

SciML 2: Neural operators

E. Franck

Master CSMI, Strasbourg University

01/09/2025

1. Operator learning: rappel

Rappel I

Definition (Nonlocal parametric operator): On se donne une fonction $\mathbf{v}_1 : \mathbb{R}^d \rightarrow \mathbb{R}^{\text{in}}$ et $\mathbf{v}_{1+1} : \mathbb{R}^d \rightarrow \mathbb{R}^{\text{out}}$. On appelle une couche non-locale l'opérateur

$$\mathbf{v}_{1+1}(\mathbf{x}) = \mathcal{B}_{\theta}(\mathbf{v})(\mathbf{x}) = \int_{\mathbf{D}} k_{\theta}(\mathbf{x}, \mathbf{y}, \mathbf{v}_1(\mathbf{x}), \mathbf{v}_1(\mathbf{y})) \mathbf{v}_1(\mathbf{y}) d\mathbf{y}$$

ou K_{θ} le noyau sera la partie apprenable. Dans la majorité des cas l'opérateur est linéaire pour généraliser les couches affines des réseaux classiques.

Definition (Local parametric operator): On se donne une fonction $\mathbf{v}_1 : \mathbb{R}^d \rightarrow \mathbb{R}^{\text{in}}$ et $\mathbf{v}_{1+1} : \mathbb{R}^d \rightarrow \mathbb{R}^{\text{out}}$. On appelle une couche non-locale l'opérateur

$$\mathbf{v}_{1+1}(\mathbf{x}) = \text{nn}_{\theta}(\mathbf{v}_1(\mathbf{x})), \forall \mathbf{x} \in \mathbb{R}^{\text{in}}$$

avec nn_{θ} un réseau de neurones ou une couche paramétrée par θ .

Rappel II

Definition (Opérateur paramétrique discret non local): On se donne $\mathbf{v} \in V \subset [L^2(\Omega)]^P$. On appelle $\mathcal{B}_\theta^h(\mathbf{v})$ un opérateur discret la transformation

$$\mathcal{B}(\mathbf{v})^h(\mathbf{x}) = \frac{1}{n} \sum_{i \in D_h} w_i k_\theta(\mathbf{x}, \mathbf{y}_i, \mathbf{v}(\mathbf{x}), \mathbf{v}(\mathbf{y}_i)) \mathbf{v}(\mathbf{y}_i).$$

Definition (discrétisation d'opérateur convergente): On se donne $\mathbf{v} \in V \subset [L^2(\Omega)]^P$. On se donne $\mathcal{B}_\theta(\mathbf{v})$ et un opérateur discret $\mathcal{B}_\theta^h(\mathbf{v})$. On dit qu'un d'opérateur discret non local est convergent si

$$\| \mathcal{B}_\theta(\mathbf{x}) - \mathcal{B}_\theta^h(\mathbf{x}) \|_{L_\infty(\Omega)} \xrightarrow[n \rightarrow \infty]{} 0, \quad \forall \mathbf{x} \in \Omega, \forall \mathbf{v} \in V$$

Definition (Opérateurs neuronaux): Un opérateur neuronaux sera une composition de discrétisation d'opérateurs paramétriques non locaux et locaux ou l'ensemble des donne une discrétisation convergente d'une composition d'opérateur.

2. Operator learning: Attention based methods

Attention I

- On va introduire la dernière architecture de réseau de neurones qui utilise des mécanismes d'attention pour améliorer la modélisation des opérateurs.
- Ici les couche nonlocale ne sont pas linéaires.

Definition (Mécanisme d'attention): Soit $Q \in \mathbb{R}^{q_i, q_o}$ une matrice de requête (query), $K \in \mathbb{R}^{k_i, k_o}$ une matrice de clé (key) et $V \in \mathbb{R}^{v_i, v_o}$ une de valeur (value). On définit le mécanisme d'attention par

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

avec d_k la dimension des clés et softmax la fonction softmax appliquée ligne par ligne. En pratique cela donne:

$$\text{attn}(Q, K, V)_i = \sum_{j=1}^n \alpha_{ij} V_j$$

avec

$$\alpha_{ij} = \frac{e^{\frac{\langle Q_i, K_j \rangle}{\sqrt{d_k}}}}{\sum_{l=1}^n e^{\frac{\langle Q_i, K_l \rangle}{\sqrt{d_k}}}}$$

Attention II

- Transformation du vecteur **valeur** en fonction de la similarité entre le vecteur de requête et le vecteur des clés.
- Plus la similarité est grande plus la composante vecteur valeur associée à la clé aura d'importance dans le résultat final.

Definition (Cross attention): Soit une matrice d'entrée X et un signal de sortie Y , la cross attention consiste à appliquer le mécanisme d'attention avec des matrices:

$$Q = YW_Q, \quad K = XW_K, \quad V = XW_V$$

avec W_Q, W_K, W_V des matrices apprennables.

- Cet opérateur va donc construire des transformations sur le signal de sortie en fonction du signal d'entrée puis construire une matrice de similarité qui va moduler les fonctions valeurs.

Definition (Self attention): Soit une matrice d'entrée X , la self attention consiste à appliquer le mécanisme d'attention avec des matrices:

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V$$

avec W_Q, W_K, W_V des matrices apprennables.

