

AI in the Sciences and Engineering 2024: Lecture 16

Siddhartha Mishra

Seminar for Applied Mathematics (SAM), D-MATH (and),
ETH AI Center,
ETH Zürich, Switzerland.

What you learnt so far

- ▶ Operator learning: Given Abstract PDE: $\mathcal{D}_a(u) = f$
- ▶ Learn **Solution Operator**: $\mathcal{G} : \mathcal{X} \mapsto \mathcal{Y}$ with $\mathcal{G}(a, f) = u$
- ▶ Approximate with **Operator Learning** Algorithms:
 - ▶ CNN/UNet
 - ▶ DeepONet
 - ▶ FNO
 - ▶ CNO
- ▶ We focus on approximating **Time-dependent problems**

Time-dependent PDEs

- ▶ Of the Abstract form:

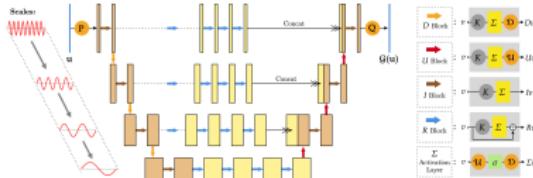
$$u_t + \mathcal{L}(t, x, u) = 0, \quad u(0) = \bar{u}.$$

- ▶ **Solution operator**: $\mathcal{S} : (0, T) \times \mathcal{X} \mapsto \mathcal{X}$; $\mathcal{S}(t, \bar{u}) = u(t)$
- ▶ Generated data is the form of **Trajectories**:

$$(u(0), u(t_1), u(t_2), \dots, u(T)) = (\bar{u}, \mathcal{S}(t_1, \bar{u}), \mathcal{S}(t_2, \bar{u}), \dots, u(T))$$

- ▶ **Learning Task**:
- ▶ Given \bar{u} + BC: generate the solution trajectory $u(t)$, for all $t \in (0, T]$
- ▶ **Direct** or **Autoregressive** Evaluation of Neural operators is not efficient.

Lead Time Conditioning



- ▶ Lead Time as an Input Channel
- ▶ CNO $(\bar{t}, u(t)) \approx \mathcal{S}(\bar{t}, u(t)) = u(t + \bar{t})$.
- ▶ Add Conditional Normalizations after each layer !!

$$\mathcal{N}(w) = g_N(t) \odot \frac{w - \mathbb{E}(w)}{\sqrt{\text{Var}(w) + \epsilon}} + h_N(t),$$

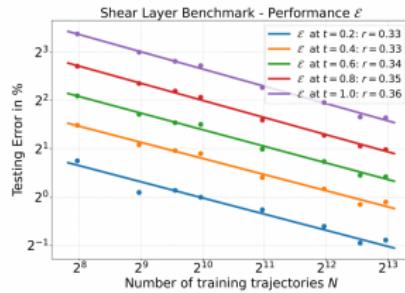
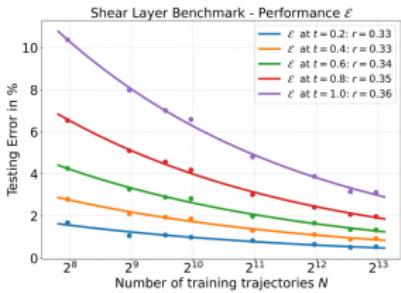
- ▶ g_N, h_N are MLPs in general.
- ▶ Instance, Batch, Layer Normalizations.

Training Strategies

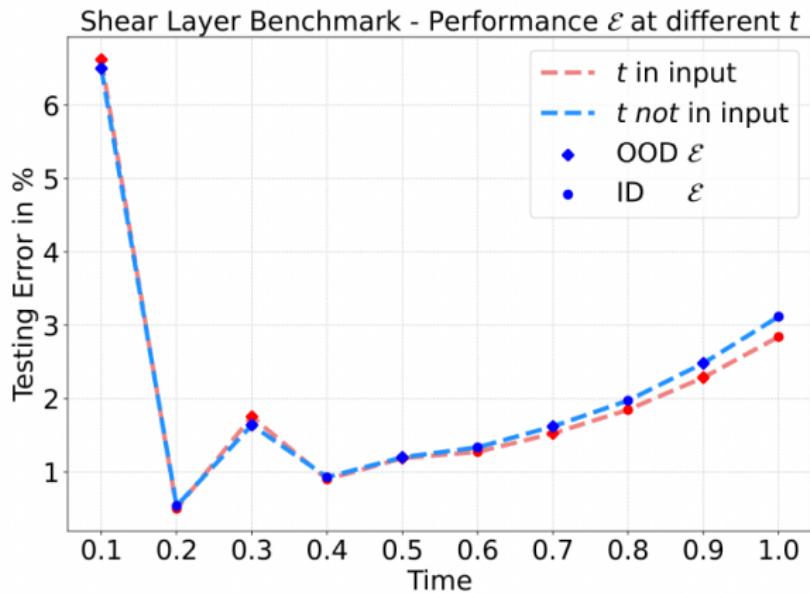
- ▶ One at a Time training based on:
- ▶ Input-Target Pairs: $\bar{u}, \mathcal{S}(\bar{u}) = u(\bar{u})$
- ▶ For $t_K = T$, K training samples per trajectory.
- ▶ all2all training based on:
- ▶ Input-Target Pairs: $u(t_i), \mathcal{S}(t_j - t_i, u(t_i)) = u(t_j), \forall i < j$
- ▶ $\frac{K^2+K}{2}$ training samples per trajectory !!
- ▶ Inference is Direct or Autoregressive
- ▶ Multiple possibilities for Autoregressive Rollouts
- ▶ Evaluation at any time $t > 0$ including Out-of-distribution times.

