

# AI in the Sciences and Engineering HS 2025: Lecture 8

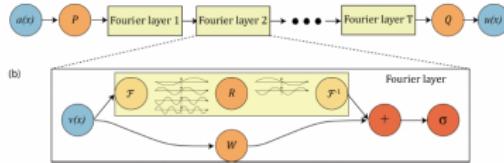
Siddhartha Mishra

Computational and Applied Mathematics Laboratory (CamLab)  
Seminar for Applied Mathematics (SAM), D-MATH (and),  
ETH AI Center (and) Swiss National AI Institute (SNAI) ,  
ETH Zürich, Switzerland.

# What we have learnt so far ?

- ▶ AIM: Learn PDEs using Deep Neural Networks
- ▶ Operator Learning: Learn the PDE Solution Operator from data.
- ▶ Examples: FNO, CNO.

# Fourier Neural Operators: Li et al, 2020



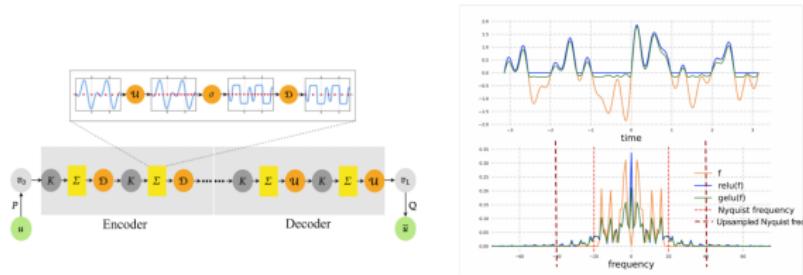
- ▶ Translation invariant Kernel  $K(x, y) = K(x - y)$
- ▶ Use Fourier and Inverse Fourier Transform to define the KIO:

$$\int_D K(x - y)v(y)dy = \mathcal{F}^{-1}(\mathcal{F}(K)\mathcal{F}(v))(x)$$

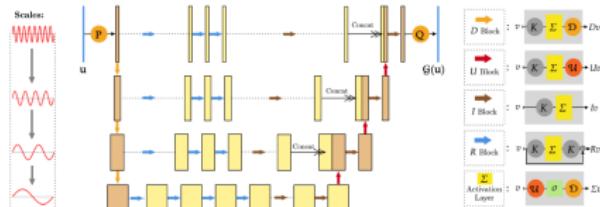
- ▶ Parametrize Kernel in Fourier space.
- ▶ Nonlocal Convolutions + Local Activations:  $[\mathcal{P}\sigma] = \sigma$
- ▶ Fast implementation through FFT
- ▶ Widely used + Backed by solid Theory <sup>1</sup>

<sup>1</sup>Kovachki, Lanthaler, SM 2020, Lanthaler, Molinaro, Hadorn, SM 2022

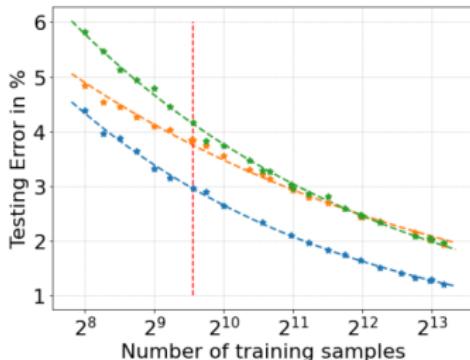
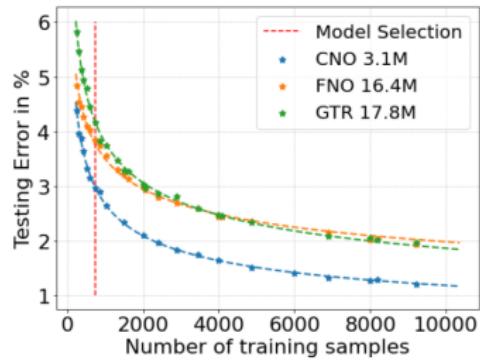
# Convolutional Neural Operators: Raonic, SM et al, 2023



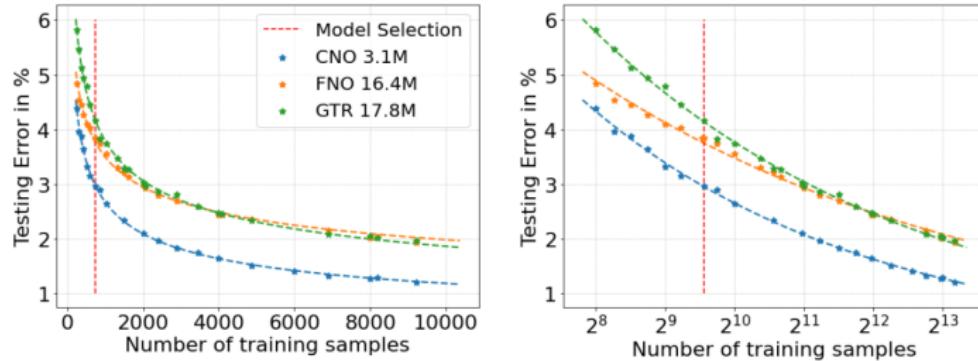
- ▶ I: Use **Convolutional Kernels in Physical space**
- ▶ II: Modulated Nonlocal activations for Alias-free processing.
- ▶ CNO instantiated as a modified **Operator UNet**