NO et Attention I

Definition (Continuous nonlocal self-attention based layer): Soit Ω un domaine (qui peut être une variété), $\mathbf{f} : \Omega \rightarrow \mathbb{R}^p$, $\mathbf{f} : \Omega \rightarrow \mathbb{R}^q$ deux fonctions sur le domaine. On définit une couche nonlocale basée sur l'attention par

$$\mathbf{g}(\mathbf{x}) = L_{\text{satt}}(\mathbf{f})(\mathbf{x}) = \int_{\Omega} k(\mathbf{x}, \mathbf{y}, ..) V \mathbf{f}(\mathbf{y}) d\mu(\mathbf{y})$$

avec

$$k_{\text{att}}(\mathbf{x}, \mathbf{y}, \mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{y})) = \frac{e^{\langle W_q \mathbf{f}(\mathbf{x}), W_k \mathbf{f}(\mathbf{y}) \rangle_{\mathbb{R}^p}}}{\int_{\Omega} e^{\langle W_q \mathbf{f}(\mathbf{x}), W_k \mathbf{f}(\mathbf{z}) \rangle_{\mathbb{R}^p}} d\mu(\mathbf{z})}$$

ou

$$k_{\text{att}}(\mathbf{x}, \mathbf{y}, \mathbf{c}(\mathbf{x}), \mathbf{f}(\mathbf{y})) = \frac{e^{\langle W_q \mathbf{c}(\mathbf{x}), W_k \mathbf{f}(\mathbf{y}) \rangle_{\mathbb{R}^p}}}{\int_{\Omega} e^{\langle W_q \mathbf{c}(\mathbf{x}), W_k \mathbf{f}(\mathbf{z}) \rangle_{\mathbb{R}^p}} d\mu(\mathbf{z})}$$

avec $W_q \in \mathbb{R}^{p \times d_k}$, $W_k \in \mathbb{R}^{p \times d_k}$ et $V \in \mathbb{R}^{p \times q}$ des matrices apprenables.

- On voit donc que l'attention permet de construire naturellement des opérateurs non locaux à noyaux.
- Premier exemple d'opérateur non local non linéaire.

Multi head attention

Definition (Multi-head attention): Soit h le nombre de têtes d'attention. Soit $\mathbf{f} : \Omega \rightarrow \mathbb{R}^p$, $\mathbf{g} : \Omega \rightarrow \mathbb{R}^q$ deux fonctions sur le domaine. On définit une couche d'attention multi-tête par

$$\mathbf{g}(\mathbf{x}) = \sum_{l=1}^L w_l L_{\text{att},l}(\mathbf{f})$$

avec $L_{\text{att},l}(\mathbf{f})$ obtenue par un module d'attention de matrice W_q^l , W_k^l et V^l et $w_l \in \mathbb{R}$ des poids apprenables.

Definition (Nonlocal self-attention based layer): Soit $G = [V, E]$ un graphe avec $|V| = n$. Soit $\mathbf{f} : V \rightarrow \mathbb{R}^p$, $\mathbf{g} : V \rightarrow \mathbb{R}^q$ deux fonctions sur le graphe. On définit une couche nonlocale basée sur l'attention par

$$\mathbf{g}(i) = \sum_{j=1}^n w_i k_{\text{att}}(i, j) V \mathbf{f}(j)$$

avec

$$k_{\text{att}}(i, j, \mathbf{f}(i), \mathbf{f}(j)) = \frac{e^{\langle W_q \mathbf{f}(i), W_k \mathbf{f}(j) \rangle_{\mathbb{R}^p}}}{\sum_{l=1}^n e^{\langle W_q \mathbf{f}(i), W_k \mathbf{f}(l) \rangle_{\mathbb{R}^p}}}, \quad \text{ou} \quad k_{\text{att}}(i, j, \mathbf{c}(i), \mathbf{f}(j)) = \frac{e^{\langle W_q \mathbf{c}(i), W_k \mathbf{f}(j) \rangle_{\mathbb{R}^p}}}{\sum_{l=1}^n e^{\langle W_q \mathbf{c}(i), W_k \mathbf{f}(l) \rangle_{\mathbb{R}^p}}}$$

avec $W_q \in \mathbb{R}^{p \times d_k}$, $W_k \in \mathbb{R}^{p \times d_k}$ et $V \in \mathbb{R}^{p \times q}$ des matrices apprenables.

Neural operator et couche locale

Proposition (Attention nonlocal neural operator layer): Soit Ω un domaine (qui peut être une variété), $f : \Omega \rightarrow \mathbb{R}^p$, $g : \Omega \rightarrow \mathbb{R}^q$ deux fonctions sur le domaine. Soit des ensemble de points de discretisation $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \Omega$ et $Y = \{\mathbf{y}_1, \dots, \mathbf{y}_m\} \subset \Omega$ obtenue à partir d'une lois de probabilité $p(y)$. Les couches de “self” et “cross” définissent des approximations convergentes des couches continues lorsque le nombre de points de discretisation tend vers l'infini a condition que

$$w_i = \frac{1}{p(\mathbf{y}_i)}$$

- Pour obtenir un opérateur neuronal complet on peut combiner les couches d'attention avec des couches locales non linéaires.
- Entre les couches d'attention les transformers utilisent des couches feedforward classiques combinées à des couches de normalisation.
- On peut imaginer des architectures similaires pour les opérateurs neuronaux basés sur l'attention.

