
Heat	A, B, C	part	1.74e-1	—	—	—	—
		full	1.45e-1	—	—	—	—
Heatsink	A, B, C	T	4.67e-1	—	—	—	—
		u	3.52e-1	—	—	—	—
		v	3.23e-1	—	—	—	—
		w	3.71e-1	—	—	—	—

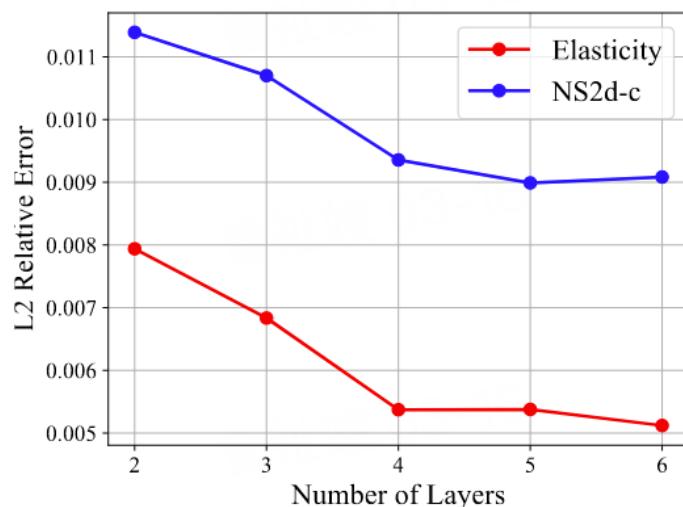
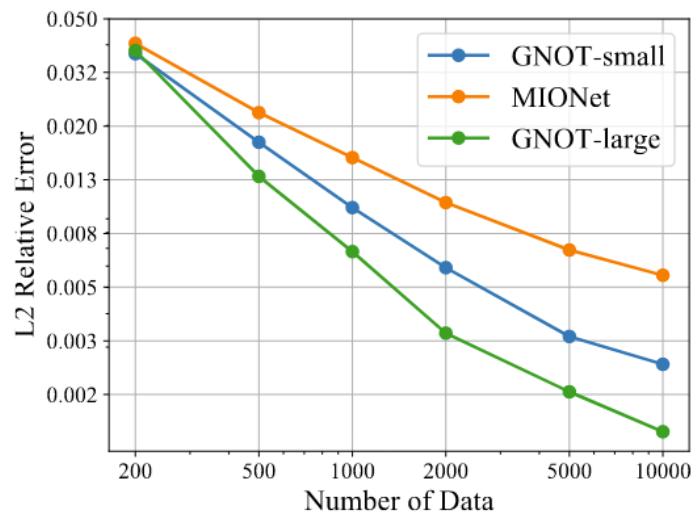


Figure 3. Results of scaling experiments for different dataset sizes (left) and different numbers of layers (right).

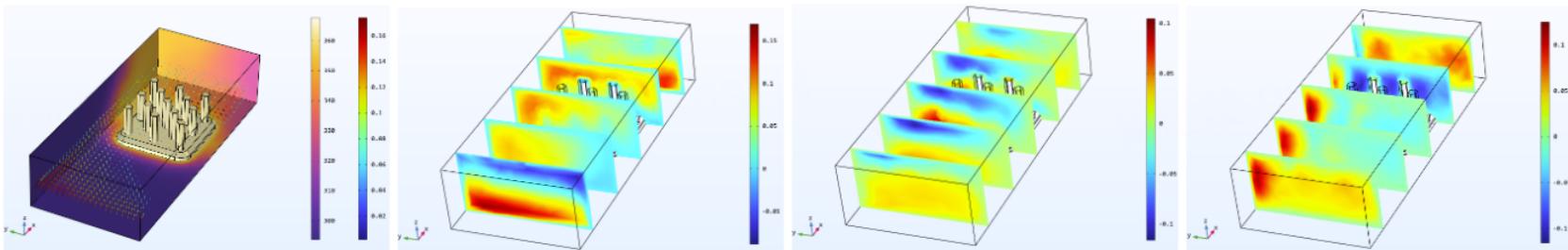


Figure 10. Visualization of T, u, v, w of Heatsink dataset.

8 Transolver: Transformer Solver for PDEs on General Geometry

Consider PDEs defined on input domain $\Omega \subset \mathbb{R}^{C_g}$, where C_g denotes the dimension of input space. For numerical calculation, Ω is firstly discretized into a finite set of N mesh points $\mathbf{g} \in \mathbb{R}^{N \times C_g}$. The task is to estimate target physical quantities based on input geometrics \mathbf{g} and quantities $\mathbf{u} \in \mathbb{R}^{N \times C_u}$ observed on \mathbf{g} . Here \mathbf{u} is optional in some PDE-governed tasks.

Physics-Aware Tokens: Given a mesh set $\mathbf{g} = \{\mathbf{g}_i\}_{i=1}^N$ with the coordinate information of N mesh points and observed quantities \mathbf{u} , we firstly embed them into deep features $\mathbf{x} = \{\mathbf{x}_i\}_{i=1}^N$ by a linear layer contains C channels, $\mathbf{x}_i \in \mathbb{R}^{1 \times C}$. Then we propose a bottom-up paradigm, that ascribes each mesh point \mathbf{g}_i to M potential slices based on its learned feature \mathbf{x}_i , which is formalized as follows:

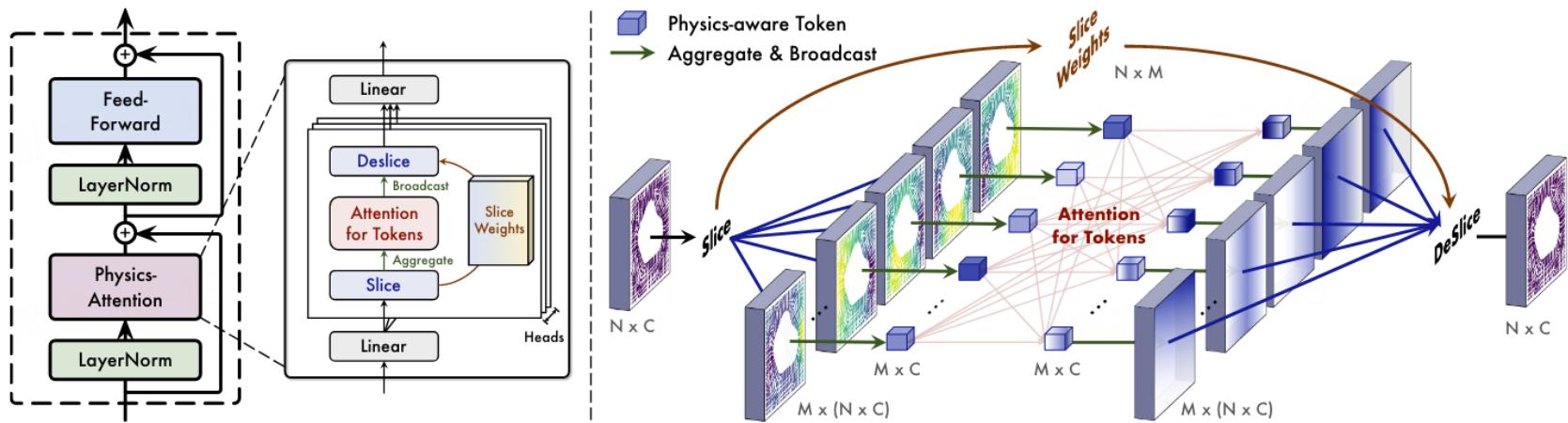
$$\begin{aligned} \{\mathbf{w}_i\}_{i=1}^N &= \{\text{Softmax}(\text{Linear}(\mathbf{x}_i))\}_{i=1}^N \\ \mathbf{s}_j &= \{\mathbf{w}_{i,j} \mathbf{x}_i\}_{i=1}^N, \end{aligned} \tag{16}$$

where $\text{Linear}()$ projects C channels into M weights and yields slice weights $\mathbf{w}_i \in \mathbb{R}^{1 \times M}$ after $\text{Softmax}()$. Specifically, $\mathbf{w}_{i,j}$ represents the degree that the i -th mesh point belongs to the j -th slices with $\sum_{j=1}^M \mathbf{w}_{i,j} = 1$. $\mathbf{s}_j \in \mathbb{R}^{N \times C}$ represents the j -th slice feature, which is a weighted combination of N mesh point features \mathbf{x} .

