



# Temporal Learning Capacity of Transformers in Non-Markovian Dynamical Systems

---

Benjamin Shih

Zhongqiang Zhang, George Em Karniadakis, Khemraj Shukla

April 26, 2024

CRUNCH Group | Brown University

## 1. Introduction and Background

Derivatives: Tempered & Fractional, Operator Learning, Attention

## 2. Methodology

Fourier Attention, Architecture, Training

## 3. Experimental Results

TF- $\{\text{Smooth, LIF, Diffusion}\}$ , Attention Matrices

## 4. Conclusion

Summary, Ongoing Work

# Introduction and Background

---

## Definition ( $\sigma$ -Tempered fractional derivative)

The  $\sigma$ -tempered fractional Caputo derivative of fractional order  $\alpha$  is defined as

$${}^C_a\mathcal{D}_t^{\alpha,\sigma} u(t) = \frac{e^{-\sigma t}}{\Gamma(n-\alpha)} \int_a^t (t-s)^{n-\alpha-1} \frac{d^n}{ds^n} (e^{\sigma s} u(s)) ds$$

Clearly, if  $\sigma = 0$ , then this reduces to the standard Caputo fractional derivative with fractional order  $\alpha$ .

## Integer order

$$\tau \frac{dV}{dt} = -(V - V_{\text{rest}}) + RI(t)$$

$R$  cell constant

$I(t)$  spiking forcing term

$\tau$  membrane time constant

$V_{\text{rest}}$  resting membrane potential

$V$  membrane potential

## $\sigma$ -Tempered Fractional

$$\tau {}_a^C \mathcal{D}_t^{\alpha, \sigma} [V](t) = -(V - V_{\text{rest}}) + RI(t)$$

Major advantage:

- Non-Markovian (memory!): captures effect of residual current

- Input:  $f(t)$ , output:  $u(t)$ . We denote the operator mapping  $f$  to  $u$  as

$$\mathcal{G}^* : \mathcal{F} \ni f \mapsto u \in \mathcal{U}$$

- On computers: use discrete samples of continuous function at “sensors.”

$$f_O = \{f(t_1), \dots, f(t_n)\}$$

- Operator learning: construction of a network

$$\mathcal{G}_\theta : \mathbb{R}^n \ni f_O \mapsto \mathcal{G}_\theta[f_O](\cdot) \in \mathcal{U}$$

such that  $\mathcal{G}_\theta$  approximates  $\mathcal{G}^*$ . We want

$$\arg \min_{\theta \in \Theta} \|\mathcal{G}_\theta[f](\cdot) - \mathcal{G}^*[f](\cdot)\|$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right) V$$
$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

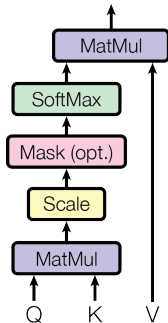


Figure 1: Scaled Dot-Product Attention<sup>1</sup>

<sup>1</sup>Ashish Vaswani et al. “Attention is All you Need”. In: Advances in Neural Information Processing Systems. Ed. by I. Guyon et al. Vol. 30. Curran Associates, Inc., 2017.

# Methodology

---

Standard:  $\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right) V$

Fourier type:  $(z_i)_j = \frac{1}{n} \sum_{s=1}^n (q_i \cdot k_s)(v^j)_s \approx \int_{\Omega} \kappa(x_i, \xi) v_j(\xi) d\xi$

# Network Architecture & Loss

$$\mathcal{L}(\theta) = \frac{1}{m} \sum_{k=1}^m \frac{\|\mathcal{G}_\theta[f_k; \alpha, \sigma] - \mathcal{G}^*[f_k]\|_{\ell^2}}{\|\mathcal{G}^*[f_k]\|_{\ell^2}}$$

where

$$\|f\|_{\ell^2} = \left( \frac{1}{n} \sum_{i=1}^n |f(t_i)|^2 \right)^{\frac{1}{2}}$$

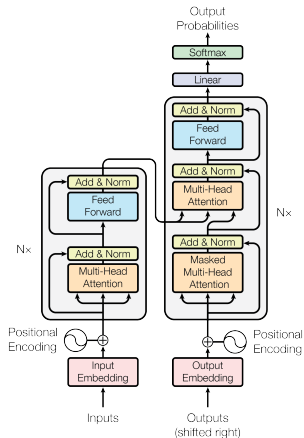


Figure 2: Transformer Architecture<sup>2</sup>

<sup>2</sup>Vaswani et al., "Attention is All you Need".

