A General Neural Operator Transformer for Operator Learning

### Self Attention

· Self Alkention is just a sequence to sequence operation:



Weighted Aug over all input  $y_i = \sum_i w_{ij} x_j$ 

· Weights are not learnable parameters: Wij = xixj ⇒ Wij = softman (Wij) -00 to 00 0 to 1

- Three things happening here:  $\rightarrow$  Each  $x_i$  is compared to every other vector to establish weights for its own output  $y_i$   $\rightarrow$  Each  $x_i$  is compared to every other vector to establish weights for other outputs  $y_i$   $\rightarrow$  Each  $x_i$  is used as a part of weighted sum to get each output vector
- · Can we somehow come up with a clever way to get these 3 things done?

### Self Attention

· Self Alkention is just a sequence to sequence operation:

Newal Nets

Weighted Ang over all input y; = \( \subseteq W\_{ij} \times\_{j} \)

Wij = xtxj Wij = Softmon (Wij) · Weights are not learnable parameters -∞ to ∞ 0 to 1

· Three things happening here: q; = Wq xi

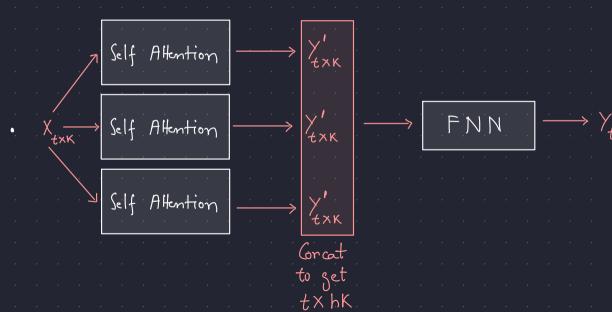
Ki = WK xi

 $\rightarrow$  Each  $x_i$  is compared to every other vector to establish weights for its own output  $y_i$   $\rightarrow$  Each  $x_i$  is compared to every other vector to establish weights for other outputs  $y_i$   $\rightarrow$  Each  $x_i$  is used as a part of weighted sum to get each output vector:  $V_i = W_{V} x_i$ 

$$W'_{ij} = q_i^T K_j$$
  $\Rightarrow$   $W'_{ij} = \frac{q_i^T K_j}{V K}$  Scaling dot product  $W_{ij} = \text{Softmox}(W_{ij})$   $\forall i = \sum_{i \in S} W_{ij} V_j$  Length

# Multi Head Attention

•  $\chi_{t \times K} \longrightarrow Self$  Altention  $\longrightarrow \chi_{t \times K}$ 



### **GNOT** Architecture

#### **GNOT: A General Neural Operator Transformer for Operator Learning**

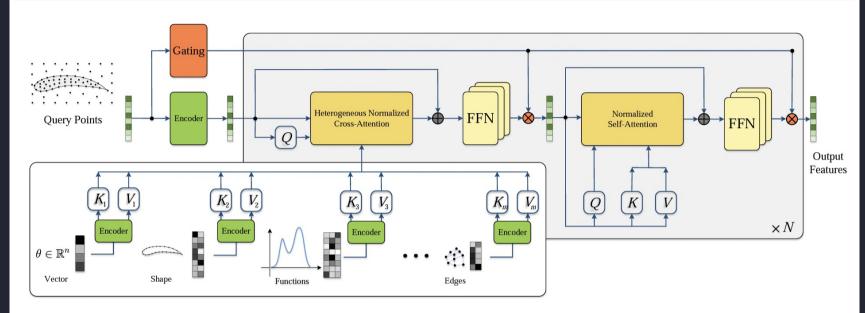


Figure 2. Overview of the model architecture. First, we encode input query points and input functions with different MLPs. Then we update features of query points using a heterogenous normalized cross-attention layer and a normalized self-attention layer. We use a gate network using geometric coordinates of query points to compute a weighted average of multiple expert FFNs. We output the features after processing them using N layers of the attention block.

Problem Formulation

· Query Points x;

· Input function 20

· Boundary Shape x;

· Final Solution U

Task is to map 20 -> 2 using domain Knowledge xi



 $(y^{\circ},x_{i}^{\prime})$ 

### **GNOT**

## General Input Encoding

- · Query Points x; 1 < i < Ng
- Input function  $u^{\circ}$   $(x_i, u_i^{\circ})$ ;  $i < N_i$
- · Boundary Shape x;
  - General Input Encoding resolves this issue
  - IXNe · x; NN q
    - Ngxne

•  $(x_i, u_i^\circ) \xrightarrow{NN} K$ 

All different dimensions

These NN are called Encoders

### **GNOT**

### Heterogeneous Normalised Attention

• From Encoder NNs we have 
$$\{q_i\}_{i \leq N}$$
  $\{K_i\}_{i \leq M}$   $\{V_i\}_{i \leq M}$ 

N ne Mine

$$\frac{Z_{t}}{i} = \sum_{\substack{e \neq h \ (q_{t}, K_{i}/T) \\ j \text{ dimensions of Embeddings}}} \frac{e \times h(q_{t}, K_{i}/T)}{i} \qquad \text{dimensions of Embeddings}$$

