

Understanding the Expressivity and Trainability of Fourier Neural Operator: A Mean-Field Perspective

Takeshi Koshizuka (The University of Tokyo)

Masahiro Fujisawa (RIKEN AIP), Yusuke Tanaka (NTT)

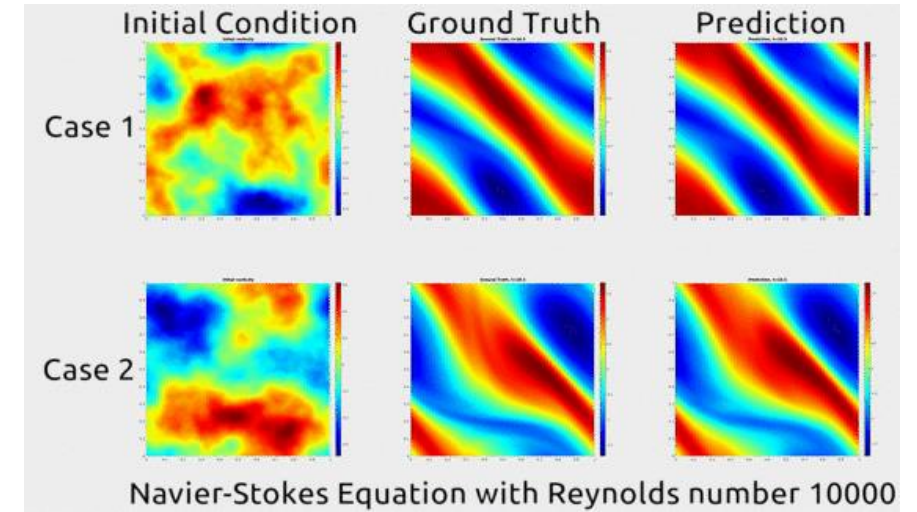
Issei Sato (The University of Tokyo)



Background

Machine Learning for PDE

- Accelerate the solution process
- Enhance the precision of solutions using observation
- Enable solutions across a wide range of conditions and parameters



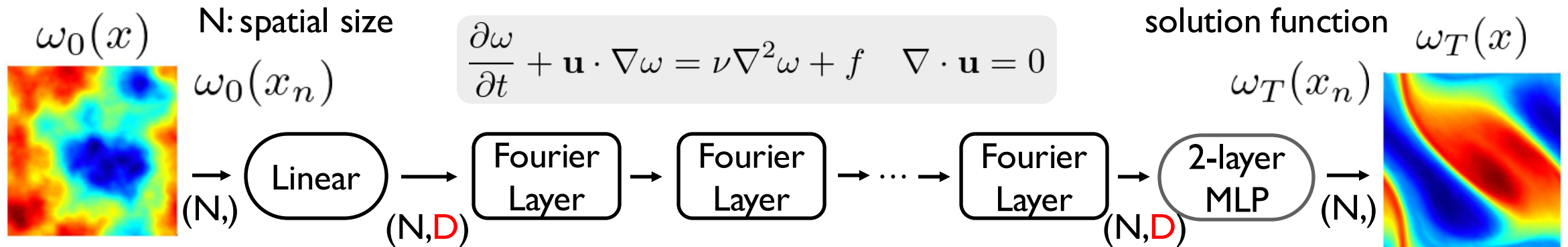
Neural Operators [Kovachiki & Li et al.'21] **discretization-invariant**

Integral Operator σ_ℓ : non-linear activation

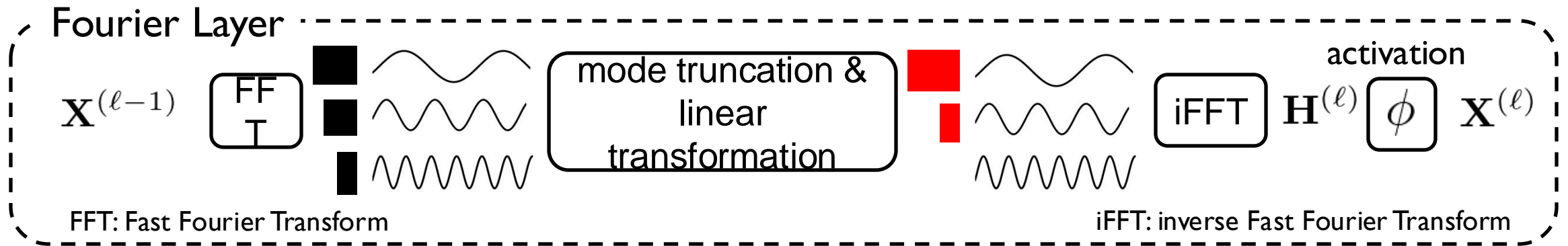
$$v_{t+1}(x) = \sigma \left(W v_t(x) + \int_D G_\theta(x, y, a(x), a(y)) v_t(y) d\nu_x(y) \right)$$