# Experimental Results

---

$$\begin{cases} {}^C_a\mathcal{D}_t^{\alpha,\sigma}[u](t) = \mu u(t) + f(t) \\ u(0) = u_0 \end{cases}$$

where  $f \sim \mathcal{GP}^3$ . Operator:

$$\mathcal{S} : (f(t), \alpha, \sigma) \mapsto u(t; \alpha, \sigma)$$

- Case 1:  $\mu = 0$ , vary  $\alpha, f$ , no tempering
- Case 2:  $\mu = -1$ , vary  $\alpha, f$ , no tempering
- Case 3:  $\mu = -1$ , vary  $\alpha, \sigma, f$

---

<sup>3</sup> $\mathcal{GP}$  specifics:  $\mathcal{N}(\mu, \sigma^2(-\Delta + \tau^2 I)^\gamma)$  with  $\mu = 1, \gamma = 2.5, \tau = 7, \sigma = 7$ ,

# Smooth TFIVP Case 1: Visually

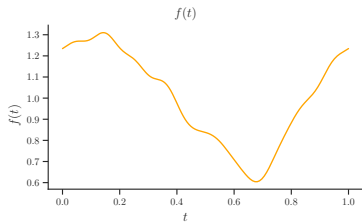


Figure 3: Input: Smooth forcing term

$\mathcal{S}$

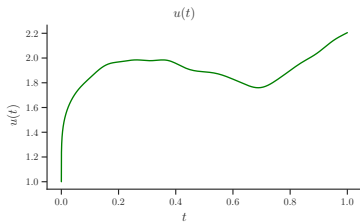


Figure 4: Output: Solution

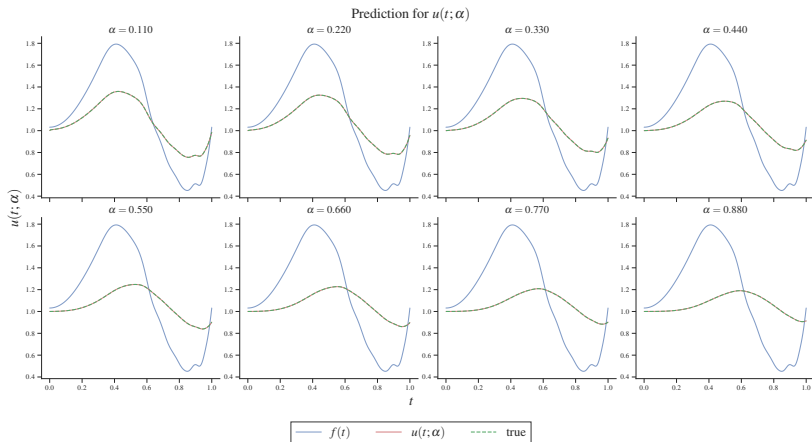


Figure 5: Results for TFIVP Case 1 (varying  $\alpha$  and  $f$ )

$$\tau_a^C \mathcal{D}_t^{\alpha, \sigma} [V](t) = -(V - V_{\text{rest}}) + RI(t)$$

Operator:

$$\mathcal{S} : (I(t), \alpha, \sigma) \mapsto V(t; \alpha, \sigma')$$

- Case 1: vary  $\alpha$ , no tempering
- Case 2: vary  $\alpha, I(t)$  (both location and intensity)
- Case 3: vary  $\alpha, \sigma$