· Here is the proposed Modification:

$$\widehat{q}_{i}^{*} = \operatorname{Softmax}(q_{i}) = \underbrace{\operatorname{exh}(q_{i})}_{\text{ch}(q_{i})}$$
;  $j \leq n_{e}$ 

$$\underbrace{\operatorname{cxh}(q_{i})}_{\text{ch}(q_{i})}$$

$$Z_{t} = \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\xi \widetilde{q}_{t} \widetilde{K}_{i}}}_{N_{0} \times m_{e} \text{ lisation}} \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{i} \cdot V_{i}}_{V_{i}} = \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\xi \widetilde{q}_{t} \cdot \widetilde{K}_{i}}}_{N_{0} \times m_{e} \text{ lisation}} \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\zeta \widetilde{q}_{t} \cdot \widetilde{K}_{i}}}_{N_{0} \times m_{e} \text{ lisation}} \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\zeta \widetilde{q}_{t} \cdot \widetilde{K}_{i}}}_{N_{0} \times m_{e} \text{ lisation}} \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\zeta \widetilde{q}_{t} \cdot \widetilde{K}_{i}}}_{N_{0} \times m_{e} \text{ lisation}} \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\zeta \widetilde{q}_{t} \cdot \widetilde{K}_{i}}}_{N_{0} \times m_{e} \text{ lisation}} \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\zeta \widetilde{q}_{t} \cdot \widetilde{K}_{i}}}_{N_{0} \times m_{e} \text{ lisation}} \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\zeta \widetilde{q}_{t} \cdot \widetilde{K}_{i}}}_{N_{0} \times m_{e} \text{ lisation}} \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\zeta \widetilde{q}_{t} \cdot \widetilde{K}_{i}}}_{N_{0} \times m_{e} \text{ lisation}} \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\zeta \widetilde{q}_{t} \cdot \widetilde{K}_{i}}}_{N_{0} \times m_{e} \text{ lisation}} \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\zeta \widetilde{q}_{t} \cdot \widetilde{K}_{i}}}_{N_{0} \times m_{e} \text{ lisation}} \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\zeta \widetilde{q}_{t} \cdot \widetilde{K}_{i}}}_{N_{0} \times m_{e} \times m_{e} \text{ lisation}} \underbrace{\underbrace{\widetilde{q}_{t} \cdot \widetilde{K}_{i}}_{\zeta \widetilde{q}_{t} \cdot \widetilde{K}_{i}}}_{N_{0} \times m_{e} \times$$

• We have Oueries, Keys and values from different embeddings. 
$$Z_{t} = \widetilde{\gamma}_{t} + \frac{1}{L} \sum_{l=1}^{L} \langle \hat{\gamma}_{t}^{l}, (\sum_{i} \tilde{\kappa}_{i}^{l} \otimes v_{i}^{l}) \rangle = 0 ((N + \sum_{l} N_{l}) n_{e}^{2})$$

### Geometric Gating

· Inspired from domain decomposition techniques where each Net is made expert in a specific domain.

Soft domain Decomposition

· Query Points Kdiff NN Gibes

$$Z_t = Z_t + \sum_{i=1}^{K} P_i(x_t) \cdot E_i(Z_t)$$

- · We can learn Gating as NN on fix them with prior Knowledge
- · Also used in XPINNs

### **Conclusions**

**GNOT: A General Neural Operator Transformer for Operator Learning** 

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	1.05e-2
NS2d	_	part	-	1.56e-1	1.40e-1	1.56e-1	1.71e-1	1.38e-1
	-	full	_	8.20e-2	7.92e-2	8.20e-2	6.46e-2	4.43e-2
Elasticity	A	-	9.65e-2	5.08e-2	2.01e-2	2.20e-2	1.83e-2	8.65e-3
NS2d-c	A, C	u	2.74e-2	6.56e-2	1.52e-2	1.41e-2	2.33e-2	6.73e-3
		v	5.51e-2	1.15e-1	3.15e-2	2.98e-2	4.83e-2	1.55e-2
		p	2.74e-2	1.11e-2	1.59e-2	1.62e-2	2.43e-2	7.41e-3
NACA	A, C	-	1.32e-1	4.21e-2	1.61e-2	1.38e-2	1.83e-2	7.57e-3
Inductor2d	A, C	$A_z$	3.10e-2	_	2.56e-1	_	2.23e-2	1.21e-2
		$B_x$	3.49e-2	_	3.06e-2	_	2.83e-2	1.92e-2
		$B_y$	6.73e-2	_	4.45e-2	_	4.28e-2	3.62e-2
Heat	A, B, C	part	1.74e-1	_	_	_	_	4.13e-2
		full	1.45e-1	_	_	_	_	2.56e-2
Heatsink	A, B, C	T	4.67e-1	_	I—I	_	_	2.53e-1
		u	3.52e-1	_	-	_	_	1.42e-1
		v	3.23e-1	_	-	_	_	1.81e-1
		w	3.71e-1	_	_	_	_	1.88e-1

Table 1. Our main results of operator learning on several datasets from multiple areas. The types like u, v are the physical quantities to predict and types like "part" denotes the size of the dataset. "-" means that the method is not able to handle this dataset. Lower scores mean better performance and the best results are **bolded**.