# How does CNO/FNO Scale ?

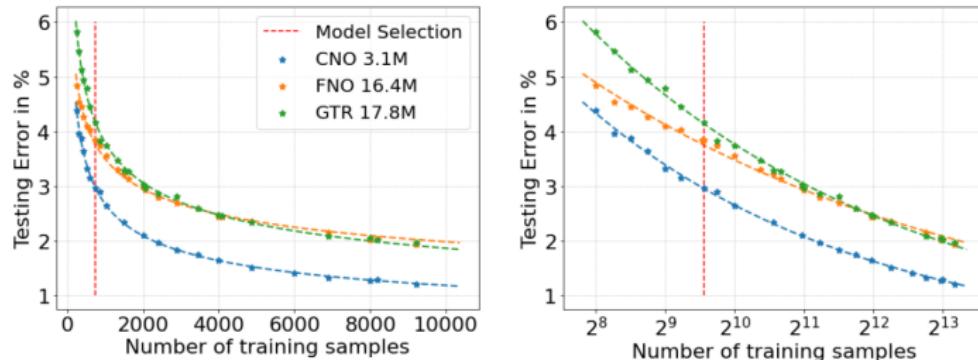


# How does CNO/FNO Scale ?



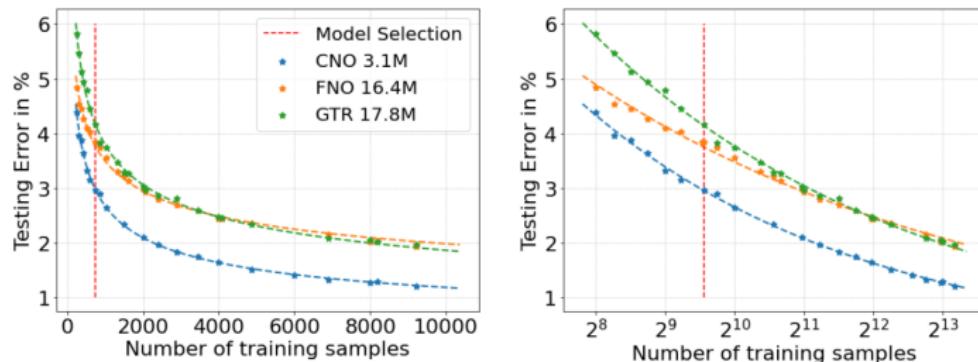
- ▶ Success is a curve, not a point !!

# How does CNO/FNO Scale ?



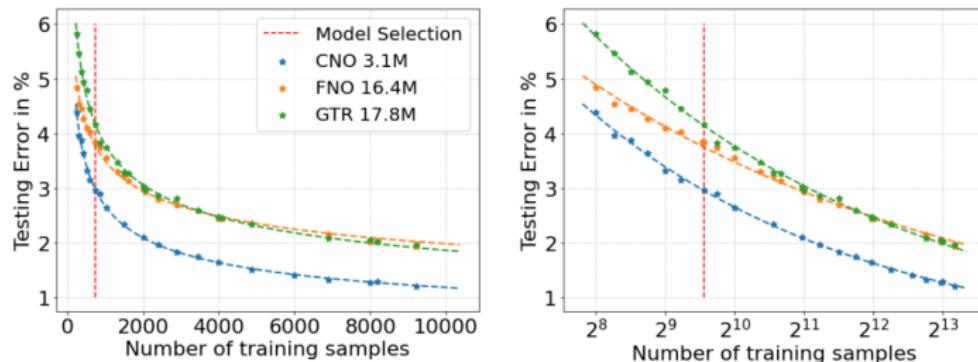
- ▶ Success is a curve, not a point !!
- ▶ Models **Scale** with sample size:  $\mathcal{E} \sim N^{-\alpha}$  but with  $\alpha$  **small**
- ▶ Theory: Lanthaler, SM, Karniadakis, De Ryck, SM,

# How does CNO/FNO Scale ?



- ▶ Success is a curve, not a point !!
- ▶ Models **Scale** with sample size:  $\mathcal{E} \sim N^{-\alpha}$  but with  $\alpha$  **small**
- ▶ Theory: Lanthaler, SM, Karniadakis, De Ryck, SM,
- ▶ ML models require **Big Data**:  $\mathcal{O}(10^3) - \mathcal{O}(10^4)$  training samples per Task
- ▶ Very Difficult to obtain Data for PDEs.

# How does CNO/FNO Scale ?



- ▶ Success is a curve, not a point !!
- ▶ Models **Scale** with sample size:  $\mathcal{E} \sim N^{-\alpha}$  but with  $\alpha$  **small**
- ▶ Theory: Lanthaler, SM, Karniadakis, De Ryck, SM,
- ▶ ML models require **Big Data**:  $\mathcal{O}(10^3) - \mathcal{O}(10^4)$  training samples per Task
- ▶ Very Difficult to obtain Data for PDEs.
- ▶ Can models **Scale** better ?