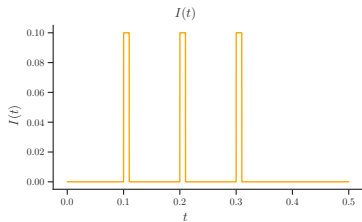


Figure 6: Input: Spiking current

$\downarrow$   
 $\sigma$

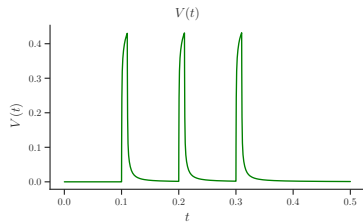


Figure 7: Output: Membrane Potential

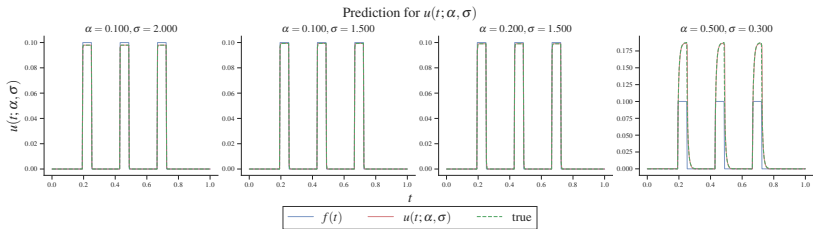


Figure 8: Results for LIF Case 3 (varying  $\alpha$  and  $\sigma$ )

$$\begin{cases} {}^C\mathcal{D}_t^{\alpha,\sigma}[u](x,t) = \mu \frac{\partial^2 u}{\partial x^2}(x,t) + f(x,t) \\ u(x,0) = u_0 \\ u(0,t) = u(1,t) = 0 \end{cases}$$

Operator:

$$\mathcal{S} : (f(x,t), \alpha, \sigma) \mapsto u(x,t; \alpha, \sigma)$$

- vary  $\alpha, f$

# TF Inhomogenous Diffusion: Visually

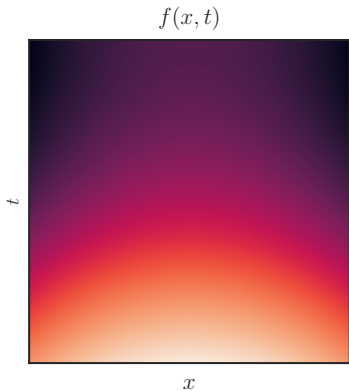


Figure 9: Input: Source term

$s$

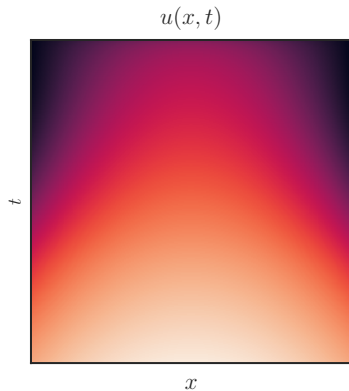


Figure 10: Output: Diffusive behavior

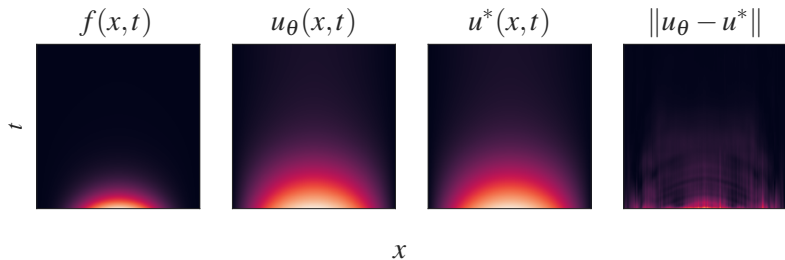


Figure 11: Results for Diffusion Case 1 (varying  $\alpha, f$ )

# Sample Attention Matrix Visualization: TF-LIF Case 1

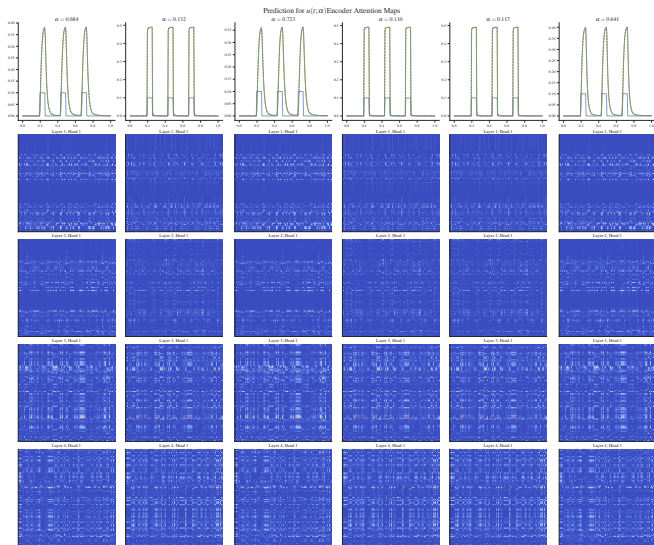


Figure 12: Transformer Encoder Attention: TF-LIF Case 1

Table 1: Experiment details and relative  $L_2$  error for each case of the tempered fractional smooth, LIF, and diffusion problems for the two different models tested.