Results for Shear Layer

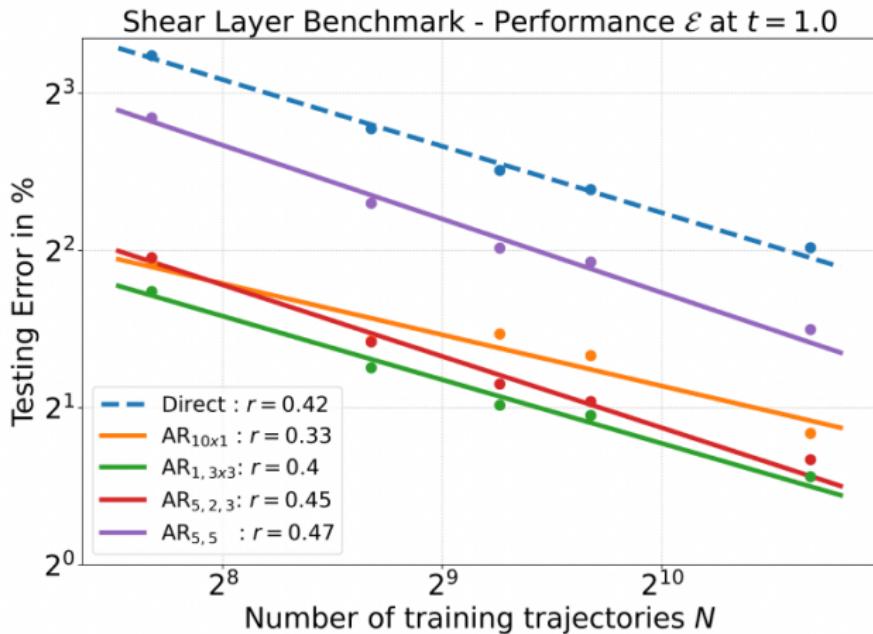
Error vs. Time



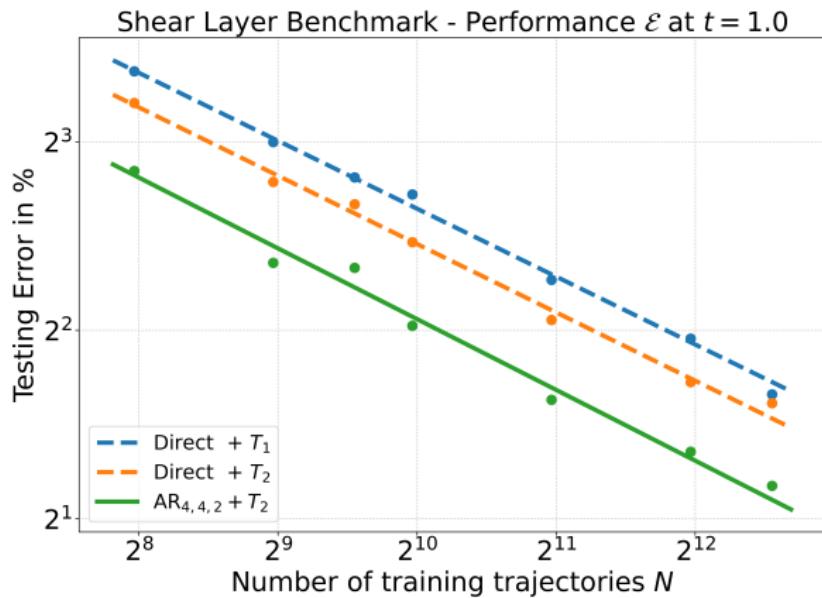
Results at OOD time levels.



