

# Accurate Fine-Tuning of Spatiotemporal Fourier Neural Operator for Turbulent Flows

<u>Shuhao Cao</u><sup>1</sup> Francesco Brarda<sup>2</sup> Ruipeng Li<sup>3</sup> Yuanzhe Xi<sup>2</sup>

<sup>1</sup>UMKC <sup>2</sup>Emory <sup>3</sup>LLNL

Efficient and Reliable Deep Learning Methods and their Scientific Applications Birs Workshop, June 2025

Question: which is more important in operator learning tasks?

Data or Model (training algorithm, neural architectures)?

# **Spatiotemporal Operator Learning**

# **Setting the Stage**

Toy model: approximating Navier-Stokes equations in a 2D periodic box

• (Velocity) Find  $\pmb{u} \in \pmb{H}^1(\mathbb{T}^2) \cap \{\pmb{v} \in \pmb{H}(\mathrm{div}): \, 
abla \cdot \pmb{v} = 0\}$ 

$$\partial_t u + u \cdot \nabla u - \nu \Delta u + \alpha u = f,$$

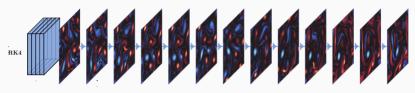
ullet (Vorticity-Streamfunction) Find  $\omega, \psi \in H^1(\mathbb{T}^2)$ 

$$\partial_t \omega + (\nabla^{\perp} \psi) \cdot \nabla \omega - \nu \Delta \omega + \alpha \omega = 0, \quad \omega + \Delta \psi = 0.$$

A long and rich history of numerical analysis for the NSE

- Projection schemes: Chorin (1968), Shen (1992).
- (Pseudo) Spectral: Orszag (1971), Orszag (1972), Tadmor (1987), Ku, Taylor, and Hirsh (1987), SHEN (1994).
- Finite element: GIRAULT and RAVIART (1986), BREZZI and FORTIN (1991), TEMAM (1995).
- Stability, convergence, time-stepping stencils: Heywood and Rannacher (1986), C. Wang and J.-G. Liu (2002), He and W. Sun (2007), GOTTLIEB et al. (2012), CHENG and C. WANG (2016).

### **Time-Marching Solvers for NSE**



• Crank-Nicolson for the diffusion term, explicit for the convection term, step size  $\tau$  has to be comparable to  $\mathcal{O}(\Delta x/\|\omega(t,\cdot)\|_{\infty})$  and  $o(\nu)$ .[1]

$$\left(I - \frac{\tau}{2}\nu L\right)\omega^{(t_{\ell+1})} = \left(I + \frac{\tau}{2}\nu L\right)\omega^{(t_{\ell})} - \tau\left(\nabla^{\perp}(-\Delta)^{-1}\omega^{(t_{\ell})}\right) \cdot \nabla\omega^{(t_{\ell})}.$$

 Aliasing error caused by the nonlinear term in pseudo-spectral methods: the modes extend "outside" of the approximation space. Possible remedies: filtering ("2/3-rule"<sup>[2]</sup>), higher-order<sup>[3]</sup>, or disspating the aliased modes.

<sup>[1]</sup> X. WANG (2012). "An efficient second order in time scheme for approximating long time statistical properties of the two dimensional Navier–Stokes equations". In: Numerische Mathematik.

<sup>[2]</sup> J. GOODMAN, T. Hou, and E. TADMOR (1994). "On the stability of the unsmoothed Fourier method for hyperbolic equations". In: Numerische Mathematik.

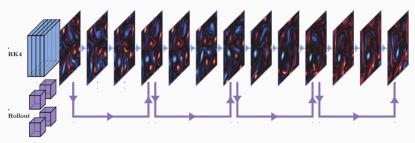
<sup>[3]</sup> S. NAGARAJAN, S. K. LELE, and J. H. FERZIGER (2003). "A robust high-order compact method for large eddy simulation". In: Journal of Computational Physics

## **Learning Operators between "Function Spaces"**

Operator-valued NN and operator learning:  $\mathcal{L}_{\theta} : \mathbb{R}^{n \times m \times d_{\text{in}}} \to \mathbb{R}^{n \times m \times d_{\text{out}}}$ .