Couche locale

Definition (Local normalization layer): Soit Ω un domaine (qui peut être une variété), $\mathbf{f} : \Omega \rightarrow \mathbb{R}^p$ une fonction sur le domaine. On définit une couche de normalisation L_{ln} locale est donnée par

$$\mathbf{g}_k(\mathbf{x}) = \gamma_k \frac{\mathbf{f}_k(\mathbf{x}) - \mu(\mathbf{f}(\mathbf{x}))}{\sqrt{\sigma^2(\mathbf{f}(\mathbf{x})) + \varepsilon}} + \beta_k, \forall k \in \{1, \dots, p\}$$

avec

$$\mu(\mathbf{f}(\mathbf{x})) = \frac{1}{p} \sum_{i=1}^p \mathbf{f}_i(\mathbf{x}), \quad \sigma^2(\mathbf{f}(\mathbf{x})) = \frac{1}{p} \sum_{i=1}^p (\mathbf{f}_i(\mathbf{x}) - \mu(\mathbf{f}(\mathbf{x})))^2$$

et $\gamma = (\gamma_1, \dots, \gamma_p)$, $\beta = (\beta_1, \dots, \beta_p)$ des paramètres apprennable et ε un petit scalaire pour la stabilité numérique.

Neural operator block

Definition (Transformer neural operator block): Soit Ω un domaine (qui peut être une variété), $\mathbf{f} : \Omega \rightarrow \mathbb{R}^p$, $\mathbf{g} : \Omega \rightarrow \mathbb{R}^q$ deux fonctions sur le domaine. On définit bloc de transformer neural operator par la succession des opérations suivantes:

$$\mathbf{g}_1(\mathbf{x}) = L_{\text{att}}(\mathbf{f}) + W_1 \mathbf{f}(\mathbf{x})$$

$$\mathbf{g}_2(\mathbf{x}) = L_{\text{ln}}(\mathbf{g}_1)(\mathbf{x})$$

$$\mathbf{g}_3(\mathbf{x}) = W_2 \mathbf{g}_2(\mathbf{x}) + L_{\text{nn}}(\mathbf{g}_2(\mathbf{x}))$$

et

$$\mathbf{g}(\mathbf{x}) = L_{\text{ln}}(\mathbf{g}_3)(\mathbf{x})$$

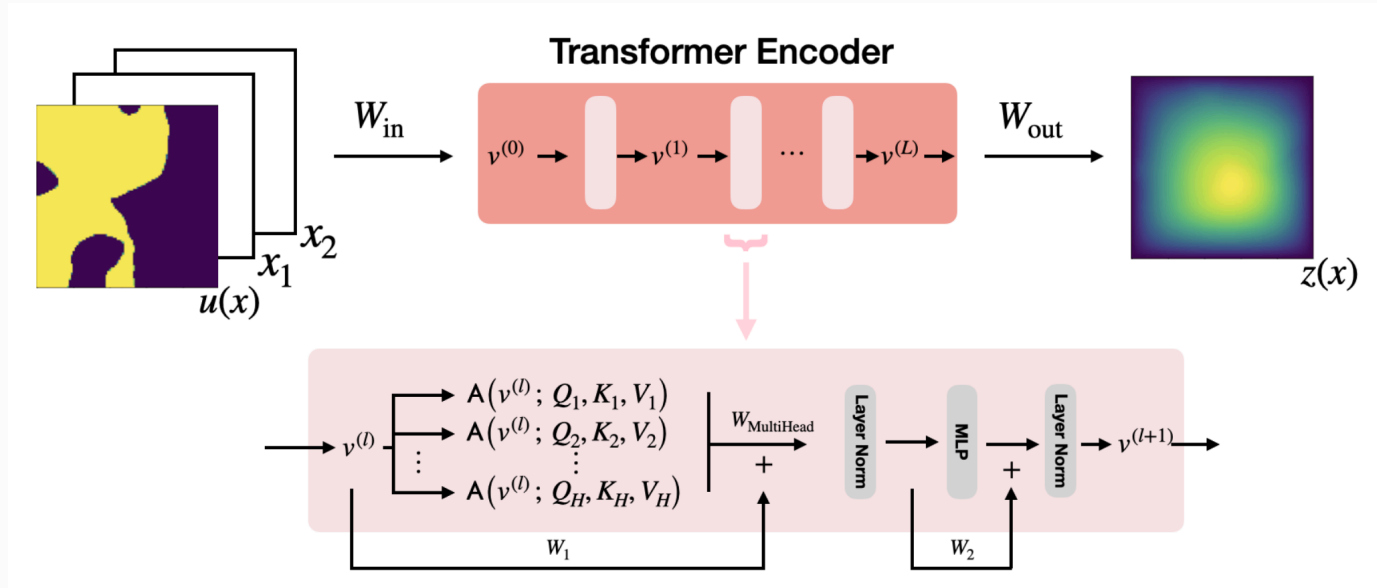
avec W_1, W_2 des matrices apprenables, L_{att} une couche d'attention (self ou cross), L_{ln} une couche de normalisation locale et L_{nn} une couche feedforward non linéaire.

Neural operator block

Definition (Transformer neural operator): Soit Ω un domaine (qui peut être une variété), $\mathbf{f} : \Omega \rightarrow \mathbb{R}^p$, $\mathbf{g} : \Omega \rightarrow \mathbb{R}^q$ deux fonctions sur le domaine. On définit un “transformer neural operator” par la transformation:

$$\mathbf{g}(\mathbf{x}) = \mathbf{Q} \circ \mathbf{L}_{1,\text{Bt}} \circ \dots \circ \mathbf{L}(L, \text{Bt}) \circ \mathbf{P}(\mathbf{f})(\mathbf{x})$$