Afterward, since each slice contains mesh points with similar geometry and physics features, we further encode them into physical-aware tokens by spatially weighted aggregation, which can be written as follows:

$$\mathbf{z}_j = \frac{\sum_{i=1}^N \mathbf{s}_{j,i}}{\sum_{i=1}^N \mathbf{w}_{i,j}} = \frac{\sum_{i=1}^N \mathbf{w}_{i,j} \mathbf{x}_i}{\sum_{i=1}^N \mathbf{w}_{i,j}}, \quad (17)$$

where $\mathbf{z}_j \in \mathbb{R}^{1 \times C}$.



8.1 Physics-Attention

After obtained physics-aware tokens $\mathbf{z} = \{\mathbf{z}_j\}_{j=1}^M \in \mathbb{R}^{M \times C}$, we employ the attention mechanism among encoded tokens to capture intricate correlations among different physical states, that is

$$\begin{aligned} \mathbf{q}, \mathbf{k}, \mathbf{v} &= \text{Linear}(\mathbf{z}) \\ \mathbf{z}' &= \text{Softmax}\left(\frac{\mathbf{q}\mathbf{k}^T}{\sqrt{C}}\right)\mathbf{v}, \end{aligned} \tag{18}$$

where $\mathbf{q}, \mathbf{k}, \mathbf{v}, \mathbf{z}' \in \mathbb{R}^{M \times C}$. Afterward, transited physical tokens $\mathbf{z}' = \{\mathbf{z}'_j\}_{j=1}^M$ are transformed back to mesh points by deslicing, which recomposes tokens with slice weights:

$$\mathbf{x}'_i = \sum_{j=1}^M \mathbf{w}_{i,j} \mathbf{z}'_j, \tag{19}$$

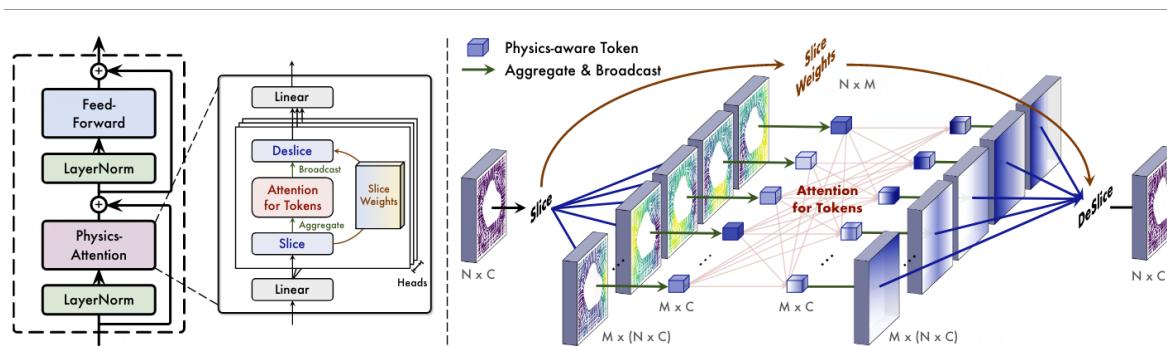
where $1 \leq i \leq N$ and each token \mathbf{z}'_j is broadcasted to all mesh points during above calculation. For clarity, we summarize the above process as $\mathbf{x}' = \text{Physics-Attn}(\mathbf{x})$, whose overall complexity is $\mathcal{O}(NMC + M^2C)$.

8.2 Overall Design

We propose the Transolver by replacing the attention mechanism with Physics-Attention. Suppose there are L layers, the l -th layer of Transolver can be formalized as follows:

$$\begin{aligned}\hat{\mathbf{x}}^l &= \text{Physics-Attn}(\text{LayerNorm}(\mathbf{x}^{l-1})) + \mathbf{x}^{l-1} \\ \mathbf{x}^l &= \text{FeedForward}(\text{LayerNorm}(\hat{\mathbf{x}}^l)) + \hat{\mathbf{x}}^l,\end{aligned}\quad (20)$$

where $l \in \{1, \dots, L\}$. $\mathbf{x}^l \in \mathbb{R}^{N \times C}$ is the output of the l -th layer. $\mathbf{x}^0 \in \mathbb{R}^{N \times C}$ represents the input deep feature, which is embedded from input geometries $\mathbf{g} \in \mathbb{R}^{N \times C_g}$ and initial observation $\mathbf{u} \in \mathbb{R}^{N \times C_u}$ by a linear embedding layer, i.e. $\mathbf{x}^0 = \text{Linear}(\text{Concat}(\mathbf{g}, \mathbf{u}))$. Here C_g is the dimension of geometry space and C_u is the number of observed physical quantities. At last, we adopt a linear projection upon \mathbf{x}^L and obtain the final output as predictions of \mathbf{u} .



8.3 Physics-Attention is equivalent to learnable integral on Ω

引理 1. The canonical attention mechanism in Transformers is a Monte- Carlo approximation of an integral operator.

证明. Given input function $u: \Omega \rightarrow \mathbb{R}^C$, the integral operation \mathcal{G} defined on the function space $\Omega \rightarrow \mathbb{R}^C$ is formalized as:

$$\mathcal{G}(u)(g^*) = \int_{\Omega} \kappa(g^*, \xi) u(\xi) d\xi, \quad (21)$$

where $g^* \in \Omega \subset \mathbb{R}^{C_g}$ and $\kappa(\cdot, \cdot)$ denotes the kernel function defined on Ω . According to the formalization of attention, we propose to define the kernel function as follows:

$$\kappa(g^*, \xi) = \left(\int_{\Omega} \exp((W_q u(\xi')) (W_k u(\xi))^T) d\xi' \right)^{-1} \exp((W_q u(g^*)) (W_k u(\xi))^T) W_v, \quad (22)$$

where $W_q, W_k, W_v \in \mathbb{R}^{C \times C}$.

□

证明. Suppose that there are N discretized mesh points $\{\mathbf{g}_1, \dots, \mathbf{g}_N\}$, where $\mathbf{g}_i \in \Omega \subset \mathbb{R}^{C_g}$. Approximating the inner-integral in Eq. (22) by Monte-Carlo, we have:

$$\int_{\Omega} \exp((\mathbf{W}_q \mathbf{u}(\xi')) (\mathbf{W}_k \mathbf{u}(\xi))^T) d\xi' \approx \frac{|\Omega|}{N} \sum_{i=1}^N \exp((\mathbf{W}_q \mathbf{u}(\mathbf{g}_i)) (\mathbf{W}_k \mathbf{u}(\xi))^T). \quad (23)$$

Applying the above equation to Eq. (21) and using the same approximation for the outer-integral, we have:

$$\mathcal{G}(\mathbf{u})(\mathbf{g}^*) \approx \sum_{i=1}^N \frac{\exp((\mathbf{W}_q \mathbf{u}(\mathbf{g}^*)) (\mathbf{W}_k \mathbf{u}(\mathbf{g}_i))^T) \mathbf{W}_v \mathbf{u}(\mathbf{g}_i)}{\sum_{j=1}^N \exp((\mathbf{W}_q \mathbf{u}(\mathbf{g}_j)) (\mathbf{W}_k \mathbf{u}(\mathbf{g}_i))^T)}, \quad (24)$$

which is the calculation of the attention mechanism with \mathbf{W}_q , \mathbf{W}_k , \mathbf{W}_v as linear layers for queries, keys and values. \square

引理 2. Suppose that Ω is a countable domain, the slice domain Ω_s is isomorphic to Ω .