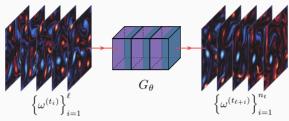
- How the "basis" (frames) are constructed and how are they aggregated?
  - o Latent frames are (nonlinear) universal approximators, and then are linearly aggregated through coefficients independent of the current latent space; Fourier Neural Operator: Z. Li, KOVACHKI, et al. (2021), and many others.
  - Frames/"basis" are linear projections of the current latent representations, then aggregated nonlinearly (input-dependent kernel integral). Nonlocal kernel: Buades, Coll, and Morel (2005), Gilboa and Osher (2007).
     Transformer/Attention: C. (2021), Z. Li, Meidani, and Farimani (2023), Hao et al. (2023), Bartolucci et al. (2024), Wu et al. (2024), Yu et al. (2024), and many others.
  - Both aggregation and frames are nonlinear. DeepONet: Lu et al. (2021), S. WANG, H. WANG, and PERDIKARIS (2021).
- Neural super-resolution operator. Spatial: Kochkov et al. (2021); temporal: Z. Sun, Yang, and Yoo (2023).

# **Learning Time-Marching with Data**



- Train an "autoregressive" operator learner  $G_{\theta}$  that maps  $\{\omega(t,\cdot)\}_{t\in(t_{\ell+1},t_{\ell+10})}$  to  $\omega(t_{\ell+11},\cdot)$  for different  $\ell$ s. This procedure repeats a "roll-out" prediction until certain steps.
- These "time-slices" are treated as the input "channels" in neural operators, such as FNO and its variants, Transformer-based NOs (Galerkin, GNOT, Oformer, DPOT, Transolver, many others).
- Time steps can be pretty big (e.g.,  $=100\tau$ ), yet the relative difference with the ground truth is horrible in evaluation ( $\approx 1\%$ ), no stability at all.

# **Learning Maps between Bochner Spaces**



#### **Problem of interest**

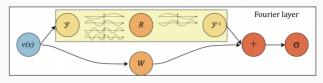
Construct an operator-valued neural network  $G_{\theta}$  to "approximate" G:

$$G: L^{2}(t_{1}, t_{\ell}; \mathcal{V}) \to L^{2}(t_{\ell+1}, t_{\ell+n_{t}}; \mathcal{V}), \quad \{\omega(t, \cdot)\}_{t \in (t_{1}, t_{\ell})} \mapsto \{\omega(t, \cdot)\}_{t \in (t_{\ell+1}, t_{\ell+n_{t}})}.$$

$$G_{\theta}: \mathcal{S}_{\ell} \to \mathcal{S}_{n_{t}}, \quad \mathbb{R}^{n \times n \times \ell} \ni \mathbf{w}_{\text{in}} \mapsto \mathbf{w}_{\text{out}} \in \mathbb{R}^{n \times n \times n_{t}}.$$

- $\mathcal{V}:=H^1(\mathbb{T}^2)$  or  $\{v\in H^1(\mathbb{T}^2):\nabla\cdot v=0\}.$
- $S_n \simeq \prod_{j=1}^n S_j$ , where  $S_j$  is a finite-dimensional approximation space of  $V_j$ .
- n,  $\ell$ ,  $n_t$  can all vary for training and evaluation.

### **Fourier Neural Operator**



Source: Figure 1 from the FNO paper by A. Stuart with his collaborators and students.<sup>[4]</sup>

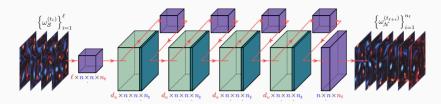
• The matrix-valued kernel matrix  $R_{ heta}(x)$  in the Fourier multiplication operator SpConv is learned from data

$$v^{\text{out}}(x) = \int_{\Omega} \kappa_{\theta}(x - y)v(y)dy = \mathcal{F}^{-1}\mathcal{F}\left(\int_{\Omega} \kappa_{\theta}(x - y)v(y)dy\right)$$
$$= \mathcal{F}^{-1}\left(\mathcal{F}(\kappa_{\theta}) \cdot \mathcal{F}(v)\right)(x) \approx : \mathcal{F}^{-1}(R_{\theta} \cdot \mathcal{F}(v))(x)$$

• Frequency truncation in parametrization: for 2D problems,  $v \in \mathbb{R}^{H_i \times n \times n}$ ,  $R_\theta \in \mathbb{C}^{H_\theta \times H_i \times d \times d}$  with  $d \ll n$ . Outputs can be viewed as learnable reduced "basis", e.g., n=128 and d=12 (# channel = # basis).