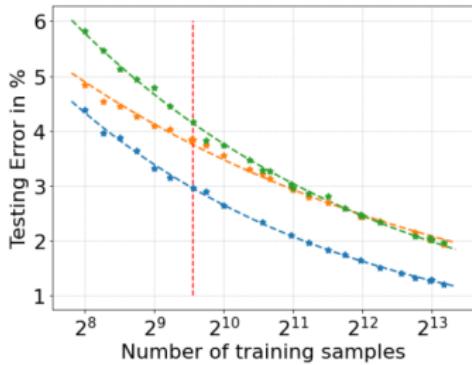
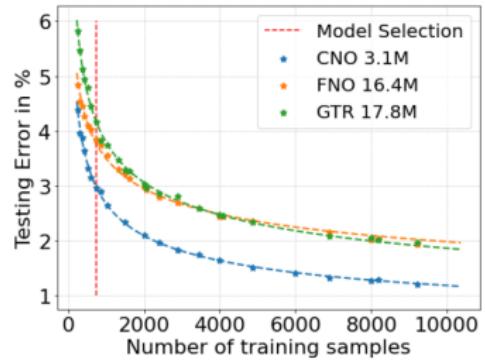
Results for Different Strategies I.



Results for Different Strategies II.



Key is Scaling



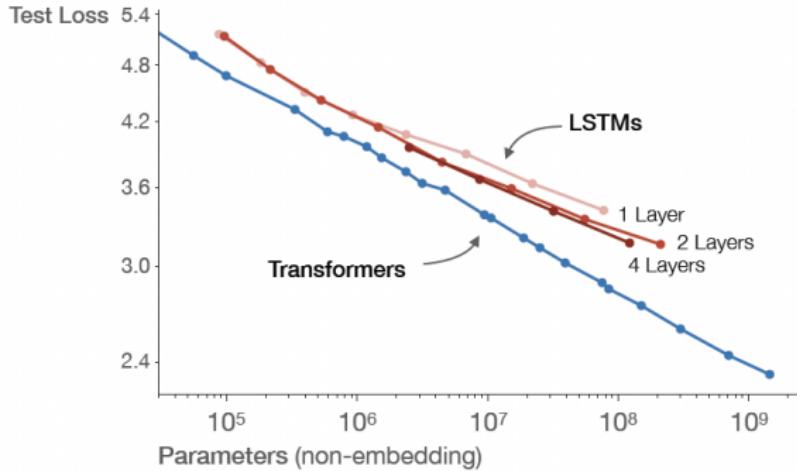
What Architectures Scale the best ?

What Architectures Scale the best ?

- ▶ What happens in **Language Modeling** ?

What Architectures Scale the best ?

- ▶ What happens in Language Modeling ?



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 $W \in \mathbb{R}^{n \times m}$
- ▶ **Permutation Invariance**

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 Δ
- ▶ 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

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

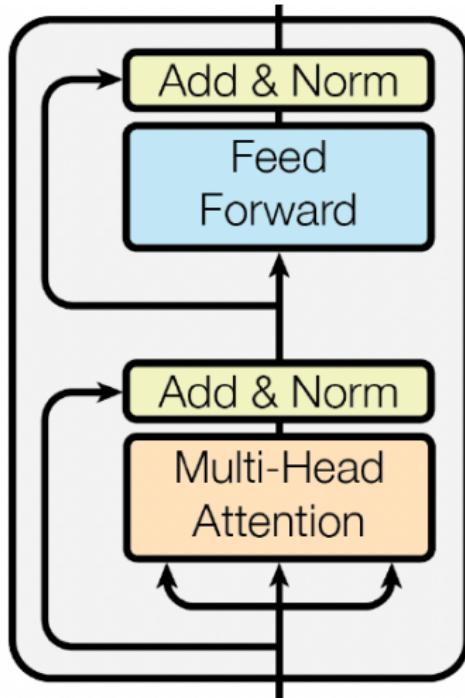
- ▶ Can be interpreted as as a **Nonlinear Kernel Neural Operator**:

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

- ▶ 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

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

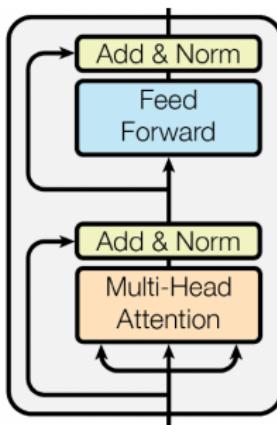
Feed Forward ?