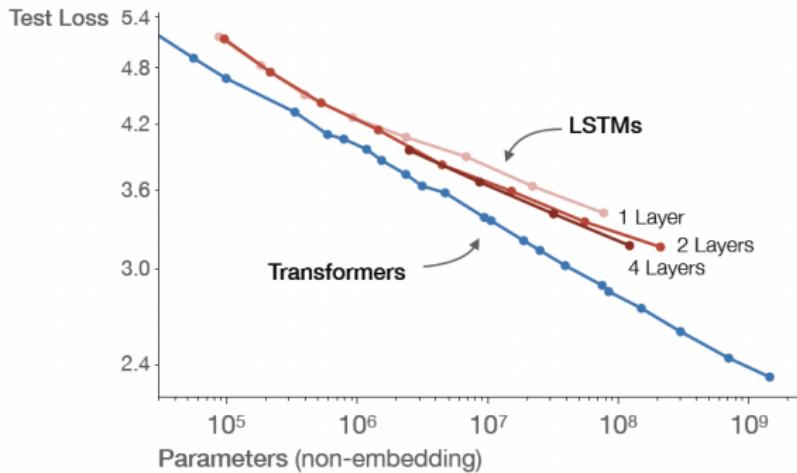
# What Architectures Scale the best ?

# What Architectures Scale the best ?

- ▶ What happens in [Language Modeling](#) ?
- ▶ Results of [Kaplan et. al., 2020.](#)

# What Architectures Scale the best ?

- ▶ What happens in **Language Modeling** ?
- ▶ Results of **Kaplan et. al., 2020.**



# Move Towards Transformers

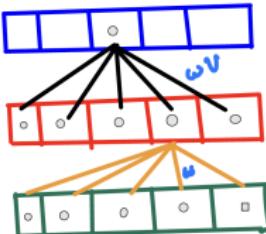


# Attention is all you need !!

- ▶ Introduced by Vaswani et. al, 2017
- ▶ Given  $K$ -inputs  $v \in \mathbb{R}^{K \times n}$  of  $n$ -features:
- ▶ **Self-Attention:**  $\mathbb{A} : \mathbb{R}^{K \times n} \mapsto \mathbb{R}^{K \times n}$  such that  $u = \mathbb{A}v$ :

$$u_k = W \sum_{j=1}^K \text{softmax}_k \left( \frac{\langle Qv_k, Kv_j \rangle}{\sqrt{m}} \right) Vv_j, (\text{softmax}(w))_i = \frac{e^{w_i}}{\sum_{\ell=1}^L e^{w_\ell}}$$

- ▶ **Query**  $Q \in \mathbb{R}^{m \times n}$ , **Key**  $K \in \mathbb{R}^{m \times n}$ , **Value**  $V \in \mathbb{R}^{m \times n}$ .
- ▶ Output Matrix  $W \in \mathbb{R}^{n \times m}$



# Multi-Head Attention

- ▶ Multi-Head Self-Attention:

$$u_k = \sum_{h=1}^H W_h \sum_{j=1}^K \text{softmax}_k \left( \frac{\langle Q_h v_k, K_h v_j \rangle}{\sqrt{m}} \right) V_h v_j$$

# Attention as a Neural Operator: I

- ▶ Let  $v \in C(D, \mathbb{R}^n)$  be the input function.
- ▶  $x_k \in D$  be sampling points on a Regular Grid with size  $\Delta$
- ▶ Apply **Self-Attention** to **Tokens**  $v_k = v(x_k)$ :

$$u_k = W \sum_{j=1}^K \text{softmax}_k \left( \frac{\langle Qv_k, Kv_j \rangle}{\sqrt{m}} \right) Vv_j, (\text{softmax } (w))_i = \frac{e^{w_i}}{\sum_{\ell=1}^L e^{w_\ell}}$$

- ▶ Passing to the limit as  $\Delta \rightarrow 0$  yields

$$u(x) = \mathbb{A}(v)(x) = W \int_D \frac{e^{\frac{\langle Qv(x), Kv(y) \rangle}{\sqrt{m}}}}{\int_D e^{\frac{\langle Qv(z), Kv(y) \rangle}{\sqrt{m}}} dz} Vv(y) dy.$$

# Attention as a Neural Operator: II

- ▶ **Operator Attention**  $\mathbb{A} : C(D, \mathbb{R}^n) \mapsto C(D, \mathbb{R}^n)$  with

$$u(x) = \mathbb{A}(v)(x) = W \int_D \frac{e^{\frac{\langle Qv(x), Kv(y) \rangle}{\sqrt{m}}}}{\int_D e^{\frac{\langle Qv(z), Kv(y) \rangle}{\sqrt{m}}} dz} Vv(y) dy.$$