<sup>[4]</sup> Z. Li, N. B. Kovachki, et al. (2021). "Fourier Neural Operator for Parametric Partial Differential Equations". In: International Conference on Learning Representations.

### **Fourier Neural Operator**



The architectural schematics of FNO for NSE: red represents fixed dimensions; blue represents dimensions that accept arbitrary-sized discretizations.  $\blacksquare$ : spectral convolution layer;  $\blacksquare$ : pointwise nn.Conv3d that works as channel expansion/reduction;  $\blacksquare$ : pointwise nonlinearity.

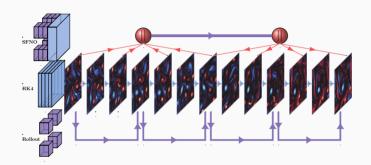
ullet The novelty of FNO for NSE: for the input tensor of dimension  $n imes n imes n_t$ 

channel dimension  $d_v \leftarrow$  time steps in the temporal dimension  $n_t$ .

However, this renders FNO unable to represent data pair in Bochner spaces. The number of parameters depends on the size of time discretization.

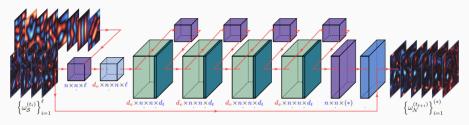
• The data for training and evaluation are prepared by applying a pointwise Gaussian normalizer that depends on the output steps  $n_t$ .

# **Toward Spatiotemporal Neural Operators**



- Goal: construct an operator-valued neural network with arbitrary spatial- or temporal-grid sizes input/output.
- Example: if the end-to-end pipeline picks the time slices in a time span of 0.5 as input, the slices in the subsequent 0.5 as the output, e.g., training data pair with the input/output tensor of dimension (64, 64, 10), then the evaluation should be able to deal with input size such as (128, 128, 40) or (512, 512, 100).

# **Spatiotemporal Fourier Neural Operator 3D**

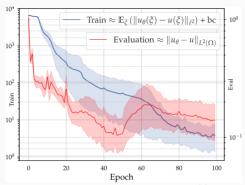


The architectural schematics of ST-FNO:  $\square$ : layer normalization after concatenation with positional encodings.

- Concatenation of random projection of the positional encodings with the input tensor act as a "channel expansion" (positional encodings ⊕ the input data, then go through a linear layer), which is a depth-wise linear combination.
- Layer normalization works as a dimension-agnostic normalizer.
- The new output operator has an extra single-channel spectral convolution layer with less truncated modes. It maps the latent time step dimension  $(d_t)$  to a given output time steps using FFT+iFFT's natural super-resolution by zero-padding the temporal steps.

# **How to Train Spatiotemporal FNO?**

### **Motivations: Function Representation**

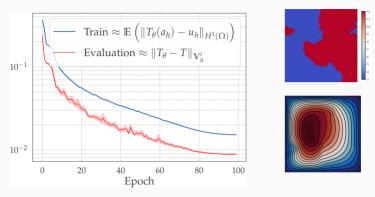


An extremely difficult problem due to the frequency principle<sup>[5]</sup> and the lottery ticket hypothesis<sup>[6]</sup>: train a function representer in low-dimensional spaces with a small number of parameters with no "labels". Trying to learn  $-\Delta u = 2\pi^2 \sum_{k \in \{1,4,16\}} \sin(k\pi x) \sin(k\pi y)$ . 1 epoch = 128 ADAM iterations. Error bars are plotted using different seeds (initialization).