定理 3. (Physics-Attention is equivalent to learnable integral on Ω) Given input function $u: \Omega \rightarrow \mathbb{R}^C$ and a mesh point $g^* \in \Omega$, Physics-Attention is to approximate the integral operator \mathcal{G} , which is defined as:

$$\mathcal{G}(u)(g^*) = \int_{\Omega} \kappa(g^*, \xi) u(\xi) d\xi, \quad (25)$$

where $\kappa(\cdot, \cdot)$ denotes the kernel function defined on $\Omega \times \Omega$.

证明. According to Lemma 2, we can obtain an isomorphic projection g between a countable input domain Ω and slice domain Ω_s . Suppose that the slice weight $w_{*,*}: \bar{\Omega} \times \bar{\Omega}_s \rightarrow \mathbb{R}$ is smooth in both $\bar{\Omega}$ and $\bar{\Omega}_s$, where the $\bar{\Omega}$ and $\bar{\Omega}_s$ denote the continuation of Ω and Ω_s respectively, we can obtain g as a diffeomorphism projection.

Then, we define the value function u_s on the physics-aware token domain Ω_s as follows:

$$u_s(\xi_s) = \left(\int_{\Omega} w_{\xi, \xi_s} u(\xi) d\xi \right) / \left(\int_{\Omega} w_{\xi, \xi_s} d\xi \right), \quad (26)$$

which corresponds to the slice token definition in Eq. (17). \square

证明. Based on the above assumptions and definitions, we have:

$$\begin{aligned}
\mathcal{G}(\mathbf{u})(\mathbf{g}) &= \int_{\Omega} \kappa(\mathbf{g}, \boldsymbol{\xi}) \mathbf{u}(\boldsymbol{\xi}) d\boldsymbol{\xi} \\
&= \int_{\Omega_s} \kappa_{ms}(\mathbf{g}, \boldsymbol{\xi}_s) \mathbf{u}_s(\boldsymbol{\xi}_s) d\mathbf{g}^{-1}(\boldsymbol{\xi}_s) \\
&= \int_{\Omega_s} \kappa_{ms}(\mathbf{g}, \boldsymbol{\xi}_s) \mathbf{u}_s(\boldsymbol{\xi}_s) |\det(\nabla_{\boldsymbol{\xi}_s} \mathbf{g}^{-1}(\boldsymbol{\xi}_s))| d\boldsymbol{\xi}_s \\
&= \int_{\Omega_s} \left(\frac{\int_{\Omega_s} w_{\mathbf{g}, \boldsymbol{\xi}'_s} \kappa_{ss}(\boldsymbol{\xi}'_s, \boldsymbol{\xi}_s) d\boldsymbol{\xi}'_s}{\int_{\Omega_s} w_{\mathbf{g}, \boldsymbol{\xi}'_s} d\boldsymbol{\xi}'_s} \right) \mathbf{u}_s(\boldsymbol{\xi}_s) |\det(\nabla_{\boldsymbol{\xi}_s} \mathbf{g}^{-1}(\boldsymbol{\xi}_s))| d\boldsymbol{\xi}_s \\
&= \underbrace{\int_{\Omega_s} w_{\mathbf{g}, \boldsymbol{\xi}'_s} \int_{\Omega_s}}_{\text{DeSlice}} \underbrace{\kappa_{ss}(\boldsymbol{\xi}'_s, \boldsymbol{\xi}_s)}_{\text{Attention among slice tokens}} \underbrace{\mathbf{u}_s(\boldsymbol{\xi}_s)}_{\text{Slice token}} |\det(\nabla_{\boldsymbol{\xi}_s} \mathbf{g}^{-1}(\boldsymbol{\xi}_s))| d\boldsymbol{\xi}_s d\boldsymbol{\xi}'_s \\
&\approx \underbrace{\sum_{j=1}^M \mathbf{w}_{i,j}}_{\text{Eq. (19)}} \underbrace{\sum_{t=1}^M \frac{\exp((\mathbf{W}_q \mathbf{u}_s(\boldsymbol{\xi}_{s,j})) (\mathbf{W}_k \mathbf{u}_s(\boldsymbol{\xi}_{s,t}))^T / \tau)}{\sum_{p=1}^M \exp((\mathbf{W}_q \mathbf{u}_s(\boldsymbol{\xi}_{s,j})) (\mathbf{W}_k \mathbf{u}_s(\boldsymbol{\xi}_{s,p}))^T / \tau)}}_{\text{Eq. (18)}} \mathbf{W}_v \left(\underbrace{\frac{\sum_{p=1}^N \mathbf{w}_{p,t} \mathbf{u}(\mathbf{g}_p)}{\sum_{p=1}^N \mathbf{w}_{p,t}}}_{\text{Eq. (17)}} \right) \\
&= \sum_{j=1}^M \mathbf{w}_{i,j} \sum_{t=1}^M \frac{\exp(\mathbf{q}_j \mathbf{k}_t^T / \tau)}{\sum_{p=1}^M \exp(\mathbf{q}_j \mathbf{k}_p^T / \tau)} \mathbf{v}_t,
\end{aligned} \tag{27} \square$$

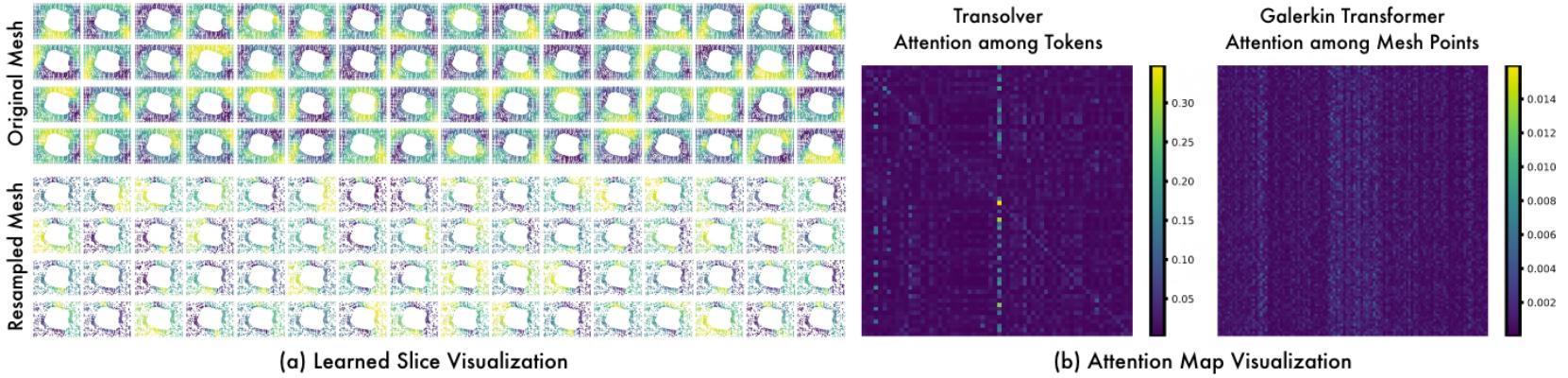


Figure 6. Physics-Attention visualization on Elasticity: (a) slice weights in the last layer of Transolver for both original and resampled meshes, (b) attention maps of the last layer in Transolver and Galerkin Transformer (Cao, 2021). See Appendix D.1 for more visualizations.