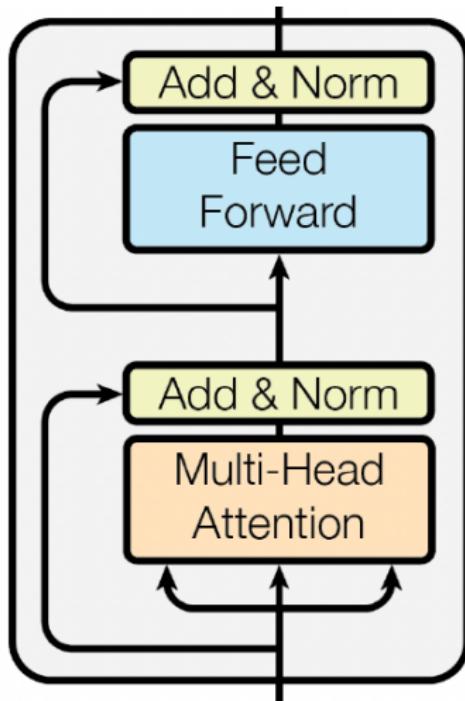
- ▶ Can be interpreted as as a **Nonlinear Kernel Neural Operator**:

$$\mathbb{A}(v)(x) = \int_D K(v(x), v(y)) v(y) dy.$$

- ▶ In contrast, FNO/CNO etc are **Linear Kernel Neural Operators**:

$$\mathbb{C}(v)(x) = \int_D K(x, y) v(y) dy = \int_D K(x - y) v(y) dy.$$

# Additional Ingredients



# Normalization ?

- ▶ Layer Normalization
- ▶ LayerNorm as a function:

$$LN : \mathbb{R}^n \mapsto \mathbb{R}^n, \quad LN(v) = \alpha \odot \frac{v - \mu_v}{\sigma_v} + \beta,$$

$$\mu_v = \frac{1}{n} \sum_{j=1}^n v_j, \quad \sigma_{bv}^2 = \frac{1}{n} \sum_{j=1}^n (v_j - \mu_{bv})^2$$

- ▶ Straightforward to realize LayerNorm as a pointwise operator

$$(LNv)(x) = LN(v(x))$$

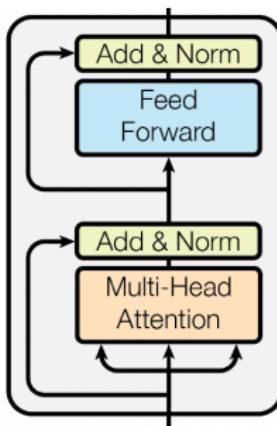
# Feed Forward ?

- ▶ Use **MLPs** pointwise
- ▶ Operator interpretation is:

$$MLP(v)(x) = \overline{W}\sigma(Wv(x) + B), \quad v \in \mathbb{R}^n$$

- ▶ with Weight Matrices  $W \in \mathbb{R}^{d \times n}$ ,  $\overline{W} \in \mathbb{R}^{n \times d}$
- ▶ and Bias vector  $B \in \mathbb{R}^d$

# An Operator Transformer Block



- ▶ A sequence of Operators of the form:
- ▶ **Multi-head Self-Attention**:  $\bar{u} = MSA(v)$
- ▶ **Residual + LayerNorm**:  $\hat{u} = LN(v + \bar{u}, \alpha_1, \beta_1)$
- ▶ **MLP**:  $u' = MLP(\hat{u})$
- ▶ **Residual + LayerNorm**:  $u = LN(u' + \hat{u}, \alpha_2, \beta_2)$

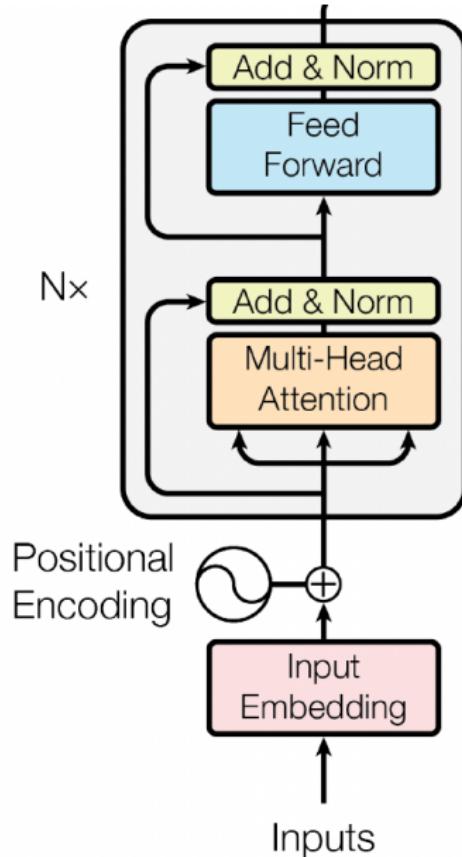
# Issue: Where is Position ?

- ▶ All operators have Permutation Invariance.
- ▶ No sense of absolute or relative Position
- ▶ Have to include Positional encodings of form:

$$\hat{v}_k = \overline{W}_v \sigma(W_v v_k) + \overline{W}_e \sigma(W_e e_k),$$

- ▶ With  $e_k \in \mathbb{R}^K$  is the  $k$ -th unit vector
- ▶ Weight matrices  
 $W_e \in \mathbb{R}^{d \times K}, \overline{W}_e \in \mathbb{R}^{n \times d}, W_v \in \mathbb{R}^{\bar{d} \times n}, \overline{W}_v \in \mathbb{R}^{n \times \bar{d}}$
- ▶ Fixed Sine-Cosine position encodings also work.