<sup>[5]</sup> B. RONEN, D. JACOBS, Y. KASTEN, and S. KRITCHMAN (2019). "The Convergence Rate of Neural Networks for Learned Functions of Different Frequencies". In: Advances in Neural Information Processing Systems.

<sup>[6]</sup> J. Frankle and M. Carbin (2019). "The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks". In: International Conference on Learning Representations.

### **Motivations: Operator Learning**



A not-so-difficult problem: "typical" convergence results for end-to-end operator-valued neural network training and evaluations<sup>[7]</sup> for a 2D benchmark problem of porous media<sup>[8][9]</sup>.

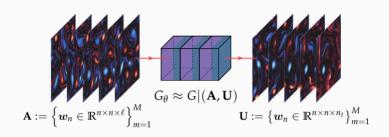
Cao-Brarda-Li-Xi

<sup>[7]</sup> C. (2021). "Choose a Transformer: Fourier or Galerkin". In: Advances in Neural Information Processing Systems (NeurIPS).

<sup>[8]</sup> A. P. ROBERTS and M. TEUBNER (1995). "Transport properties of heterogeneous materials derived from Gaussian random fields: bounds and simulation". In: *Physical Review E.* 

<sup>[9]</sup> N. H. NELSEN and A. M. STUART (2021). "The random feature model for input-output maps between Banach spaces". In: SIAM Journal on Scientific Computing.

### **Are Neural Operators Really Learning Operators?**



Given the training data pairs  $\{(a_m, u_m)\}_{m=1}^M$ , the operator learning problem is a Bayesian inverse problem with a linear or nonlinear operator as the unknown object to be inferred from data.

- DE HOOP, KOVACHKI, NELSEN, and STUART (2023)[10].

• With the "right" data, the Bayesian inverse problem is well-conditioned.

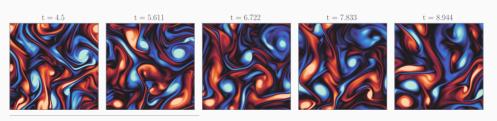
[10] M. V. DE HOOP, N. B. KOVACHKI, N. H. NELSEN, and A. M. STUART (2023). "Convergence rates for learning linear operators from noisy data". In: SIAM/ASA Journal on Uncertainty Quantification.

Cao-Brarda-Li-Xi Fine-tuning Spatiotemporal FNO

## Statistical Property of NSE: Energy Cascade

- Inverse cascade: a flux of energy from smaller scales (high frequency) to larger scales (low frequency)[11][12].
- Direct cascade: if the energy dissipation is caused by larger scale vortices inducing vortex stretching via viscosity, then a stationary regime can be established<sup>[13]</sup>.

$$\mathcal{E}(t,k) = \sum_{k-\delta k \leq |k| \leq k+\delta k} |\widehat{\nabla \times u}(t,k)|^2 \sim \mathcal{O}(k^{-\beta}) \quad \text{ after } t > t_0$$

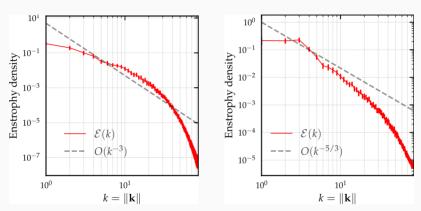


[11] A. N. KOLMOGOROV (1941). "The local structure of turbulence in incompressible viscous fluid for very large Reynolds". In: Numbers. In Dokl. Akad. Nauk SSSR.

[12] D. KOCHKOV et al. (2021). "Machine learning-accelerated computational fluid dynamics". In: Proceedings of the National Academy of Sciences

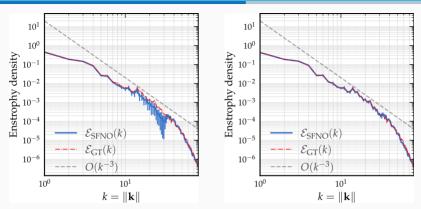
[13] J. C. McWilliams (1984). "The emergence of isolated coherent vortices in turbulent flow". In: Journal of Fluid Mechanics.