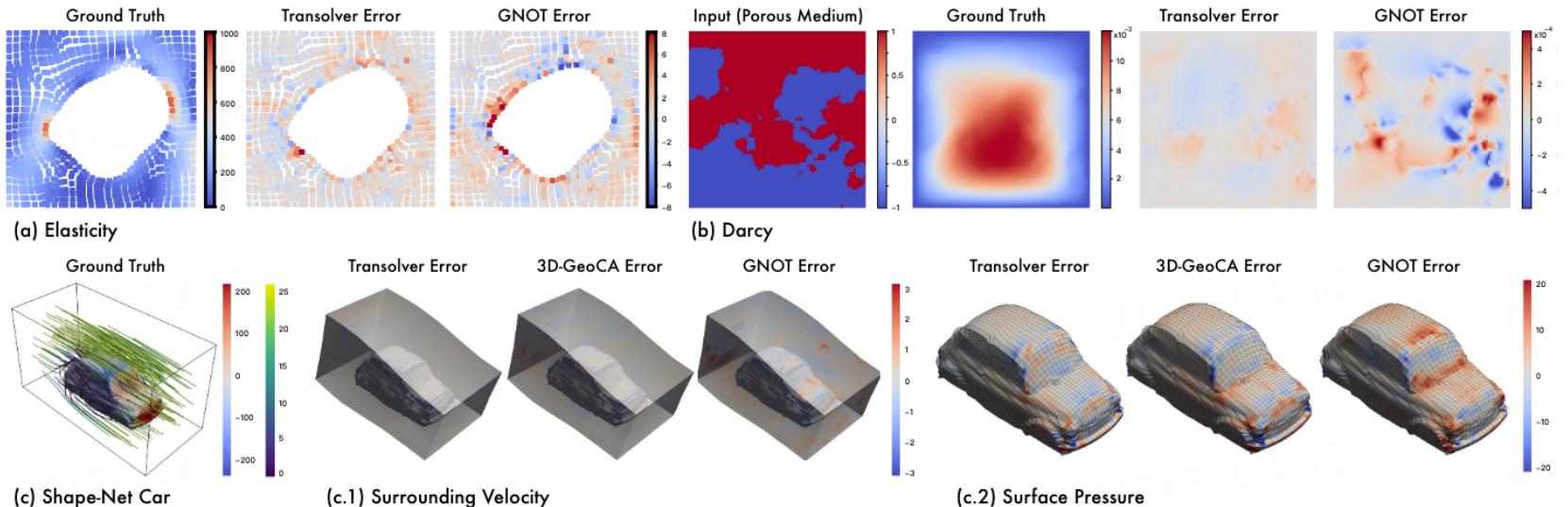


Figure 7. Case study on error maps of different models. Notably, Shape-Net Car requires to predict the surrounding velocity and surface pressure simultaneously. For clearness, we plot the error on volume and surface separately. See Appendix D.2 for more showcases.

MODEL*	SHAPE-NET CAR				AIRFRANS			
	VOLUME ↓	SURF ↓	C_D ↓	ρ_D ↑	VOLUME ↓	SURF ↓	C_L ↓	ρ_L ↑
SIMPLE MLP	0.0512	0.1304	0.0307	0.9496	0.0081	0.0200	0.2108	0.9932
GRAPH SAGE (HAMILTON ET AL., 2017)	0.0461	0.1050	0.0270	0.9695	0.0087	0.0184	0.1476	0.9964
POINTNET (QI ET AL., 2017)	0.0494	0.1104	0.0298	0.9583	0.0253	0.0996	0.1973	0.9919
GRAPH U-NET (GAO & JI, 2019)	0.0471	0.1102	0.0226	0.9725	0.0076	0.0144	0.1677	0.9949
MESHGRAPHNET (PFAFF ET AL., 2021)	0.0354	0.0781	0.0168	0.9840	0.0214	0.0387	0.2252	0.9945
GNO (LI ET AL., 2020A)	0.0383	0.0815	0.0172	0.9834	0.0269	0.0405	0.2016	0.9938
GALERKIN (CAO, 2021)	0.0339	0.0878	0.0179	0.9764	0.0074	0.0159	0.2336	0.9951
GEO-FNO (LI ET AL., 2022)	0.1670	0.2378	0.0664	0.8280	0.0361	0.0301	0.6161	0.9257
GNOT (HAO ET AL., 2023)	0.0329	0.0798	0.0178	0.9833	0.0049	0.0152	0.1992	0.9942
GINO (LI ET AL., 2023A)	0.0386	0.0810	0.0184	0.9826	0.0297	0.0482	0.1821	0.9958
3D-GEOCA (ANONYMOUS, 2023A)	0.0319	0.0779	0.0159	0.9842	/	/	/	/
TRANSOLVER (OURS)	0.0207	0.0745	0.0103	0.9935	0.0037	0.0142	0.1030	0.9978

MODEL*: The second best error / on each benchmark. / means that the baseline cannot apply to this benchmark.

MODEL	POINT CLOUD		STRUCTURED MESH			REGULAR GRID	
	ELASTICITY	PLASTICITY	AIRFOIL	PIPE	NAVIER-STOKES	DARCY	
FNO (LI ET AL., 2021)	/	/	/	/	0.1556	0.0108	
WMT (GUPTA ET AL., 2021)	0.0359	0.0076	0.0075	0.0077	0.1541	0.0082	
U-FNO (WEN ET AL., 2022)	0.0239	0.0039	0.0269	0.0056	0.2231	0.0183	
GEO-FNO (LI ET AL., 2022)	0.0229	0.0074	0.0138	0.0067	0.1556	0.0108	
U-NO (RAHMAN ET AL., 2023)	0.0258	0.0034	0.0078	0.0100	0.1713	0.0113	
F-FNO (TRAN ET AL., 2023)	0.0263	0.0047	0.0078	0.0070	0.2322	0.0077	
LSM (WU ET AL., 2023)	0.0218	0.0025	0.0059	0.0050	0.1535	0.0065	
GALERKIN (CAO, 2021)	0.0240	0.0120	0.0118	0.0098	0.1401	0.0084	
HT-NET (LIU ET AL., 2022)	/	0.0333	0.0065	0.0059	0.1847	0.0079	
OFORMER (LI ET AL., 2023C)	0.0183	0.0017	0.0183	0.0168	0.1705	0.0124	
GNOT (HAO ET AL., 2023)	0.0086	0.0336	0.0076	0.0047	0.1380	0.0105	
FACTFORMER (LI ET AL., 2023D)	/	0.0312	0.0071	0.0060	0.1214	0.0109	
ONO (ANONYMOUS, 2023B)	0.0118	0.0048	0.0061	0.0052	0.1195	0.0076	
TRANSOLVER (OURS)	0.0064	0.0012	0.0053	0.0033	0.0900	0.0057	
RELATIVE PROMOTION	25.6%	29.4%	10.2%	29.7%	24.7%	12.3%	

9 Multiple Physics Pretraining

Problem: Let S be an arbitrary physics-driven spatiotemporal dynamical systems, either described by a parameterized family of PDEs with fixed parameters.