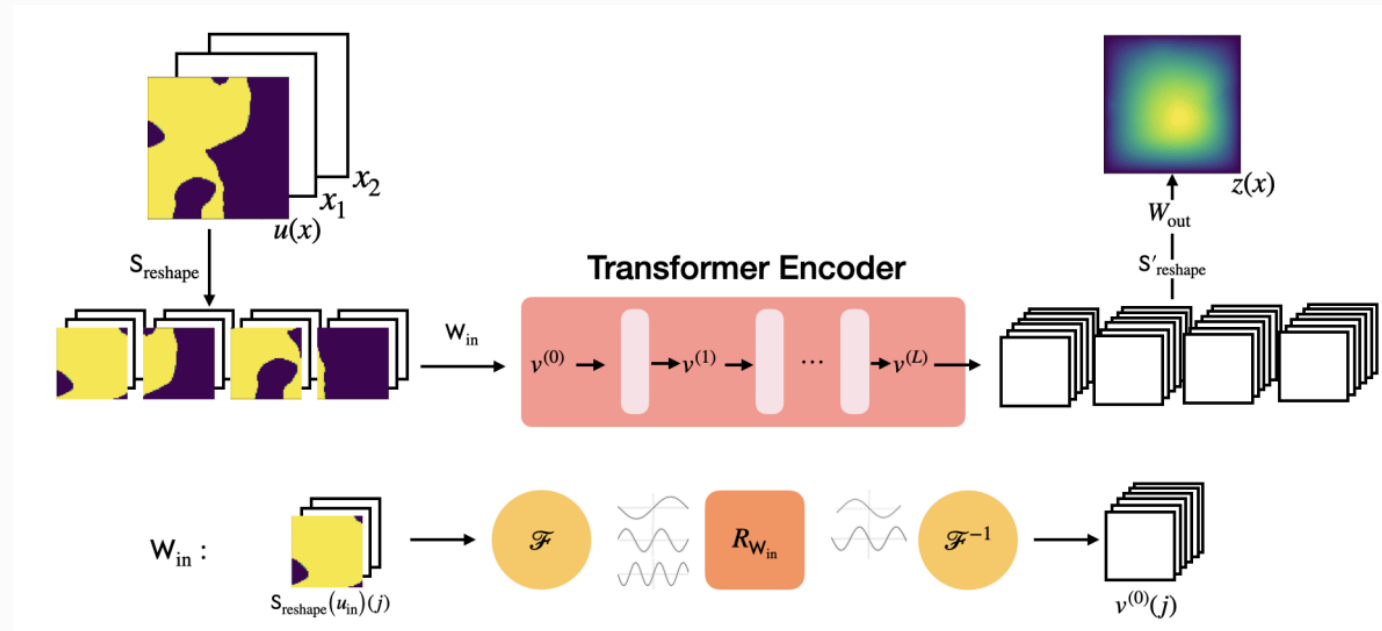
avec $\mathbf{L}(i, \text{Bt})$ des blocs de type transformer, \mathbf{P} et \mathbf{Q} des couches linéaires locales d’adaptation de dimension.



Patch Transformer neural operator

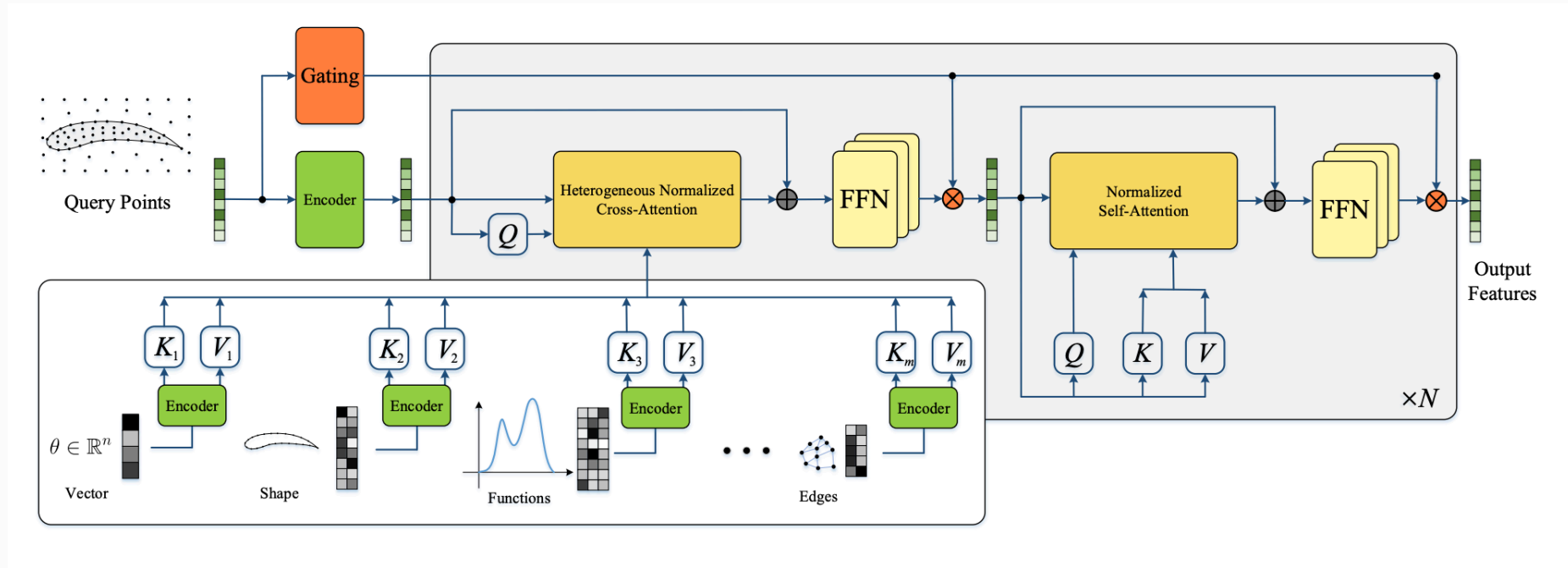
Idea: Pour réduire le coût computationnel on peut imaginer découper le domaine en patches/sous domaines et appliquer un transformer sur les représentations locales de chaque patch.