# Statistical Property of NSE: Energy Cascade



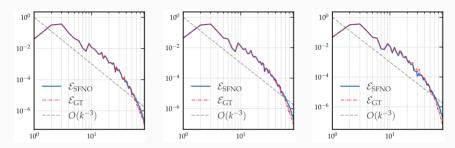
The plots of  $\mathcal{E}(t,k)$  for all the training samples in the case of the direct cascade (left) and the inverse cascade (right). Both examples' initial conditions are sampled from fixed random fields, respectively. The error bars are plotted with +/-10 times the standard deviation from the mean to boost the visibility.

### **Neural Operators Learn Fast**



**1 epoch** vs **10 epochs**: Enstrophy spectrum density comparisons to illustrate the "convergence" of ST-FNO's training. There are 10 training runs on a  $64 \times 64$  grid starting from 10 different seeds. The evaluation is on a  $256 \times 256$  grid for a fixed randomly chosen sample. The error bars are plotted with +/- 10 times the standard deviation from the mean to boost the visibility.

# Neural Operators Learn Low-Frequency Fast like really fast but not that accurate



"Spectral bias/frequency principle"<sup>[14]</sup>: not only neural operators can capture the pattern of the low-frequency modes of the data well, they can also learn the low modes really fast. The plots show the enstrophy spectrum density comparison for a randomly selected trajectory for a given trained SFNO after only 10 epochs. However, the magnitude of the error/difference still dominates in low-frequency for modeling turbulence<sup>[15]</sup>.

<sup>[14]</sup> J. ZHI-QIN Xu et al. (2020). "Frequency Principle: Fourier Analysis Sheds Light on Deep Neural Networks". In: Communications in Computational Physics.

<sup>[15]</sup> P. LIPPE et al. (2023). "PDE-Refiner: Achieving Accurate Long Rollouts with Neural PDE Solvers". In: Thirty-seventh Conference on Neural Information Processing Systems

### Do we really need 500 epochs of training?

Table 1: Benchmarks on Navier Stokes (fixing resolution  $64 \times 64$  for both training and testing)

Config	Parameters	Time per	$\begin{array}{c} \nu = 1\mathrm{e}{-3} \\ T = 50 \end{array}$	$\begin{array}{c} \nu = 1\mathrm{e}{-4} \\ T = 30 \end{array}$	$\begin{array}{c} \nu = 1\mathrm{e}{-4} \\ T = 30 \end{array}$	$\begin{array}{c} \nu = 1\mathrm{e}{-5} \\ T = 20 \end{array}$
		epoch	N = 1000	N = 1000	N = 10000	N = 1000
FNO-3D	6,558,537	38.99s	0.0086	0.1918	0.0820	0.1893
FNO-2D	414,517	127.80s	0.0128	0.1559	0.0834	0.1556
U-Net	24,950,491	48.67s	0.0245	0.2051	0.1190	0.1982
TF-Net	7,451,724	47.21s	0.0225	0.2253	0.1168	0.2268
ResNet	266, 641	78.47s	0.0701	0.2871	0.2311	0.2753

Source: Table 1 from the original FNO paper. FNO3d is trained 500 epochs ( $500 \times 128$ mini-batch ADAM iterations). However, SFNO reaches  $1 \times 10^{-2}$  relative difference evaluated on a 256  $\times$  256 grid with the ground truth in 5 to 15 epochs, and  $6\times10^{-3}$  after 500 epochs. More sophisticated Transformers with multilevel feature aggregation<sup>[16]</sup> can drive the error down further to  $4 \times 10^{-3}$  level but not any further.