# Final version of a Transformer Block



# Caveat: Computational Complexity

- ▶ Computational Cost is Quadratic in # (Tokens) !!

$$\text{Compute} \sim \mathcal{O}(mnK^2)$$

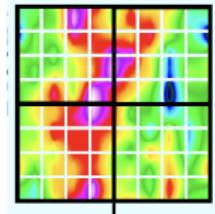
- ▶ With  $K$ - Input Length,  $n$  Input features and  $m$  hidden dimension.
- ▶ But lots of possible Parallelism in Computation
- ▶  $\mathcal{O}(1)$  sequential operations.
- ▶  $\mathcal{O}(1)$  Path Length.
- ▶ Nevertheless, Infeasible for 2 or 3-d inputs.

# Possible Solution

- ▶ Vision Transformers (ViT) of Dosovitskiy et. al.
- ▶ Key Ingredient: Patching (Patchification):
- ▶ Given Image on resolution  $H \times W$ , i.e.  $v \in \mathbb{R}^{H \times H \times C}$
- ▶ Divide into  $N = \frac{H^2}{p^2}$  patches, each of size  $p \times p$
- ▶  $v \sim [v_1, v_2, \dots, v_N]$  with  $v_k \in \mathbb{R}^{p^2 \times C}$ ,  $1 \leq k \leq N$
- ▶ Introduce Patch Embeddings  $E \in \mathbb{R}^{n \times (p^2 \cdot C)}$
- ▶ Input is a Sequence of Patch Tokens:

$$\hat{v} = [Ev_1, Ev_2, \dots, Ev_N]$$

- ▶  $N$  Tokens, each with  $n$  features are fed into a transformer !!



# Operator Version

- ▶ Patch Formation + Embeddings can be written as operator:

$$\hat{E}(v)(x) = \sum_{k=1}^N F \left( \int_{D_k} W(x)v(x)dx \right) \mathbb{I}_{D_k}(x).$$

- ▶ With  $D = \cup_{k=1}^N D_k$  non-overlapping and equal partition.
- ▶ With learnable  $F \in \mathbb{R}^{n \times C}$
- ▶ Weight Function  $W(x) = \sum_{1 \leq i, j \leq H} W_{ij} \delta_{z_{ij}}$
- ▶ With uniform grid points  $z_{ij}$  and discrete weights defined by,

$$\begin{aligned} W_{ij} &= \omega_{ij} \quad \text{if } 1 \leq i, j \leq p \\ &= \omega_{i \bmod p, j \bmod p}, \quad \text{otherwise.} \end{aligned}$$

# ViT operator Block

- ▶ For  $D \subset \mathbb{R}^2$  + input  $v \in C(D, \mathbb{R}^C)$
- ▶ A sequence of Operators of the form:
- ▶ Patch Embeddings+ Positional Encoding:  $\hat{v} = \hat{E}(v) + E_{pos}(v)$
- ▶ LayerNorm + MSA+ Residual:  $\bar{u} = \hat{v} + MSA(LN(v))$
- ▶ LayerNorm + MLP+ Residual:  $u = \bar{u} + MLP(LN(\bar{u}))$

# Computational Complexity

- ▶ Given an Image at resolution  $H \times W$
- ▶ Standard Transformer needs

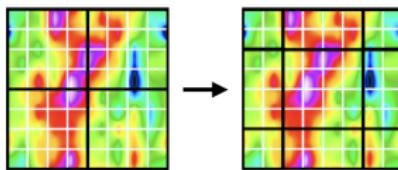
$$\text{Compute} \sim \mathcal{O}((HW)^2)$$

- ▶ ViT needs

$$\text{Compute} \sim \mathcal{O}\left(\frac{(HW)^2}{p^4}\right)$$

- ▶ Still not scalable for small patch size  $p$

# Another Idea: Windowed Attention of Liu et. al., 2022

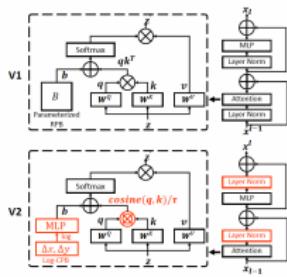


## ► Windowed Attention:

$$\mathbb{A}_W(v)(x) = W \int_{D_{q_x}^{w,\ell}} \frac{e^{\frac{\langle Qv(x), Kv(y) \rangle}{\sqrt{m}}}}{\int_{D_{q_x}^{w,\ell}} e^{\frac{\langle Qv(z), Kv(y) \rangle}{\sqrt{m}}} dz} Vv(y) dy.$$

- With  $M$ -Windows, Compute  $\sim \mathcal{O}\left(\frac{HWM^2}{p^2}\right)$
- How to Attend to Patches outside the Window ?
- Solution: Window shifts across Layers !!