- Ils peuvent être recombinaés entre eux par les couches de multi-head attention.



GNOT I

- On va maintenant introduire une autre architecture basée sur l'attention: le GNOT (General Neural Operator Transformer).



- On cherche un opérateur:

$$\mathcal{G} : \mathcal{A} \rightarrow \mathcal{H}$$

avec $\mathcal{A} = \mathcal{H} \times \dots \times \mathcal{H} \times \mathbb{R}^p$ qui va contenir les **points du domaine**, les fonctions d'entrées et de sortie du modèle (source, condition aux limites, etc) et les paramètres physiques.

GNOT II

- L'objectif est de pouvoir traiter des entrées variées: les points du domaine, les conditions aux limites, les sources, les paramètres physiques.
 - Soit $\beta \in \mathbb{R}^p$. On le transforme avec un encodeur MLP $f_{\theta_1}(\beta) \in \mathbb{R}^{1,n_e}$
 - Soit $\mathbf{x}_b \in \mathbb{R}^{n_x,d}$ les points de bord du domaine. On les transforme avec un encodeur MLP $f_{\theta_2}(\mathbf{x}) \in \mathbb{R}^{n_x,d}$
 - Soit $\mathbf{x}, \mathbf{a} \in \mathbb{R}^{n_x,d+p}$ les points du domaine et les valeurs des fonctions d'entrée à ses points. On les transforme avec un encodeur MLP $f_{\theta_3}(\mathbf{x}, \mathbf{a}) \in \mathbb{R}^{n_d,d}$

Remark: On peut encoder d'autres quantités comme des quantités aux arrêtes etc.

- On va encoder (\mathbf{x}, \mathbf{a}) avec un encodeur MLP $f_{\theta_e}(\mathbf{x}, \mathbf{a}) \in \mathbb{R}^{n,n_e}$. On le nomme X
- L'encodage des quantités supplémentaires, comme les paramètres, les points de bords ou autre est nommé $Y \in \mathbb{R}^{n,n_e}$

Idea: On va utiliser une architecture basée sur l'attention cross entre X et Y pour construire un opérateur non local entre les deux.

GNOT III

Idea: Après l'étape d'encodage on va appliquer des opérateurs d'attention.

- Cross-attention sur X et Y:

$$\mathbf{z}_i = \sum_{j=1}^{n_d} \frac{e^{\mathbf{q}_i \cdot \mathbf{k}_j}}{\sum_{l=1}^n e^{\mathbf{q}_i \cdot \mathbf{k}_l}} \mathbf{v}_j$$

avec $\mathbf{k}_j = W_k \mathbf{y}_j$ et $\mathbf{v}_j = W_v \mathbf{y}_j$

Il est proposé une simplification:

- ▶ On normalise:

$$\tilde{\mathbf{q}}_i = \frac{e^{\mathbf{q}_i j}}{\sum_{l=1}^{n_e} e^{\mathbf{q}_i l}}, \quad \tilde{\mathbf{k}}_j = \frac{e^{\mathbf{k}_j i}}{\sum_{l=1}^{n_e} e^{\mathbf{k}_l i}}$$

- ▶ La nouvelle attention est donnée attention est donnée par:

$$\mathbf{z}_i = \sum_{j=1}^{n_d} \frac{\tilde{\mathbf{q}}_i \cdot \tilde{\mathbf{k}}_j}{\sum_{l=1}^{n_e} \tilde{\mathbf{q}}_i \cdot \tilde{\mathbf{k}}_l} \mathbf{v}_j = \alpha_i \sum_{j=1}^{n_d} \tilde{\mathbf{q}}_i \cdot \tilde{\mathbf{k}}_j \mathbf{v}_j, \quad \text{avec} \quad \alpha_i = \frac{1}{\sum_{l=1}^{n_e} \tilde{\mathbf{q}}_i \cdot \tilde{\mathbf{k}}_l} = \alpha_i \tilde{\mathbf{q}}_i \cdot \sum_{j=1}^{n_d} \tilde{\mathbf{k}}_j \otimes \mathbf{v}_j$$

GNOT III

Definition (cross-attention layer of GNOT): On définit notre signal d'entrée après encodage $X \in \mathbb{R}^{N, n_e}$ et les signaux supplémentaires $Y_{1 \leq l \leq L} \in \mathbb{R}^{N, n_e}$. La couche d'attention cross de GNOT est donnée par

$$\mathbf{z}_i = \tilde{\mathbf{q}}_i + \frac{1}{L} \sum_{l=1}^L \alpha_{i,l} \tilde{\mathbf{q}}_{i,l} \cdot \left(\sum_{j=1}^N \tilde{\mathbf{k}}_j \otimes \mathbf{v}_j \right)$$

- Pour traiter les aspects multi-échelles les auteurs proposent de construire des pondérations dépendantes de l'espace de notre état latent modifié par plusieurs MLP.
- On a donc

$$\mathbf{z}_i = \mathbf{z}_i + \sum_{k=1}^K p(\mathbf{x}_i) \text{MLP}_i(\mathbf{z}_i)$$

avec

$$p(\mathbf{x}_i) = \frac{\exp G_k(\mathbf{x}_i)}{\sum_{l=1}^K \exp G_l(\mathbf{x}_i)}, \quad G : \mathbb{R}^d \rightarrow \mathbb{R}^K$$