Cao-Brarda-Li-Xi

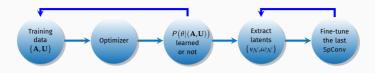
# Fine-tuning/Post-processing

### **Fine-tuning to Improve Accuracy**

- Make certain parameters in a spectral convolution layer in FNO correspond to the coefficients of a basis in the spectral domain (the last single-channel spectral conv layer).
- Use a stronger norm to train to learn low-frequency modes ( $L^2$ -based loss with aggressive frequency truncation), use a weaker norm to post-process/fine-tune (negative Sobolev norm with less frequency truncation).
  - $\circ$  The frequency truncation in FNO  $\approx$  Tikhonov regularization (HANSEN, NAGY, and O'LEARY 2006).
  - Negative Sobolev norm as an implicit regularization in inverse problems:
     OSHER, SOLÉ, and L. VESE (2003) and LIEU and L. A. VESE (2008), ZHU, HU, LOU, and YANG (2024).
- The negative Sobolev norm of  $f \in L^2(\mathbb{T}^2)/\mathbb{R}$  can be handily computed using FFT, and this negative norm happens to be the correct functional norm to measure the residual:  $R(\omega_{\mathcal{N}})(v) := \langle R(\omega_{\mathcal{N}}), v \rangle$ , denote  $R(\omega_{\mathcal{N}})(t, \cdot) =: r(t, \cdot)$

$$||r||_{\mathcal{H}'} \simeq |r|_{-1} := \sum_{k \in 2\pi \mathbb{Z}_v^2 \setminus \{0\}} |k|^{-2} |\hat{r}(k)|^2.$$

### **New Paradigm**



- Train the neural operator only for a few epochs (e.g., 10) using a strong norm as the loss until it learns the frequency signature of the data (e.g., the energy cascade of the NSE).
- Extract the latent tensor up to the **last** channel-reduction **LINEAR SpConv** layer  $Q_{\theta}(\cdot)$ , such that for  $i=1,\cdots,n_t$

Neural representation of 
$$\omega^{(t_{\ell+i})} = \omega_{\mathcal{S}}^{(t_{\ell})} + Q_{\theta}(\mathbf{v}_{\mathrm{latent}})$$
.

- Remove the frequency truncation and fine-tune this layer  $Q_{\theta}$  **ONLY** using an optimizer with a weaker norm.
- $\partial_t \omega_{\mathcal{N}}$  is approximated using an IMEX solver implemented as an nn.Module.
- Searching for the best possible  $\theta^*$  under a functional norm (negative norm) is equivalent to solving a variational problem in the positive Sobolev norm.

#### Theorem (Reliable and efficient a posteriori error estimation)

Under certain smoothness assumptions for  $\omega \in \mathcal{H}$ , if the fine-tuning problem can be solved exactly, for  $m=\ell+1,\cdots,\ell+n_t-1$ ,  $\mathcal{T}_m:=[t_m,t_{\ell+n_t}]$ , the PDE residual

$$R(\omega) := \operatorname{rot} f - \partial_t \omega - (\nabla^{\perp} (-\Delta)^{-1} \omega \cdot \nabla) \omega + \nu \Delta \omega$$

is efficient in representing the error:

$$\|R(\omega_{\mathcal{N}})\|_{L^{2}(\mathcal{T}_{m};\mathcal{V}')}^{2} \lesssim \|\omega - \omega_{\mathcal{N}}\|_{L^{2}(\mathcal{T}_{m};\mathcal{V})}^{2} + \|\partial_{t}(\omega - \omega_{\mathcal{N}})\|_{L^{2}(\mathcal{T}_{m};\mathcal{V}')}^{2} + \text{h.o.t.}$$

When  $\omega_N$  is sufficiently close to  $\omega$ ,

$$\|\omega - \omega_{\mathcal{N}}\|_{L^{\infty}(\mathcal{T}_m; L^2)}^2 + \|\omega - \omega_{\mathcal{N}}\|_{L^2(\mathcal{T}_m; H^1)}^2 \leq \|(\omega - \omega_{\mathcal{N}})(t_m, \cdot)\|_{H^1}^2 + C \int_{\mathcal{T}_m} \|R(\omega_{\mathcal{N}})(t, \cdot)\|_{\mathcal{V}'}^2 dt.$$

 The optimization of this is not tied to local adaptive mesh refinement, the loss can be simply evaluated globally in the spectral domain to refine the (reduced) spectral basis, similar to ROM (PATERA and ROZZA 2007).