		Experiment settings							Results	
Model	Experiment	Case	$\alpha$	$\sigma$	$n$	Learning Rate	Batch Size	Iterations	$\mathcal{L}(\theta)$	
Adam/cycle	Smooth	1		—					$4.7 \times 10^{-4}$	
		2	[0.11, 0.99]	—	201	$1.0 \times 10^{-3}$	32	$1.0 \times 10^6$	$9.2 \times 10^{-4}$	
		3		[0.11, 2]					$6.5 \times 10^{-4}$	
	LIF	1			—	204				$3.9 \times 10^{-4}$
		2	[0.11, 0.99]		—	200	$1.0 \times 10^{-4}$	32	$2.0 \times 10^6$	$1.3 \times 10^{-2}$
		3			[0.11, 0.99]	314				$1.3 \times 10^{-3}$
Diffusion	1	[0.11, 0.99]		—	$257^2$	$1.0 \times 10^{-4}$	32	$5.0 \times 10^4$	$6.1 \times 10^{-3}$	
Lion/polynomial	Smooth	1		—					$5.5 \times 10^{-5}$	
		2	[0.11, 0.99]	—	201	$1.0 \times 10^{-4}$	2000	$1.0 \times 10^6$	$7.9 \times 10^{-5}$	
		3		[0.11, 0.99]					$1.1 \times 10^{-2}$	
	LIF	1			—	204				$3.9 \times 10^{-5}$
		2	[0.11, 0.99]		—	200	$1.0 \times 10^{-4}$	1000	$2.0 \times 10^5$	$3.0 \times 10^{-1}$
		3			[0.11, 2]	314				$3.4 \times 10^{-2}$
Diffusion	1	[0.11, 0.99]		—	$257^2$	$1.0 \times 10^{-4}$	32	$5.0 \times 10^4$	$8.6 \times 10^{-2}$	

# Sample Loss Trajectory Comparison (Smooth Case 2)

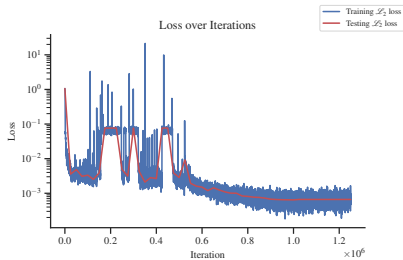


Figure 13: Adam/1cycle loss curve

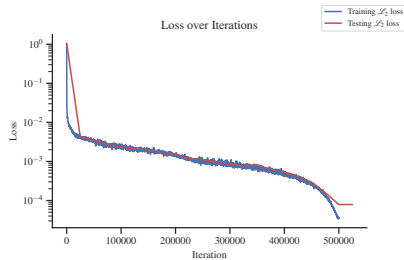


Figure 14: Lion/polynomial loss curve

## Conclusion

---

- Novel application of the transformer architecture to the operator learning problem; specifically to tempered fractional operators
- Success on 3 different problems that vanilla DeepONet was unable to learn
- Visualization and interpretation of the attention in each of these models

- RiemannONet: success with massive jump discontinuities by decomposition in the trunk net:

$$\mathcal{L}^t(\Theta) = \left\| \sum_l T_l(\Theta) \mathcal{A}_l - U^* \right\|_{L_2}$$

↓

$$Q^* R^* = \text{qr}(T(\Theta))$$

- Comparison between different neural operators
- Ablation studies

Thank you!



Vaswani, Ashish et al. “Attention is All you Need”. In: Advances in Neural Information Processing Systems. Ed. by I. Guyon et al. Vol. 30. Curran Associates, Inc., 2017.