Model **the integral kernel (Green's function) G_θ** with a neural network and stack multiple layers

Fourier Neural Operator (FNO) [Zongyi et al '20]



D : number of hidden feature (width)



The deep FNO has poor performance.

e.g. instability of initial training
[Lu et al. '20, Tran et al. '23]

Analyze **expressivity** & **trainability** from a mean-field perspective

Infinite-width FNO at initialization

Fourier Layer $\mathbf{X}^{(\ell-1)} \mapsto \mathbf{X}^{(\ell)}$ concatenation learnable parameters

$$\mathbf{X}^{(\ell)} = \phi(\mathbf{H}^{(\ell)}), \mathbf{H}^{(\ell)} = \text{iFFT} \left(\left\|_{k=0}^{N-1} \text{FFT}(\mathbf{X}^{(\ell-1)})_{k,:} \left(\Theta^{(\ell,k)} + \sqrt{-1} \Xi^{(\ell,k)} \right) \right\| + \mathbf{1}_N \mathbf{b}^{(\ell)} \right)$$

linear transformation tunable hyperparameters

Initialization

$$\Theta_{i,j}^{(\ell,k)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N} \left(0, \frac{\sigma^2}{2D} \right), \Xi_{i,j}^{(\ell,k)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N} \left(0, \frac{\sigma^2}{2D} \right) \quad (0 \leq k \leq K-1), b_i^{(\ell)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N} \left(0, \sigma_b^2 \right)$$

D: number of hidden feature (width)

Infinite width

central limit theorem

$$\mathbf{H}_{:,1}^{(\ell)}, \mathbf{H}_{:,2}^{(\ell)}, \dots, \mathbf{H}_{:,D}^{(\ell)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N} \left(\mathbf{0}, \Sigma^{(\ell)} \right)$$

signals

$$\frac{\partial \mathcal{L}}{\partial \mathbf{H}_{:,1}^{(\ell)}}, \frac{\partial \mathcal{L}}{\partial \mathbf{H}_{:,2}^{(\ell)}}, \dots, \frac{\partial \mathcal{L}}{\partial \mathbf{H}_{:,D}^{(\ell)}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N} \left(\mathbf{0}, \tilde{\Sigma}^{(\ell)} \right)$$

Mode truncation

$$\Theta^{(\ell,k)} = \Xi_{i,j}^{(\ell,k)} = 0 \quad \text{for } K \leq k \leq N-1-K$$