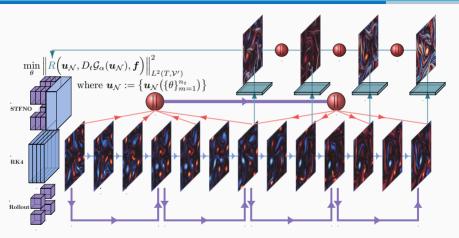
## **New Paradigm**

- (Step 1) Train  $G_{\Theta} := Q_{\theta} \circ F_{\Theta \setminus \theta}$  using  $\sum_{\omega \sim p} \big\{ \|G_{\Theta}(\omega_{\mathcal{S}}) \omega_{\mathcal{S}}\|_0^2 + \gamma \|\Theta\|_0^2 \big\}.$
- (Step 2) Stop training after the relative difference reached a threshold on the validation sets.
- (Step 3) Fine-tune (post-process) the **last** layer  $Q_{\theta}$  where  $\theta_{\text{train}} \in \mathbb{C}^{1 \times H \times d \times d \times d_t}$  and  $\theta_{\text{ft}} \in \mathbb{C}^{1 \times (H+1) \times \frac{2n}{3} \times \frac{2n}{3} \times n_t}$  for n the number of grid points in each axis of an evaluation sample,  $n_t$  evaluation time steps (e.g., 40).
  - Approximating the time derivative by  $\psi^{(m)}=(-\Delta_h)^{-1}\omega^{(m)}$  with  $\{\omega^{(m)}\}$  being the output of  $G_{\theta}$ , and

$$\frac{3\omega^{(m+1)} - 4\omega^{(m)} + \omega^{(m-1)}}{2\left(\delta t\right)} + \nabla^{\perp}(2\psi^{(m)} - \psi^{(m-1)}) \cdot \nabla\left(2\omega^{(m)} - \omega^{(m-1)}\right) - \nu\Delta\omega^{(m+1)} = f^{(m+1)}.$$

- Apply an optimizer to find  $\theta_{\mathrm{ft}}$  by minimizing  $\sum_{m\in\{\ell+1,\dots,\ell+n_t\}} \left\| R(\omega^{(m)}(\theta_{\mathrm{ft}})) \right\|_{-1}^2$ .
- (Output) A trajectory that satisfies the statistical signature and minimizes the residual while not having to march  $\mathcal{O}(1/\tau)$  steps.

# **New Paradigm**



Schematics comparison: the new method a shares striking mathematically resemblance with a parallel-in-time two-grid method using a learned reduced basis. On a very "coarse" temporal grid, the spatiotemporal surrogate model gives a very good "initial" guess to perform implicit residual smoothing.

# Numerical Experiments

#### **Taylor-Green Vortex**

$$u(t, x, y) = e^{-2\kappa^2 \nu t} \begin{pmatrix} \sin(\kappa x) \cos(\kappa y) \\ -\cos(\kappa x) \sin(\kappa y) \end{pmatrix} \text{ with } \nu = 10^{-2}.$$

- Training dataset: 10 samples with  $\kappa=1,\ldots,10$ ; evaluation: 1 sample with  $\kappa=11$ .
- Trajectories are randomly sampled before the breakdown phase.
- Ground truth: pseudo-spectral discretization, second-order Runge-Kutta for the explicit part, Crank-Nicolson for diffusion part with  $\Delta t = 4 \times 10^{-3}$ .
- Training is done on  $64\times64$  for 5 epochs, the evaluation is on  $256\times256$ , time step is 10 for training, 40 for evaluation. ST-FNO's parameters have capacity to represent the discrete solutions in the eval set exactly.

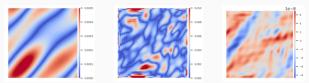
Results for Taylor-Green vortex example  $\epsilon:=\omega_{\mathsf{true}}-\omega_{\mathcal{N}}$ , the errors at the final time step.

	Evaluation after training		After fine-tuning		
	$\ \varepsilon\ _{L^2}$	$  R  _{-1,n}$	$\ \varepsilon\ _{L^2}$	$  R  _{-1,n}$	
FNO3d	$1.84 \times 10^{-1}$	$5.40 \times 10^{-1}$	N/A	N/A	
ST-FNO3d	$1.94 \times 10^{-1}$	$2.18 \times 10^{-1}$	$1.24 \times 10^{-6}$	$3.21 \times 10^{-7}$	
PS+RK2 (GT)	$5.91 \times 10^{-6}$	$1.16 \times 10^{-5}$	N/A	N/A	

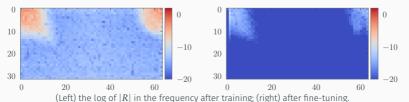
#### Forced Turbulence: FNO Data

Results for forced turbulence (original toy model example from the FNO paper,  $\nu=1\times10^{-3}$ ).