3. Operator learning: Time problem extension

Apprentissage d'opérateurs temporel

- Soit un problème elliptique de la forme

$$\begin{cases} \partial_t u - \Delta u = f(t, \mathbf{x}), & \mathbf{x} \in \Omega \\ u(t, \mathbf{x}) = 0 & \mathbf{x} \in \partial\Omega \\ u(0, \mathbf{x}) = u_{\{0\}}(\mathbf{x}) & \mathbf{x} \in \Omega \end{cases}$$

- On voudrait prédire rapidement une solution pour un $f(t, \mathbf{x})$ ou un $u_0(\mathbf{x})$ généraux.
- Si la solution existe, on peut formellement dire qu'il existe $G : \mathcal{H}_x \rightarrow \mathcal{H}_y$ tel que

$$u(.) = G^{-1}(f(.), u_0(.))$$

- Comme précédemment peut généraliser en cherchant à approcher $G^+(\mathbf{a}(.)) = u(.)$ avec \mathbf{a} l'ensemble des fonctions qui sont des entrées du problème (sources, conditions limites, conditions initiales, coefficients de l'EDP).
- Si on regarde précisément le problème:

$$u(t, .) = G^{-1}(u_0(.))$$

- On voit qu'il s'agit en fait d'apprendre le flot d'EDP.

Remark (Deux cas): flot discret ou flot continu en temps

Opérateurs neuronaux continus en temps

- Peut t'on généraliser la théorie de Green ? **Oui**

$$u(t, \mathbf{x}) = \int_{\Omega} G(t, 0, \mathbf{x}, \mathbf{y}) u_0(\mathbf{y}) d\mathbf{y} + \int_0^t \int_{\Omega} G(t, \tau, \mathbf{x}, \mathbf{y}) f(\tau, \mathbf{y}) d\mathbf{y} d\tau \\ + \int_0^t \int_{\partial\Omega} [u(t, \mathbf{x}) \nabla G(t, 0, \mathbf{x}, \mathbf{y}) - G(t, 0, \mathbf{x}, \mathbf{y}) \nabla u(t, \mathbf{x})] d\mathbf{y} d\tau$$

avec

$$(\partial_t - \Delta)G(t, \tau, \mathbf{x}, \mathbf{y}) = \delta(\mathbf{x} - \mathbf{y})\delta(t - \tau)$$

et des conditions initiales/bords nulles.

Definition (GreenNet temporal operator): Il s'agit d'un réseau à une **un couche** sans couche d'extrapolation et de projection ou la couche principale est $G_{\theta} : H^{s(\Omega)^p} \rightarrow H^{s(\Omega)^p}$ est paramétrée de la façon suivante:

$$G_{\theta}(\mathbf{f}, \mathbf{u}_0) = \int_{\Omega} k_{\theta}(t, 0, \mathbf{x}, \mathbf{y}) \mathbf{u}_0(\mathbf{y}) d\mathbf{y} + \int_0^t \int_{\Omega} k_{\theta}(t, \tau, \mathbf{x}, \mathbf{y}) f(\tau, \mathbf{y}) d\mathbf{y} d\tau$$

avec k_{θ} un réseau de type MLP, Fourier etc.

Les discretisations convergentes sont similaires au cas stationnaire.

Low rank approximation et deepOnet

Definition (GreenNet temporel de Bas rang): Soit f une fonction uniquement spatiale (condition initiale) ou dépendent du temps (source). On propose de prendre:

$$G_{\theta}(f) = \sum_{k=1}^R \langle f, \psi_i \rangle_{L^2} \varphi_{\theta,k}(\mathbf{x}, t)$$

avec

$$\langle f, \psi_i \rangle_{L^2} = \frac{1}{m} \sum_{j=1}^m \psi_{\theta,i}(\mathbf{x}_j, t_j) f(\mathbf{x}_j, t_j), \quad \text{ou} \quad \langle f, \psi_i \rangle_{L^2} = \frac{1}{m} \sum_{j=1}^m \psi_{\theta,i}(\mathbf{x}_j, t_j) f(\mathbf{x}_j),$$

Definition (DeepOnet temporel): Soit \mathbf{a} une fonction uniquement spatiale (condition initiale) ou dépendent du temps (source). On propose de prendre:

$$G_{\theta}(\mathbf{a}(t_1, \mathbf{x}_1), \dots, \mathbf{a}(t_m, \mathbf{x}_m)) = \sum_{i=1}^R \alpha_i b_{\theta,i}(\mathbf{x}, t)$$

avec b_{θ} un réseau de type MLP de \mathbb{R}^d dans \mathbb{R}^R et

$$t_{\theta}(\mathbf{a}(t_1, \mathbf{x}_1), \dots, \mathbf{a}(t_m, \mathbf{x}_m)) \rightarrow (\alpha_1, \dots, \alpha_K)$$

un réseau de $\mathbb{R}^{q \times m}$ dans \mathbb{R}^R

FNO continue en temps

- Une couche de Fourier est donnée par:

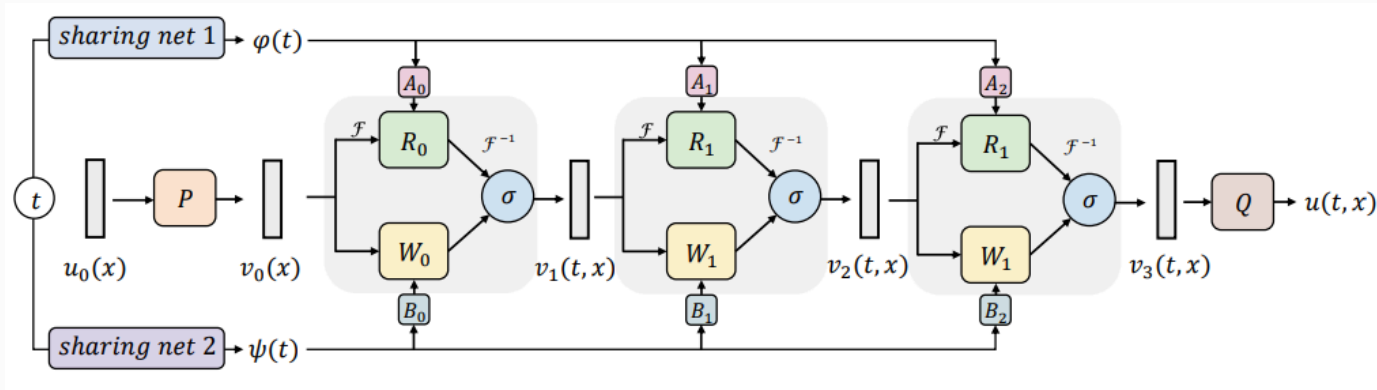
$$\mathbf{v}_{l+1} = \sigma\left(W_l \mathbf{v}_l + \left(\mathcal{F}^{\{-1\}} \circ \hat{\kappa}_{\theta_l} \circ \mathcal{F}\right)(\mathbf{v}_l)\right)$$

- Pour l'étendre au temporel il est proposé une **modulation temporelle des poids** de chaque couche:

$$\mathbf{v}_{l+1} = \sigma\left(\psi_{\theta_t}(t) W \mathbf{v}_l + \left(\mathcal{F}^{-1} \circ \varphi_{\theta_t}(t, \mathbf{k}) \hat{\kappa}_{\theta_l} \circ \mathcal{F}\right)(\mathbf{v}_l)\right)$$

avec \mathbf{k} les fréquences.

- Il est proposé de prendre $\varphi_{\theta_t}(t, \mathbf{k}) = \varphi(t) A_l(\mathbf{k})$ avec $A(\mathbf{k})$ une matrice entraînable.
- Il est proposé de prendre $\varphi_{\theta_t}(t, \mathbf{k}) = \text{diag}(\psi(t) B_l)$ avec B_l une matrice entraînable.
- les fonctions $\psi(t)$ et $\varphi(t)$ sont communes à toutes les couches.



Transformeur continue en temps

Idea (Transformer): pour prendre un transformer temporel la méthode la plus utilisé consiste a moduler les couches locales par des fonctions temporelles

- On peut imaginer faire dependre les MLP locales du temps.
- On peut moduler les couches linéaires:

$$\mathbf{g}(\mathbf{x}, t) = \psi_{\theta_t}(t)W$$

- On peut moduler les couches de normalisation locale:

$$\left(F_{\text{LayerNorm}}(\mathbf{v}; \gamma, \beta,)(\mathbf{x})\right)_k = \gamma_{k(t)} \cdot \frac{(\mathbf{v}(\mathbf{x}))_k - m(\mathbf{v}(\mathbf{x}))}{\sqrt{\sigma^2(\mathbf{v}(\mathbf{x})) + \varepsilon}} \beta_{k(t)},$$

avec

- $\beta(t) = \beta_1 + \beta_2 t$ et $\gamma(t) = \gamma_1 + \gamma_2 t$
- $\beta(t), \gamma(t)$ des réseaux de neurones type MLP.

Remark: On pourrait moduler les matrices d'attention par une fonction $\psi_{\theta_t}(t)$.

Propriété de semi-groupe

Remark (Propriété de semi-groupe): Si on cherche à prédire la solution à partir de la condition initiale sans source. Le flot qui sera approché par l'opérateur neuronale satisfait la **structure de semi-groupe**.

- On se sait pas imposer cette structure dans les architectures actuelles. On peut imposer cette structure faiblement.
 - La fonction de coût classique est donnée par:

$$\mathcal{L}(\theta) := \frac{1}{M(K+1)} \sum_{i=1}^M \sum_{k=0}^K \| G(t_k, a_i(\mathbf{x})) - u_i(t_k, \mathbf{x}) \|_{L^p(D)}^p,$$

avec $u_i(t_k, \mathbf{x})$ la solution associée à la condition initiale $a_i(\mathbf{x})$;

- On impose faiblement la structure de semi-groupe:

$$\hat{\mathcal{L}}(\theta) := \frac{1}{M\hat{K}} \sum_{i=1}^M \sum_{k, \hat{k}=0, k \leq \hat{k}}^K \| u_i(t_{\hat{k}}, \mathbf{x}) - G_{\theta}(t_{\hat{k}} - t_k, u_i(t_k, \mathbf{x})) \|_{L^p(D)}^p,$$

Remark: Une autre approche est de considérer une architecture de type EDO

$$\frac{d}{dt} \mathbf{u}(t) = \mathcal{N}_{\theta}(\mathbf{u}(t))$$

avec $\mathcal{N}_{\theta}(\mathbf{u}(t))$ un opérateur neuronale.

Opérateur neuronaux discret en temps I

Idea: On demande à l'opérateur neuronale de prédire la solution à des temps discret.