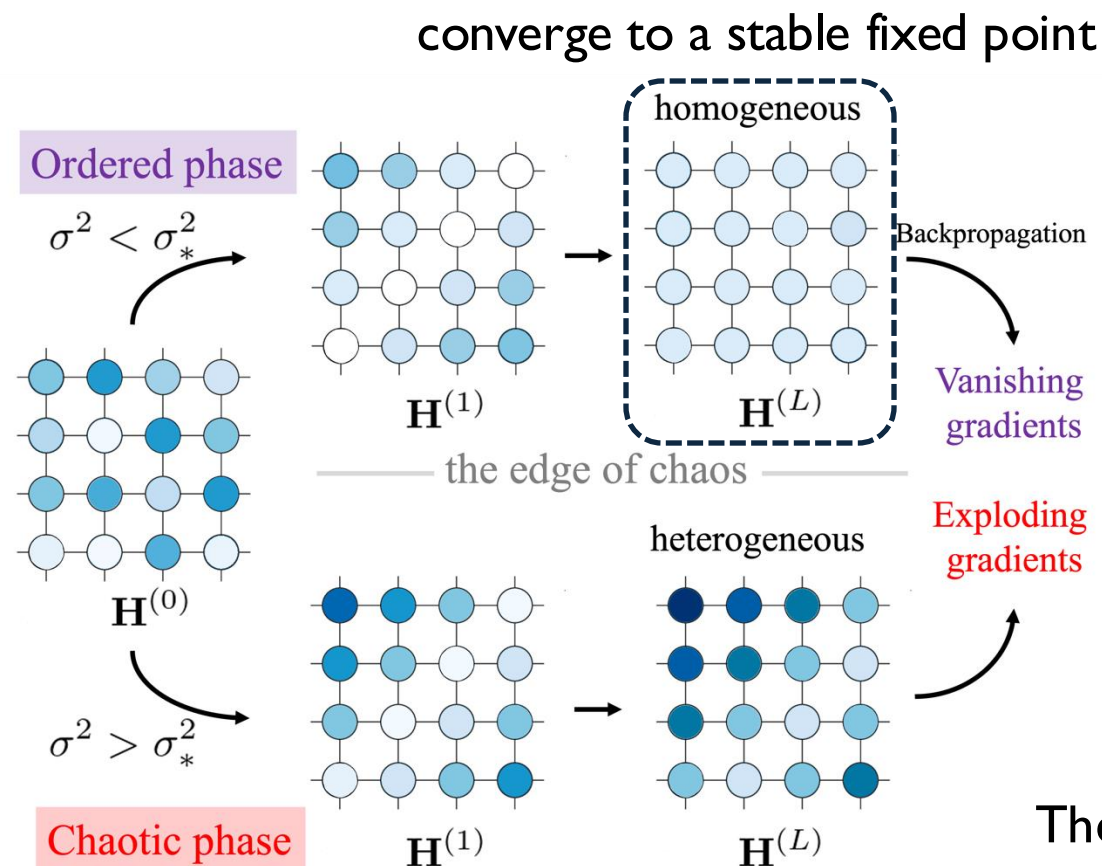
Analyze the dynamics

$$\Sigma^{(1)} \xrightarrow{c} \dots \xrightarrow{c} \Sigma^{(\ell-1)} \xrightarrow{c} \Sigma^{(\ell)}$$

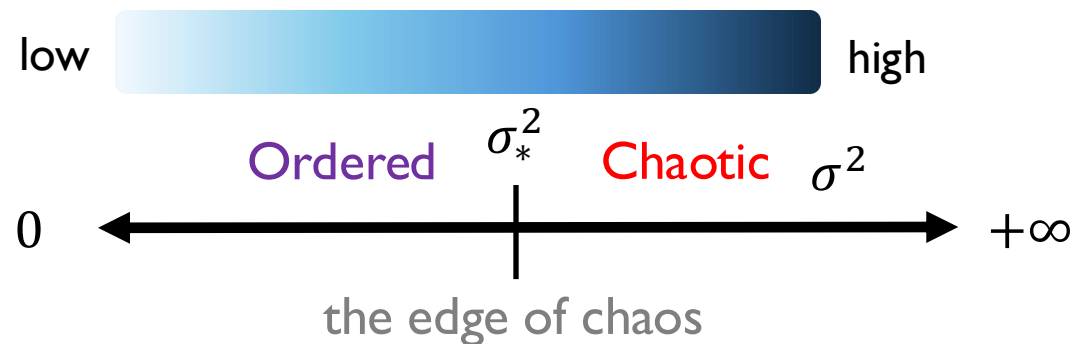
$$\tilde{\Sigma}^{(\ell)} \xleftarrow{\tilde{c}} \tilde{\Sigma}^{(\ell+1)} \xleftarrow{\tilde{c}} \dots \xleftarrow{\tilde{c}} \tilde{\Sigma}^{(L)}$$

Main Results

Ordered-Chaos phase transition for the weight initialization parameter σ^2



Expressivity $\text{diag } \Sigma^{(\ell)} \sim e^{\ell \log C \sigma^2}$, $C : \text{const}$



Trainability $\text{diag } \tilde{\Sigma}^{(\ell)} \sim e^{(L-\ell) \log C \sigma^2}$

The constant C is determined by the activation function.
 The transition point σ_*^2 is obtained by $\sigma_*^2 = 1/C$.