	Evaluation after training		After fine-tuning	
	$\ \varepsilon\ _{L^2}$	$  R  _{-1,n}$	$\ \varepsilon\ _{L^2}$	$  R  _{-1,n}$
FNO 100 ep ST-FNO 10 ep + <i>L</i> <sup>2</sup> FT	$1.31 \times 10^{-2}$ $1.02 \times 10^{-2}$		N/A 2.82×10 <sup>-4</sup>	N/A 2.78×10 <sup>-5</sup>
ST-FNO 10 ep + $H^{-1}$ FT	-	-	$3.16 \times 10^{-4}$	$4.59 \times 10^{-7}$

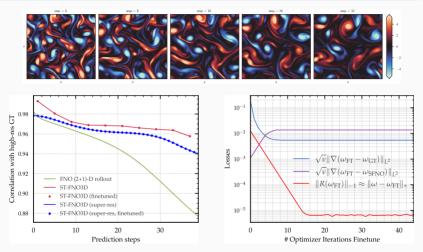


Pointwise residuals for the ground truth streamfunction, ST-FNO inference, and ST-FNO inference after fine-tuning.



Fine-tuning Spatiotemporal FNO

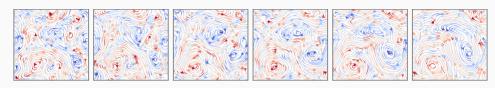
#### Isotropic Turbulence: Direct Cascade (McWilliams 1984)



(Left) the correlation with the highest-resolution DNS ground truth, one step  $=55\Delta t$ . (Right) Fine-tuning convergence.

# Torch-CFD

#### github.com/scaomath/torch-cfd



A native PyTorch port of Google Research's Jax-CFD (Jax-CFD has not been in active development since Aug 2023) with many new and enhanced features.

- Fully implemented as "tensor-in-tensor-out" with a batch dimension (b, c, \*), more user-friendly for tensor-tracking debugger (VSCode).
- Gridded variables implemented natively as a subclass of torch.Tensor with their ops using \_\_torch\_function\_\_, perfect for staggered grids.
- Modular designs: time-marching schemes, multigrid solvers, different flux schemes, also implemented as nn.Module with learnable components (Butcher's tableau, smoothers, weights for fluxes interpolations). Taking advantage of handy methods such as \_forward\_hooks and \_backward\_hooks.

10

13

14

17

```
class MultigridSolver(nn.Module):
 def __init__(self) -> None:
    super(). init ()
    self.smoothers: nn.ModuleList = ...
    self.operators: nn.ModuleList = ...
  def forward(self. f. u. lvl):
     u = self.smoothers[lvl](f, u)
     r = f - self. apply op(u. self.operators[lvl])
      if lvl == self.levels - 1:
          e = self. coarse solve(r)
     else:
          rc = self.restrict(r)
          ec = self.forward(rc, 0, lvl=lvl + 1)
          e = self.prolong(ec)
     u += e
      return self.smoothers[lvl](f, u)
```

### **Acknowledgments**

- Organizers: Jack Xin (UC Irvine), Gitta Kutyniok (LMU Munich), Stanley Osher (UCLA), Bao Wang (University of Utah), Andrea Bertozzi (UCLA).
- Co-authors, as well as Long Chen (UC Irvine), Ludmil Zikatanov (back at Penn State), Ari Stern (WashU).
- Source codes and data to replicate the experiments are available at
- github.com/scaomath/torch-cfd/examples 🚱 🚭 10.57967/hf/2470
- References
  - C., F. BRARDA, R. Li, and Y. Xi (2025). "Spectral-Refiner: Accurate Fine-Tuning of Spatiotemporal Fourier Neural Operator for Turbulent Flows". In: The Thirteenth International Conference on Learning Representations (ICLR) cs.LG:2405.17211.
- This research is supported in part by NSF award DMS-2309778.