- ▶ 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**: $u = LN(u' + \bar{u}, \alpha_2, \beta_2)$

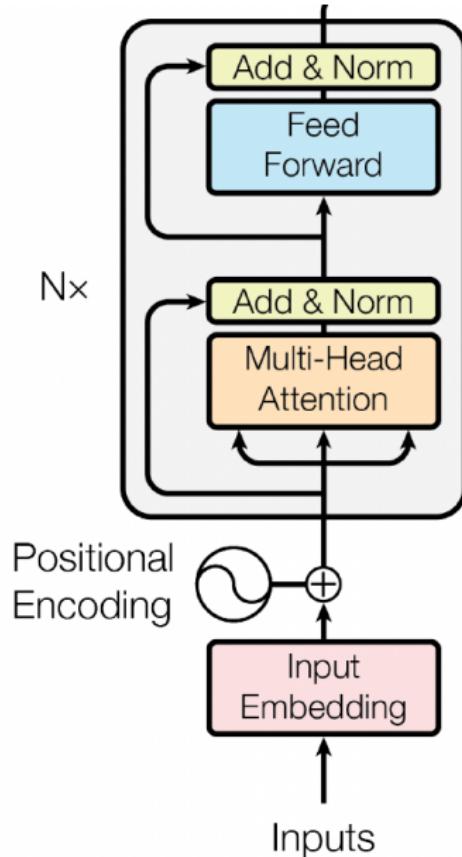
Issue: Where is Position ?

- ▶ 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,

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

Issue: Where is Position ?

- ▶ Need to add **Positional Encodings**
- ▶ Input:
- ▶ Sequence of Patch Tokens + Learnable Positional Encodings

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

- ▶ With $E_{pos} \in \mathbb{R}^{n \times N}$
- ▶ Can be viewed as an operator:

$$E_{pos}v(x) = \sum_{k=1}^N w_k \mathbb{I}_{D_k}(x), \quad w_k \in \mathbb{R}^n, \quad \forall 1 \leq k \leq N.$$

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

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

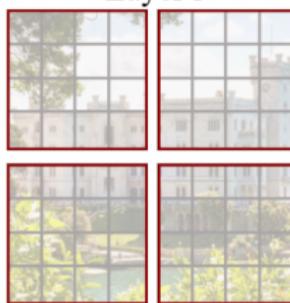
- ▶ ViT needs

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

- ▶ Still not scalable for small patch size p

Another Idea: Windowed Attention

- ▶ Introduced in Liu et. al.
- ▶ Use **Windowed Attention**:



- ▶ With M -Windows,

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

Operator version of Windowed Attention

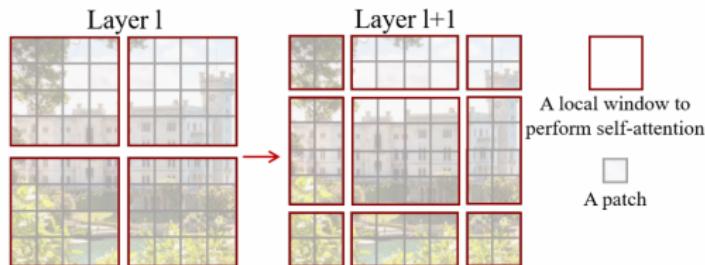
- ▶ For layer ℓ , Assume $D = \bigcup_{q=1}^M D_q^{w,\ell}$
- ▶ With non-overlapping Windows.
- ▶ **Windowed Attention** is instantiated as Operator:

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

- ▶ Where $1 \leq q_x \leq M$ such that $x \in D_{q_x}^{w,\ell}$

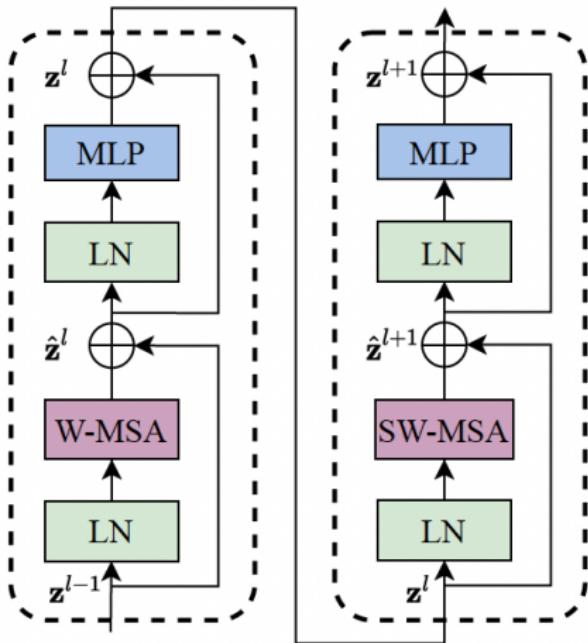
Shifting the Windows

- ▶ How to Tokens outside the Window ?
- ▶ Solution: Window shifts across Layers !!

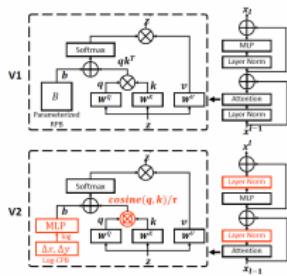


- ▶ Reduces Path Length across Tokens.

Swin Transformer Block



Modifications for Scalability



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

scOT: scalable Operator Transformer