System: To simplify notation, we discuss systems with a single state variable in one spatial dimension. A continuous state variable for system S is represented as $u^S(x, t) : [0, L_S] \times [0, \infty) \rightarrow \mathbb{R}$. We discretize the system uniformly in space and time at resolutions N_S, T_S respectively.

Data: A snapshot $u_t^S \in \mathbb{R}^{N_S}$ represents the value of state variable u^S at all N_S spatial discretization points at time t .

Task: learn a single model \mathcal{M} that can take a uniformly spaced sequence of T_S snapshots $u_t^S = [u_{t-T_s\Delta t_S}^S, \dots, u_t^S]$ from system S sampled from some distribution over systems and predict $\mathcal{M}(u_t^S)$ such that $\mathcal{M}(u_t^S) \approx u_{t+\Delta t_S}^S$.

9.1 Observation

1. Learning partially overlapping physics is beneficial for transfer learning
2. Single models can simultaneously learn many types of physics

Let $\psi(x, t)$ be a scalar defined on a periodic spatial domain, v a constant one-dimensional velocity coefficient and δ a constant diffusion coefficient, then:

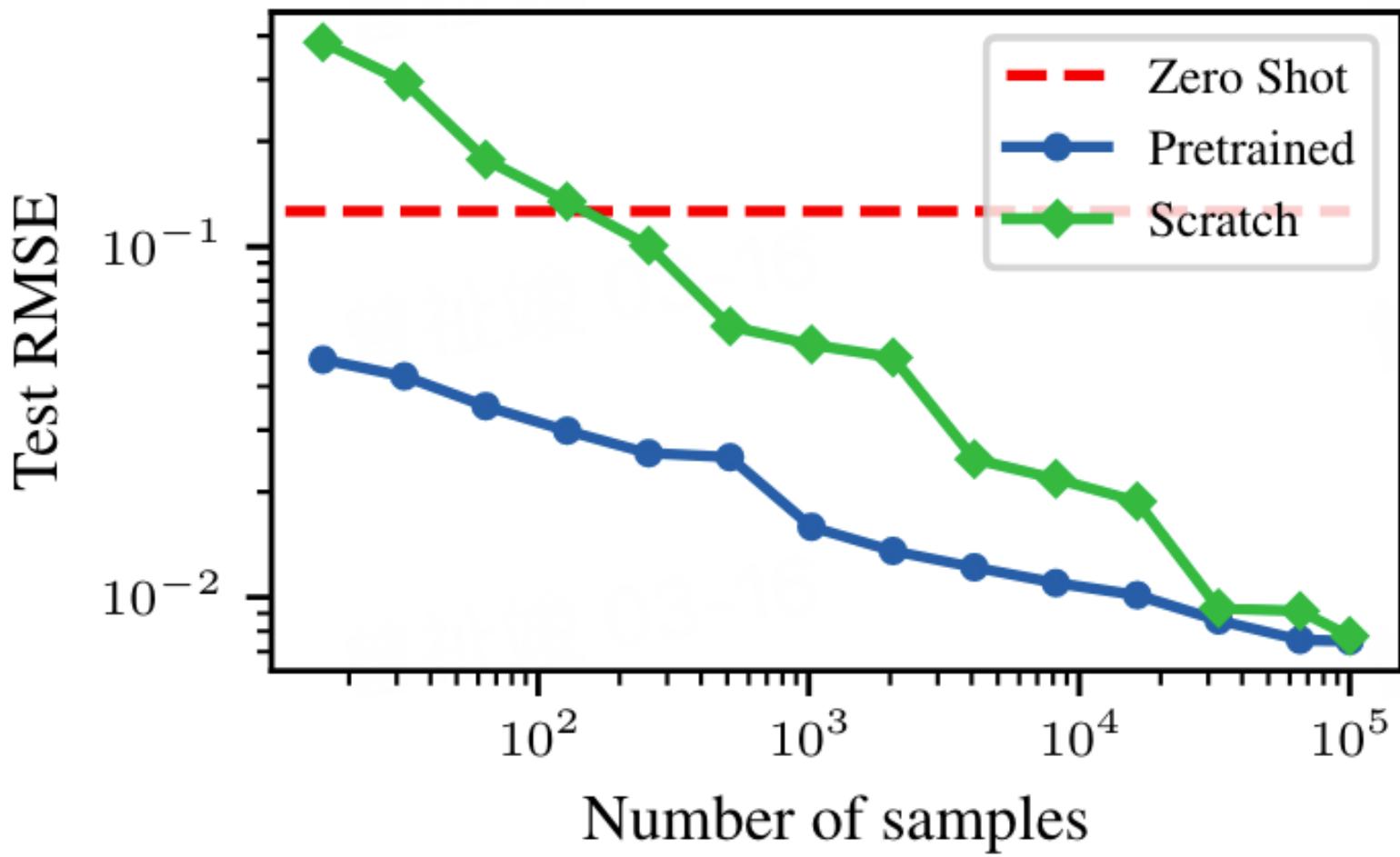
Advection:
$$\frac{\partial \psi}{\partial t} + \nabla \cdot (v \psi) = 0 \quad (28)$$

Diffusion:
$$\frac{\partial \psi}{\partial t} + \nabla \cdot (-\delta \nabla \psi) = 0 \quad (29)$$

Advection-Diffusion:
$$\frac{\partial \psi}{\partial t} + \nabla \cdot (v \psi - \delta \nabla \psi) = 0. \quad (30)$$

Experiment:

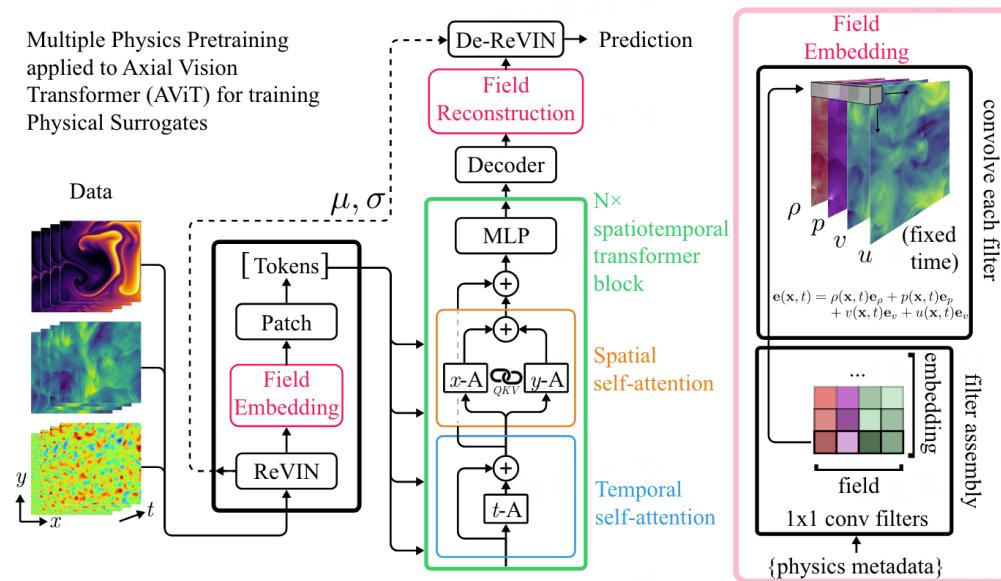
We pretrain a spatiotemporal transformer model on a large amount of trajectories (100,000 each) with uniformly sampled coefficients ($v \in [-3, 3]$, $\delta \in [10^{-3}, 1.]$) generated from the advection and diffusion equations while finetuning on restricted samples from advection-diffusion simulations.