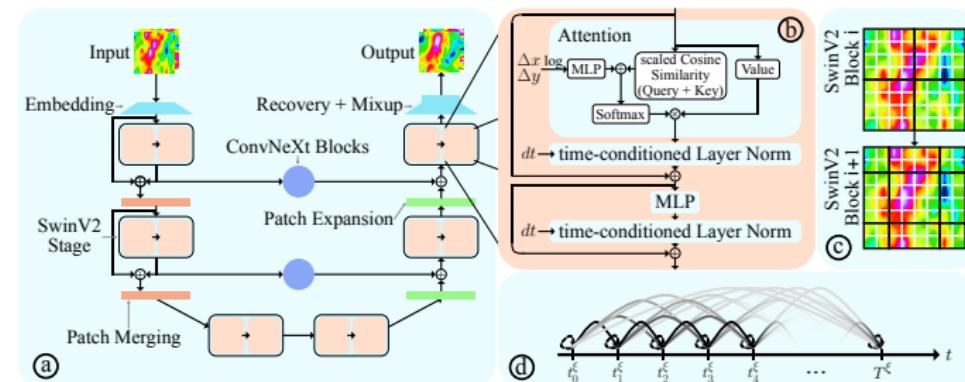
# Modifications for Scalability



- ▶ Replace scaled dot product with Scaled Cosine
- ▶ Use MLPs on Relative Position Coordinates to generate positional encodings.

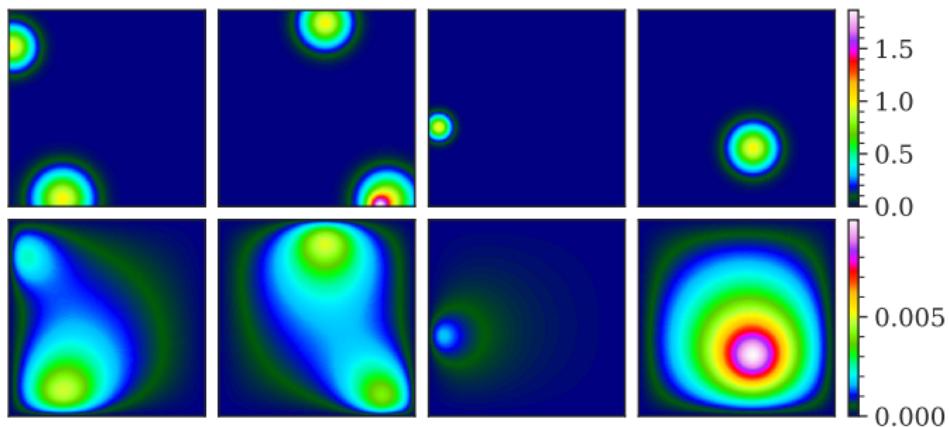
# scOT: scalable Operator Transformer

- ▶ Proposed in Herde,SM et. al., 2024.

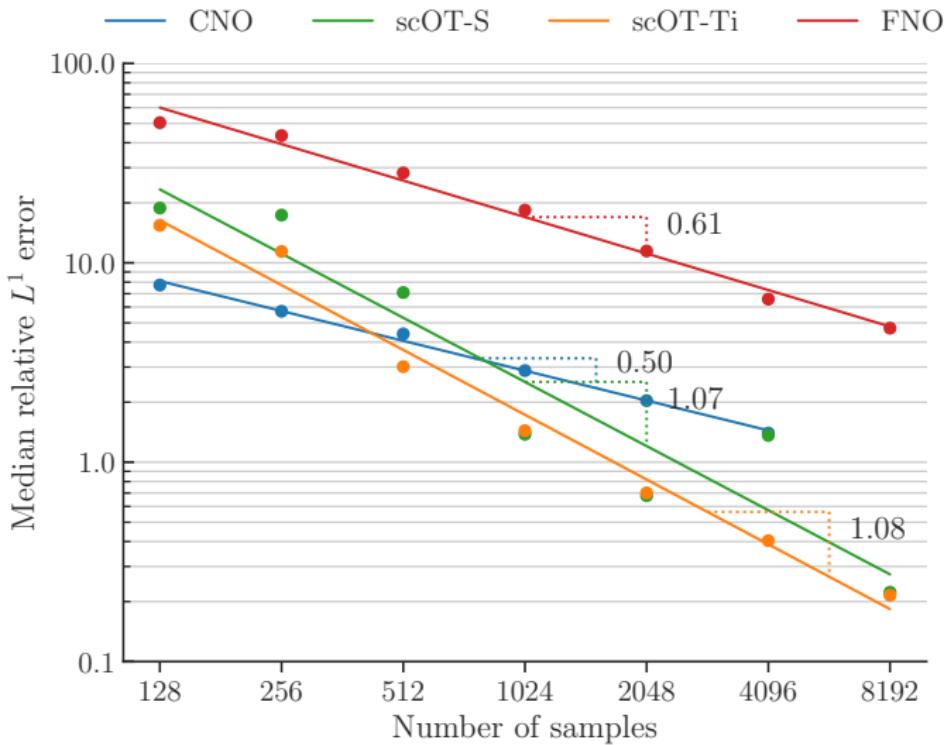


- ▶ **scalable Operator Transformer:** Based on SWIN Attention.
- ▶ Universal Approximator of Continuous Operators