- On se donne un opérateur neuronale $G_{(\theta(u))}$
- Plusieurs options:
 - On prédit un temps avec un temps: $u_{t_{k+1}}(\mathbf{x}) = G_{\theta}(u_{t_k}(\mathbf{x}))$
 - On prédit plusieurs temps avec un temps: $(u_{t_{k+m}}(\mathbf{x}), \dots, u_{t_{k+1}}(\mathbf{x})) = G_{\theta}(u_{t_k}(\mathbf{x}), \dots, u_{t_{k-m}}(\mathbf{x}))$

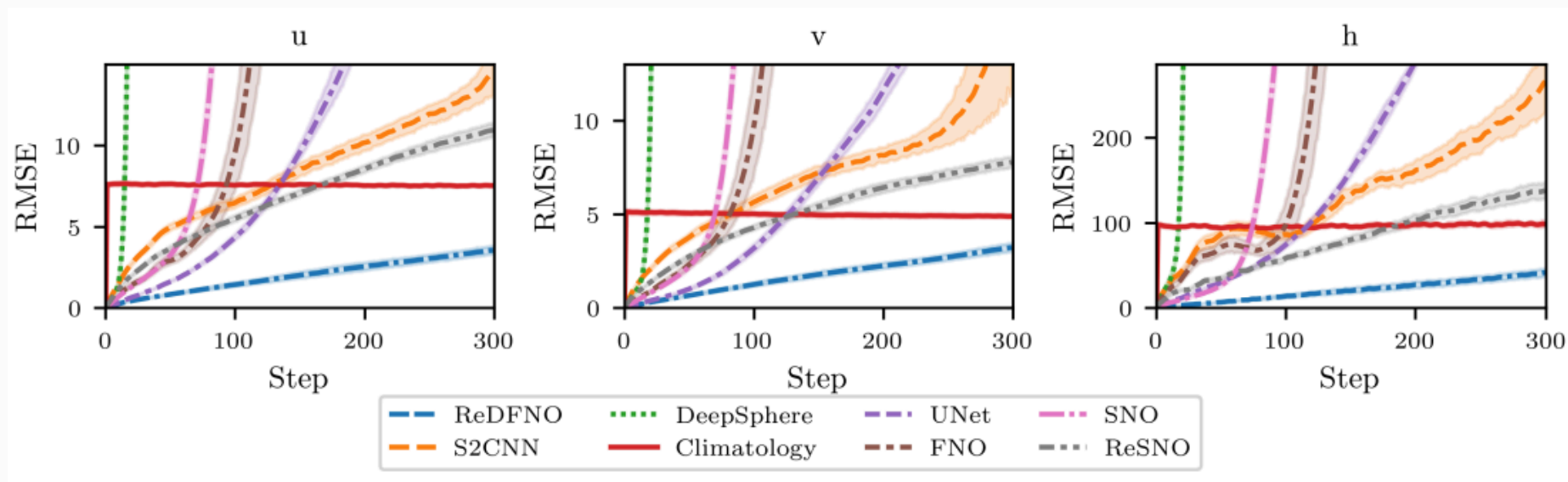
Definition (NO auto-régressif): On utilise donc les opérateurs neuronaux de façon auto-régressif:

$$u_{t_n}(\mathbf{x}) = (G_{\theta} \circ \dots \circ G_{\theta})(u_0(\mathbf{x}))$$

Remark (Inconvénient): Est ce ce type de prédiction est **stable** en temps ?

Opérateur neuronaux discret en temps II

- Problème de stabilité en temps long.

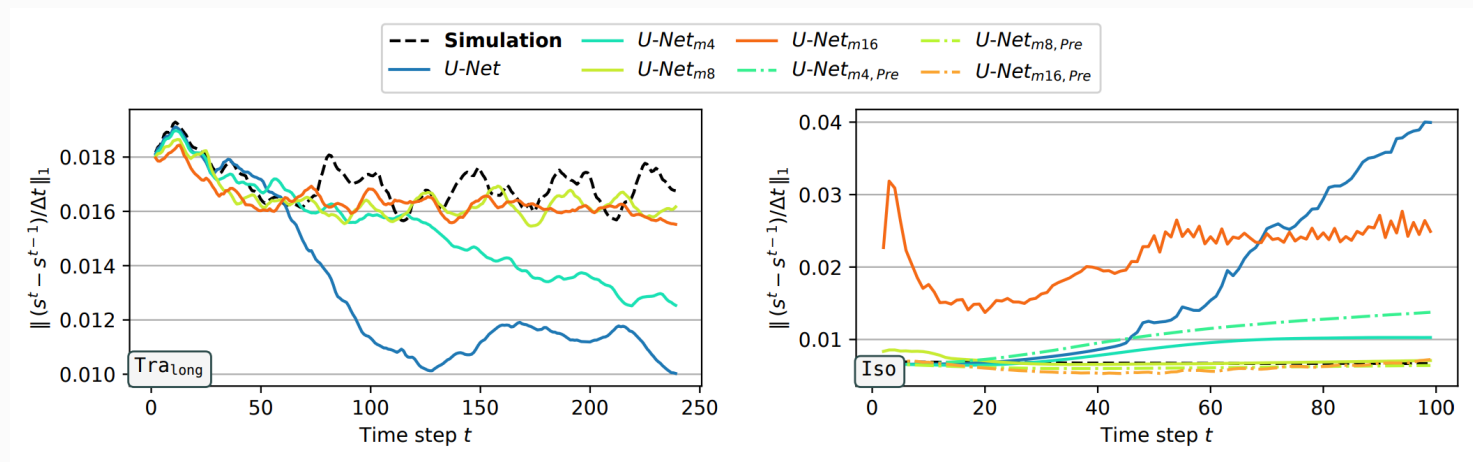


Idea: Un entraînement sur plusieurs pas de temps (unrolled) peut stabiliser les opérateurs neuronaux.

- Références:
 - ▶ “Benchmarking Autoregressive Conditional Diffusion Models for Turbulent Flow Simulation”. N. Thuerey and al.
 - ▶ “Towards Stability of Autoregressive Neural Operators”. J. Brown and al.

Opérateur neuronal discret en temps II

- Comparaison de different unrolling pour un Unet



- Coût mémoire et CPU de “l’unrolling”