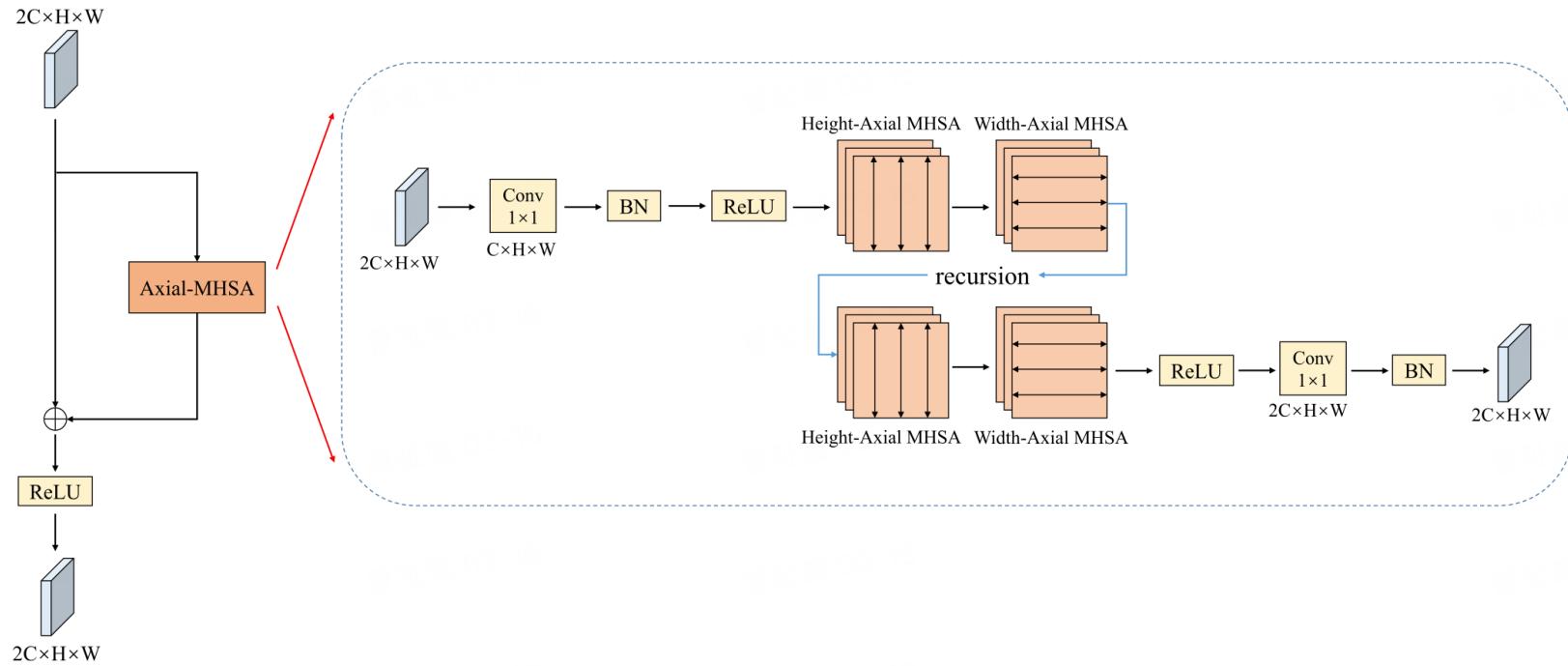


9.2 Axial Attention

For a (2+1)-dimensional system with $T \times H \times W$ tokens, conventional dense attention attends over all tokens simultaneously and has cost $O((HWT)^2)$. Axial attention instead performs a series of attention operations over each axis in turn, limiting the cost to $O(H^2 + W^2 + T^2)$.

While we perform attention on each axis independently, spatial attention utilizes one set of linear projections for both the height (y) and width (x) axes.





9.3 Field Embedding and Normalization

Normalization: To unify the magnitudes, we utilize reversible instance normalization . We compute the mean and standard deviation of each channel over the space- time dimensions and use them to normalize the input fields.

Embedding: the data is projected into a shared embedding space. This is the only component with weights that are unique to each source system. Given a system S with state variables $u(x, t), v(x, t), p(x, t) \in \mathbb{R}$, we project each sample point or “pixel” into a space of dimension D^{emb} :

$$\mathbf{e}(x, t) = u(x, t)\mathbf{e}_u + v(x, t)\mathbf{e}_v + p(x, t)\mathbf{e}_p \quad (31)$$

where \mathbf{e} are embedding vectors in $\mathbb{R}^{D^{\text{emb}}}$.

Prediction: Prediction are reconstructed by sequence of transposed conv blocks and inner-product projection

$$u(x, t + \Delta t) = \langle \mathbf{e}(x, t + \Delta t), \mathbf{r}_u \rangle$$

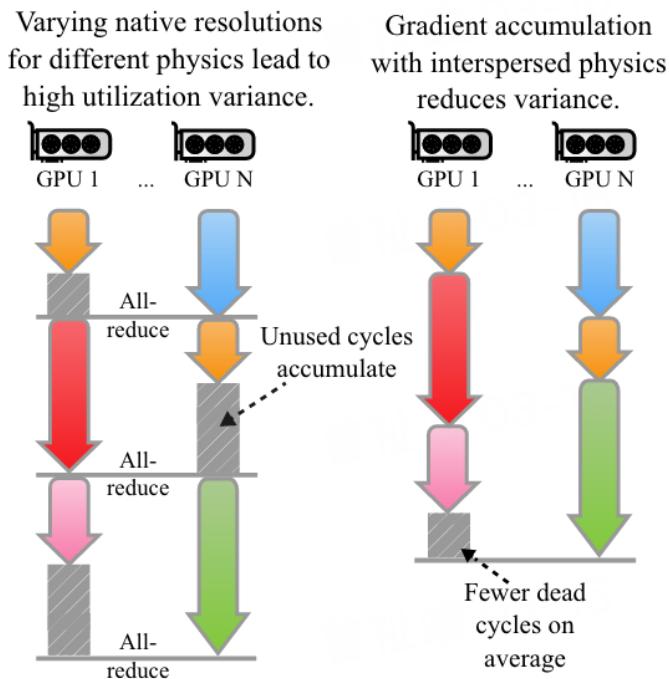
Position Biases: Using RoPE to embedding position inform. For periodic B.C we change the calculation of distance via periodic property.

9.4 Blancing Objectives

Our pretraining procedure operates on multiple levels of sampling. The task distribution varies in system S , spatial resolution N_S , and time resolution T_S .

Gradient accumulation utilizes multiple backward passes per synchronization step.

We therefore sample a single system S uniformly from \mathcal{S} for each *micro-batch*. With m micro-batches per synchronization step, we reduce the work-per-GPU variance σ_B^2 to $\frac{1}{m} \sigma_B^2$, significantly reducing the average lost cycles due to work discrepancies.



9.5 Experiment

Can large transformer models learn the dynamics of multiple physical systems simultaneously?

Table 1: NRMSE comparison between MPP-pretrained models and dedicated baselines. MPP-pretrained models learn multiple physical systems at least as well as standard baselines. Top performing within size range and overall are bolded. Dashes indicate precision not available. [†] While the PINN is much smaller, these models are fit per-example.