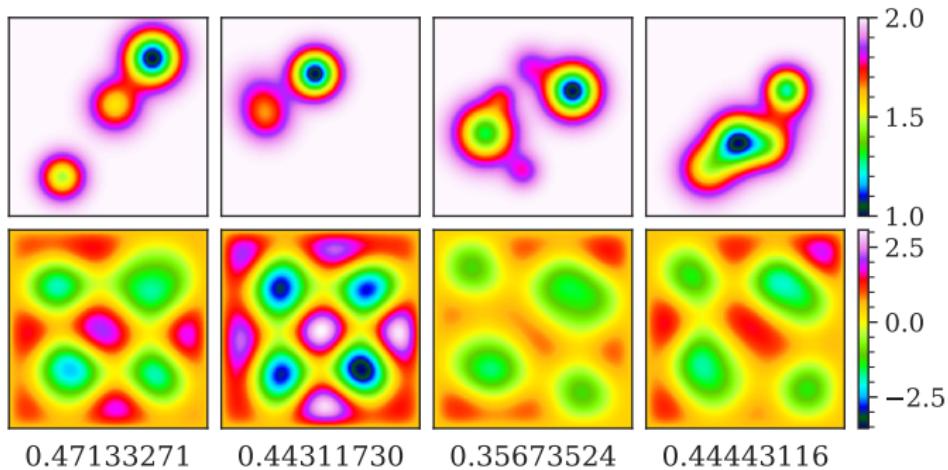
# Poisson with Gaussian Sources

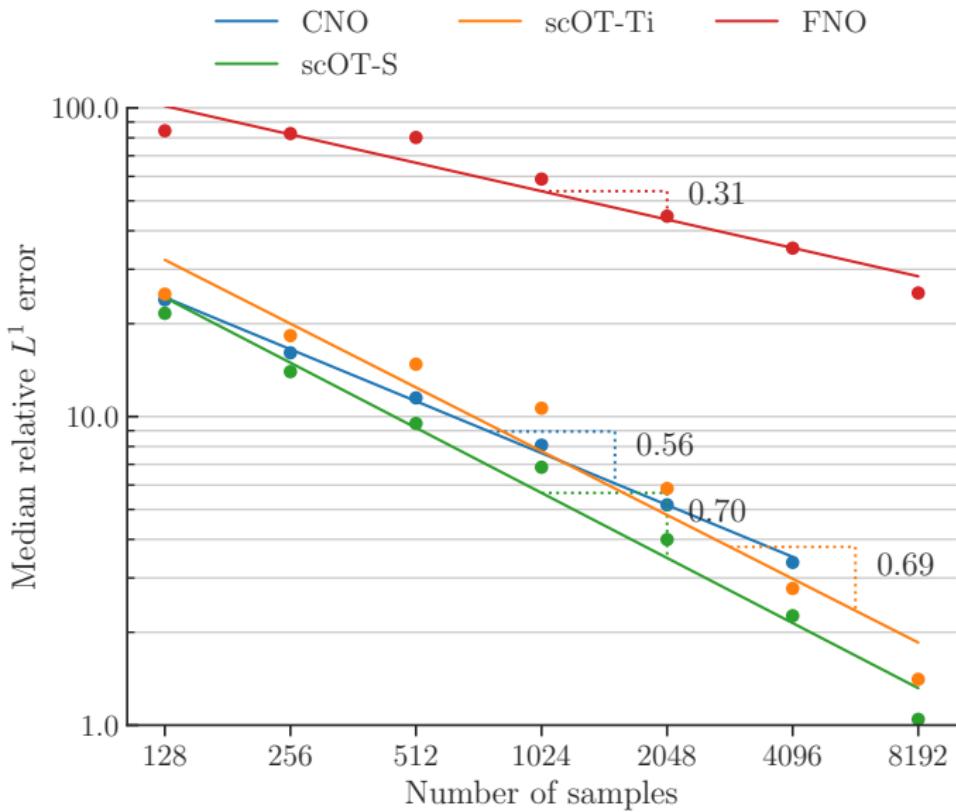


# Poisson with Gaussian Sources



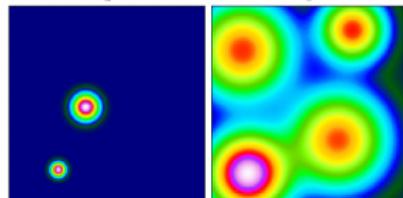
# Helmholtz



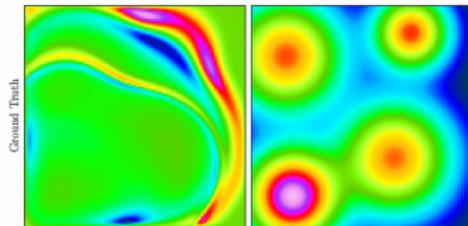


# Wave Equation

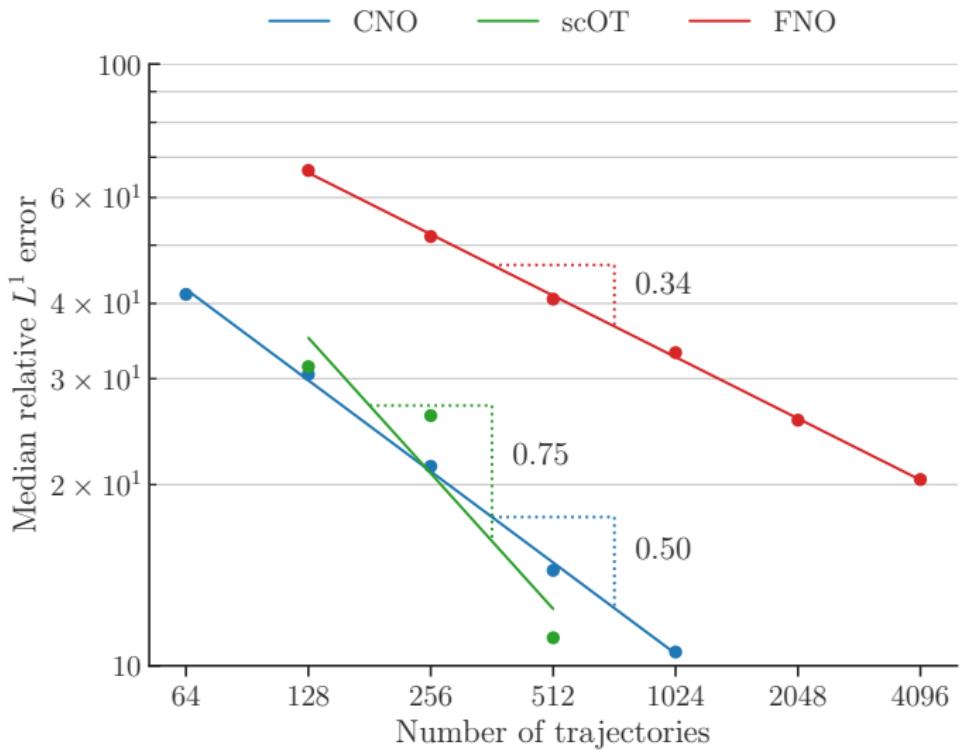
- Input ( $T = 0$ ):



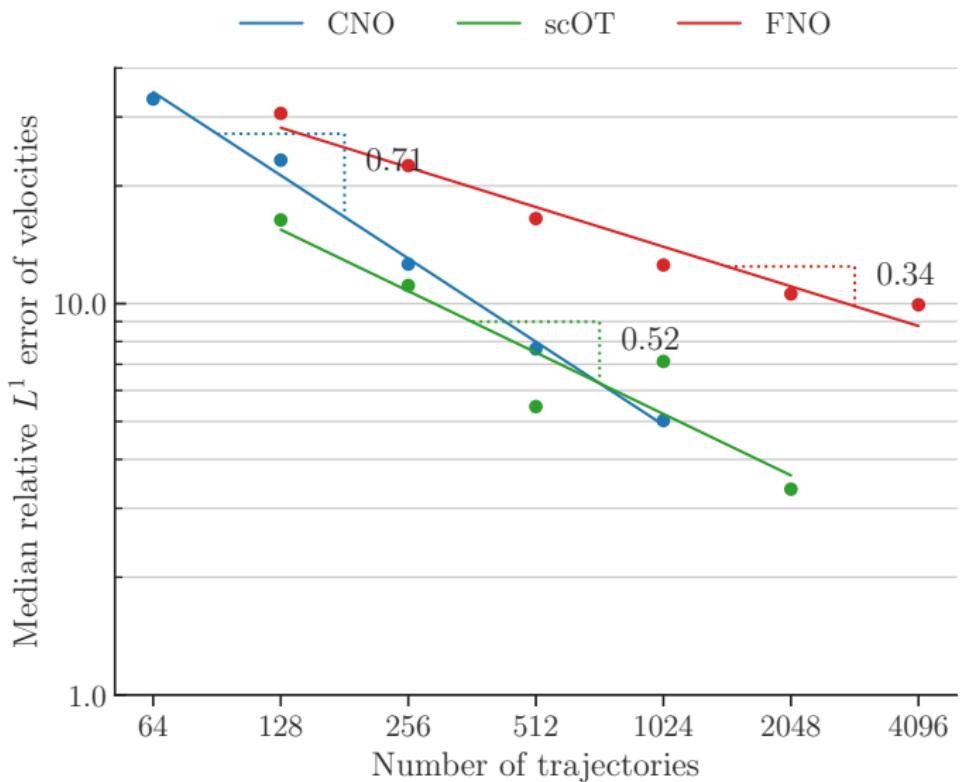
- Output t ( $T = 1$ ):



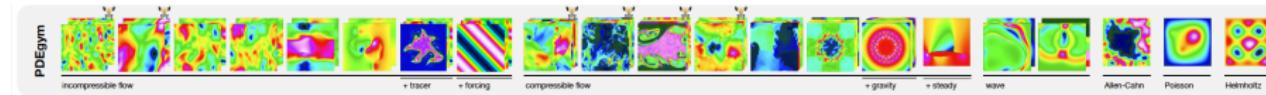
# Wave Equation



# Navier-Stokes: Shear Layer



# CNO vs. scOT vs. FNO



- ▶ Performance on PDEgym Dataset.
- ▶ Set of 15 PDE operator learning tasks.

Model	Median (EG)	Mean (AG)	N(EG)	N(AG)
FNO	1	1	-	-
CNO	4.6	2.61	5	6
scOT	5.4	2.56	4	8

- ▶ Sample Complexity remains a problem !!