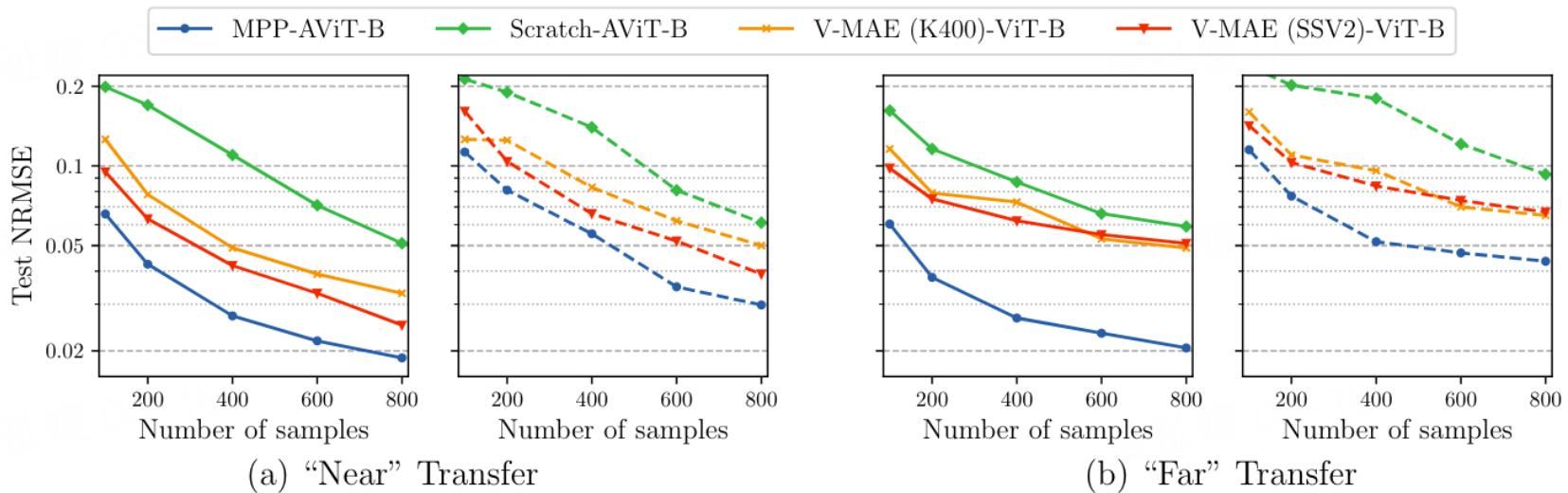
MODEL	#PARAM	SWE	DIFFRE2D	CNS M1.0	CNS M0.1
MPP-AViT-TI	7.6M	0.0066	0.0168	0.0442	0.0312
UNET	7.7M	0.083-	0.84-	0.4725	1.6650
FNO	466K	0.0044	0.12-	0.1685	0.2425
PINN	8.5K [†]	0.017-	1.6—	—	—
ORCA-SWIN-B	88M	0.0060	0.82-	—	—
MPP-AViT-B	116M	0.0024	0.0106	0.0281	0.0172
MPP-AViT-S	29M	0.0039	0.0112	0.0319	0.0213
MPP-AViT-L	409M	0.0022	0.0098	0.0208	0.0147

The improvement over the baselines is nearly an order of magnitude in NRMSE and the performance improves with scale.

Does MPP provide a finetuning advantage over existing spatiotemporal foundation models for new autoregressive prediction tasks?

We remove all compressible fluid data from the training corpus and pretrain on the three remaining spatiotemporal systems. We evaluate transfer to two specific compressible Navier-Stokes datasets:

- “Near”: $M = 0.1$, viscosity= 10^{-2} , Random Periodic Initial Conditions
- “Far”: $M = 1.0$, viscosity= 10^{-8} , Turbulent Initial Conditions



9.6 Broader Usage Of Pretrained Representation

Forcing Identification for Incompressible Navier-Stokes

We attempt to identify the constant forcing term used in the incompressible Navier-Stokes simulation from an input trajectory

Buoyancy for Incompressible Navier-Stokes

PDEArena includes an incompressible Navier-Stokes simulation with variable buoyancy. Since this set was not used during training, we take 1,000 randomly sampled trajectories for train, 100 for validation, and a further 1,000 for testing.

Table 3: RMSE for inverse problem tasks. Error from constant prediction included for context.

TRAINING	FORCING	BUOYANCY
MPP	$0.20 \pm .008$	$0.78 \pm .006$
SCRATCH	$0.43 \pm .012$	$0.77 \pm .005$
BEST CONSTANT	$1.00 \pm .000$	$0.77 \pm .000$

10 Flexible ML Models for Scientific Computing at Scale

10.1 Background

\mathcal{E} denotes set of parametrized PDEs describing a spatialtemporal dynamical system \mathcal{S} . We represent the snapshots of the observed physical field within the system by \mathcal{P} . A continuous variable for \mathcal{S} is defined as $u^{\mathcal{S}}(x, t) : \Omega \times [0, \infty] \rightarrow \mathbb{R}$. We discretized the system uniformly in space and time with resolution N_S and T_S then a snapshot $u_t^{\mathcal{S}} \in \mathbb{R}^{N_S}$ represent the value of state at time t.

1. Forward Problem: given equation \mathcal{E} , to predict the state of dynamical system \mathcal{S} at any coordinate (x, t) .
2. Inverse Problem: Given a set of PDEs $\tilde{\mathcal{E}}$, or in some case no explicit PDE information, along with observation comprising T_S snapshot $U_t^{\mathcal{S}} = [u_1^{\mathcal{S}}, \dots, u_t^{\mathcal{S}}]$, to infer the missing component(BC,IC,form of PDE) and reconstruct $\tilde{\mathcal{E}}$.

10.2 Method

10.2.1 Flexible Multiple Physics Pretraining

The framework is composed of a tokenization and encoding scheme, a state representation strategy for multichannel PDEs, an auto-regressive training mechanism, and a specialized optimization criterion.

Tokenization and Axis Encoding: Let \mathcal{F} denotes a physical field, which is a function of space x and time t with channel C , an axis encoder $E: \mathbb{R}^M \rightarrow \mathbb{R}^H$ transform the physical field with dimension M to hiddien field with dimension H .

$$h_{x,t} = E(\mathcal{F}(x, t))$$

The hidden vectors $h_{x,t}$ are flattened to form the input embeddings.

Channel-wise State Representation: For a PDE with multiple channels, we represent the state at t as a sequence grouped tokens T_t which includes C tokens

$$\mathbf{T}_t = [h_{1,t}; h_{2,t}; \dots; h_{C,t}]$$

This allows for the expression of interchannel and intra- channel relationships and dependencies on previous timesteps.

Auto-regressive Training with Shift-Right: For a sequence of states $\mathbf{T}_{1:T}$ the model predicts the next state token $\hat{\mathbf{T}}_{t+1}$

$$\hat{\mathbf{T}}_{t+1} = G(\mathbf{T}_{1:T})$$

Optimization via Batch-wise nRMSE:

Let $\hat{\mathcal{F}}$ denote the predicted physical field, the NRMSE of batch B is

$$\text{nRMSE}(B) = \frac{1}{|B|} \sqrt{\sum_{(\mathbf{x}, t) \in B} \left(\frac{\mathcal{F}(\mathbf{x}, t) - \hat{\mathcal{F}}(\mathbf{x}, t)}{\sigma_{\mathcal{F}}} \right)^2}$$

where $\sigma_{\mathcal{F}}$ is the std of this batch.

Decoder with Axis Numerical Embedding: We utilize the inverse of the axis encoding function $E^{-1}: \mathbb{R}^H \rightarrow \mathbb{R}^D$ to reconstruct the physical field from hidden representation.

$$\hat{\mathcal{F}}(\mathbf{x}, t) = E(\hat{\mathbf{h}}_{\mathbf{x}, t})$$

10.3 Physics-informed Reinforcement Learning

Let $\mathcal{D} = \{(x_i, t_i)\}_{\{i=1\}}^N$ be the dataset containing pairs of physical field data and their corresponding weakly annotated PDE captions, here x_i represents a frame of physical field data and t_i is the textual equation, boundary conditions, and other related annotations in natural language.

Physics-informed Scorer Model S: A clip-style model S takes a pair of x, t and output a similarity score $s \in \mathbb{R}$ indicating the match between the data and the textual information $S(x, t; \theta_S) \rightarrow s$

Pretrained Model G: Let G be a model that takes a sequence of recent time steps of physical fields $X_{(0:t)} = x_0, x_{\delta_t}, \dots, x_t$ as input and predict physical field \hat{x}_{t+1}

$$G(X_{(0:t)}; \theta_G) \rightarrow \hat{x}_{t+1}$$

Fine-Tuning with Physics-informed Rewards: During fine-tuning, input sequence $X_{(0:t)}$ and its caption t , the model G generates \hat{x}_{t+1} , the scorer model S evaluate the reward $r = S(\hat{x}_{t+1}, t; \theta_S)$, the reward signal is used to update θ_G

$$\theta_G^* = \arg \max_{\theta_G} \mathbb{E}(X_{(0:t)}), t \sim \mathcal{D}[r]$$

10.4 Implementation

1. For generator, we use LLaMA model. For the 2D and 3D tasks, we apply a parallel convolutional encoder with the axis numerical encoder for better high-frequency pattern capturability. For the decoder, we use the sum of logits from both transposed convolutional decoder and the weights-binned inverse axis numerical decoder.
2. Equation Rewriting: We apply mathematical identities to modify the equation, ensuring the core properties remain intact.
3. Linear Combination: For systems of equations, we derive new variants through linear combinations, enriching the dataset without altering the system's nature.
4. Symbol Substitution: We systematically swap variables with alternative symbols, such as replacing x with ξ , to maintain consistency and avoid ambiguity.
5. Physical Checking: A panel of GPT-4-based experts evaluates the augmented equations, filtering out those that do not align with physical principles.

A.2.1. BURGERS 1D

- Original form:

$$\partial_t u(t, x) + \partial_x(u^2(t, x)/2) = \nu/\pi \partial_{xx} u(t, x), \quad x \in (0, 1), t \in (0, 2],$$

$$u(0, x) = u_0(x), \quad x \in (0, 1),$$

- After augmented:

$$0.77 \int \left(\frac{\partial}{\partial t} v(t, x) + \frac{\partial}{\partial x} \frac{v^2(t, x)}{2} \right) dt = \frac{0.77\nu \int \frac{\partial^2}{\partial x^2} v(t, x) dt}{\pi}$$

$$0.73t v(0, x) = 0.73t v_0(x)$$

- Explanation: We replace u with v and ∂_t with $\frac{\partial}{\partial t}$. We integrate and multiply some factors on both sides of the equation at the same time.

A.2.3. CFD 1D

- Original form:

$$\partial_t \rho + \nabla \cdot (\rho \mathbf{v}) = 0,$$

$$\rho(\partial_t \mathbf{v} + \mathbf{v} \cdot \nabla \mathbf{v}) = -\nabla p + \eta \Delta \mathbf{v} + (\zeta + \eta/3) \nabla(\nabla \cdot \mathbf{v}),$$

$$\partial_t \left[\epsilon + \frac{\rho v^2}{2} \right] + \nabla \cdot \left[\left(\epsilon + p + \frac{\rho v^2}{2} \right) \mathbf{v} - \mathbf{v} \cdot \boldsymbol{\sigma}' \right] = 0,$$

- After augmented:

$$\varrho(t, x) \frac{\partial}{\partial x} \mathbf{w}(t, x) + \frac{\partial}{\partial t} \varrho(t, x) = 0$$

$$0.61 \left(\mathbf{w}(t, x) \frac{\partial}{\partial x} \mathbf{w}(t, x) + \frac{\partial}{\partial t} \mathbf{w}(t, x) \right) \varrho(t, x) = 0.61\eta \frac{\partial^2}{\partial x^2} \mathbf{w}(t, x) + 0.61 \left(\chi + \frac{\eta}{3} \right) \frac{\partial^2}{\partial x^2} \mathbf{w}(t, x) - 0.61 \frac{\partial}{\partial x} p(t, x)$$

$$\int \left(\frac{\partial}{\partial t} \left(\varepsilon + \frac{\mathbf{w}^2(t, x) \varrho(t, x)}{2} \right) + \frac{\partial}{\partial x} \left(-\tau \mathbf{w}(t, x) + \left(\varepsilon + \frac{\mathbf{w}^2(t, x) \varrho(t, x)}{2} + p(t, x) \right) \mathbf{w}(t, x) \right) \right) dt = 0$$

- Explanation: We replaced many symbols, such as replacing ∇ with ∂_t and Δ with $\frac{\partial^2}{\partial x^2}$,. We integrate and multiply some factors on both sides of the equation at the same time. We also swapped the order of some items, such as $\zeta + \eta/3$.

The text encoder leverages the pre-trained albert-math model, which is adept at processing LaTeX-encoded PDE captions due to its extensive training on a large corpus of LaTeX data.

For the physics encoder, we implement a streamlined one-channel Vision Transformer (ViT) model tailored to both 1D and 2D PDEs, and we adopt a large-batch contrastive learning approach, akin to the CLIP framework.

PDEs		FNO	PINNs	U-Net	MPP-L	OmniArch	+ PIRL	Lead %
1D	CFD	1.4100	-	2.6700	-	0.0981	0.0392	97.2 ↑
	ReacDiff	0.0005	0.2140	0.0026	-	0.0004	0.0002	64.0 ↑
	Advection	0.0091	0.8130	0.7760	-	0.0081	0.0045	49.9 ↑
	Burgers	0.0174	0.9450	0.3200	-	0.0067	0.0032	81.9 ↑
	Diff_Sorp	0.0017	0.2200	0.1500	-	0.0019	0.0006	62.2 ↑
2D	CFD	0.2060	-	1.0700	0.0178	0.0994	0.0153	14.0 ↑
	SWE	0.0044	0.0170	0.0830	0.0022	0.0031	0.0015	3.32 ↑
	ReacDiff	0.1200	1.6000	0.8400	0.0098	0.0993	0.0084	14.0 ↑
	NS-Incom	0.2574	-	1.1200	-	0.1494	0.0827	67.9 ↑
3D	CFD	0.3050	-	0.6150	-	0.6600	0.5072	66.3 ↓

Table 3: The Peroformace on Zero-Shot PDEs.

Methods	Shock	KH	OTVortex
FNO	0.7484	1.0891	0.5946
U-Net	1.6667	0.1677	0.4217
MPP	0.3243	1.3261	0.3025
OmniArch	0.2126	0.5432	0.1718

- **OTVortex:** The Orszag-Tang Vortex system is a compressible flow problem that generates highly complex vortex structures through the careful selection of initial conditions. The dataset includes one exsample, which is a 1024×1024 resolution physical field evolved over 101 time steps with a time interval of 0.01.
- **2D Shock:** Shock waves are characterized by abrupt changes in flow properties resulting from sudden discontinuities in fluid flow, such as rapid changes in pressure, temperature, and density. The dataset includes one exsample, which is also a 1024×1024 resolution physical field evolved over 101 time steps with a time interval of 0.01.
- **2D KH:** The Kelvin-Helmholtz instability is a fluid instability that occurs at the interface between two fluid layers with different velocities or densities. This dataset consists of seven exsamples generated based on different parameters M , dk , and Re . Each is a 1024×1024 resolution physical field evolved over 51 time steps with a time interval of 0.1. We conducted experiments on all samples and averaged the results.