

# Data Continuity Matters

Improving Sequence Modeling with Lipschitz Regularizer

Eric Qu<sup>1,2</sup> Xufang Luo<sup>1</sup> Dongsheng Li<sup>1</sup>



<sup>1</sup>Microsoft Research

<sup>2</sup>Duke Kunshan University

May 3, 2023



## Sequence Models Works Well On Specific Tasks

Transformers

Text



Gene



State Space Models

Audio Time-series



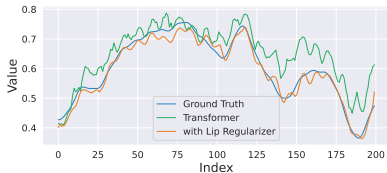
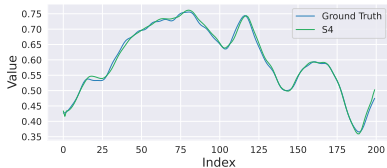
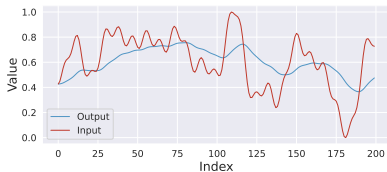
But Why?



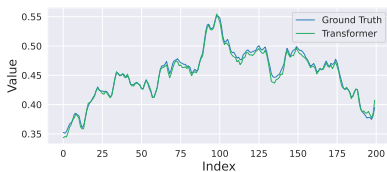
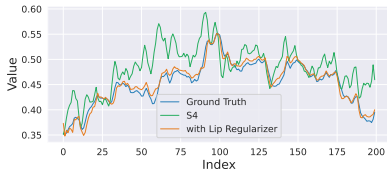
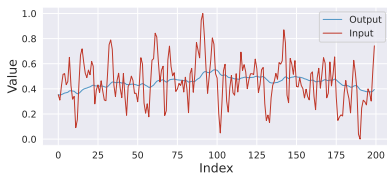
Sequence models have preferences in **Data Continuity**

- ▶ Transformers → Discrete Data
- ▶ State Space models → Continuous Data
- ▶ Proof of Concept Experiment
  - ▶ Generate discrete and continuous input sequence
  - ▶ Map them to the output sequence with the same function
  - ▶ Use Transformer and S4 Model to learn this mapping

# Motivation



(a) High Continuity



(b) Low Continuity



Sequence Models + Unpreferred Data Continuity



Deteriorated Performance

- ▶ Solution: a regularizer that alters input data continuity!
- ▶ Apply the regularizer to the input embedding



- ▶ Continuity Measure: Lipschitz Constant
- ▶ For a sequence  $x_0, x_1, \dots, x_n$ , view it as a sample of  $f(t)$ :

$$f(t_0) = x_0, f(t_1) = x_1, \dots, f(t_n) = x_n,$$

where  $t_0, t_1, \dots, t_n$  are time steps.

- ▶ The Lipschitz constant  $L_f$  of  $f(t)$  is

$$L_f = \max_{t_i, t_j \in \{0, 1, \dots, n\}} \frac{|f(t_i) - f(t_j)|}{|t_i - t_j|} = \max_{i, j \in \{0, 1, \dots, n\}} \frac{|x_i - x_j|}{|i - j|}.$$

- ▶ By Mean Value Theorem

$$L_f = \max_{i, j \in \{0, 1, \dots, n\}} \frac{|x_i - x_j|}{|i - j|} = \max_{k \in \{0, 1, \dots, n-1\}} |x_{k+1} - x_k|.$$



- ▶ To help with optimization, we introduce a surrogate:
  - ▶ Max  $\rightarrow$  Mean
  - ▶ L1 norm  $\rightarrow$  L2 norm

## Definition 1: Lipschitz Regularizer

Suppose the sequence is  $x_0, x_1, \dots, x_n$ . We define the Lipschitz Regularizer as follows:

$$\mathcal{L}_{\text{Lip}} = \frac{1}{n} \sum_{i=0}^{n-1} (x_{i+1} - x_i)^2 \quad (1)$$



- ▶ State Space models prefer continuous input
  - ▶ Assumption: input is a discrete sample of a continuous func
  - ▶ Higher input continuity  $\rightarrow$  Lower HiPPO Leg-S error rate
- ▶ Use LipReg to make input continuous
  - ▶ Introduce a 1D Convolution embedding layer
  - ▶ Apply LipReg on the embedding

$$\mathcal{L}(y, \hat{y}, \hat{l}) = \mathcal{L}_{S4}(y, \hat{y}) + \lambda \mathcal{L}_{\text{Lip}}(\hat{l})$$

	ListOps	Text	Retrieval	Image	Image-c	Path	Path-c	PathX	PathX-c
S4	59.53	86.51	91.07	88.54	84.27	<b>94.02</b>	89.11	<b>96.03</b>	92.41
S4 + Emb	58.94	87.12	90.28	87.25	85.13	92.37	90.32	93.87	92.81
<b>S4 + Emb + Lip</b>	<b>61.37</b>	<b>89.74</b>	<b>93.83</b>	<b>89.19</b>	<b>88.43</b>	93.52	<b>91.39</b>	95.72	<b>94.36</b>





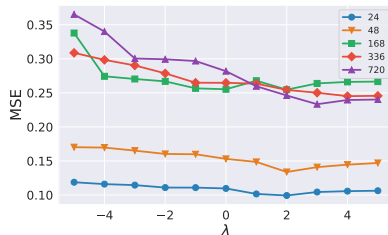
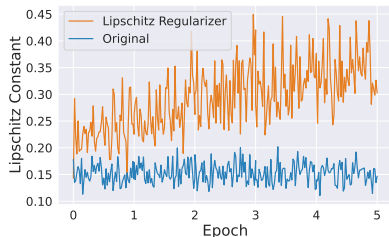
- ▶ Transformers prefer discrete input

$$\mathcal{L}(y, \hat{y}, \hat{l}) = \mathcal{L}_{\text{Transformer}}(y, \hat{y}) - \lambda \mathcal{L}_{\text{Lip}}(\hat{l})$$

Methods		Transformer		Transformer + Lip		Informer		Informer + Lip		Autoformer		Autoformer + Lip	
Metric		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETT <sub>h1</sub>	24	0.07047	0.20586	<b>0.07019</b>	<b>0.20570</b>	0.09842	0.24747	<b>0.08882</b>	<b>0.23674</b>	0.05567	0.18596	<b>0.05504</b>	<b>0.18495</b>
	48	0.18902	0.37046	<b>0.16716</b>	<b>0.34974</b>	0.15845	0.31907	<b>0.12615</b>	<b>0.28333</b>	0.07860	0.22324	<b>0.07422</b>	<b>0.21398</b>
	168	0.39773	0.55569	<b>0.30811</b>	<b>0.48183</b>	0.18314	0.34619	<b>0.10579</b>	<b>0.25552</b>	0.09232	0.24037	<b>0.08983</b>	<b>0.23544</b>
	336	0.41523	0.56902	<b>0.41324</b>	<b>0.56402</b>	0.22164	0.38720	<b>0.11810</b>	<b>0.26959</b>	0.10462	0.25484	<b>0.10461</b>	<b>0.25483</b>
	720	0.65586	0.75324	<b>0.62233</b>	<b>0.73160</b>	0.26883	0.43506	<b>0.13131</b>	<b>0.28731</b>	<b>0.12069</b>	<b>0.27791</b>	0.12394	0.27833
ETT <sub>h2</sub>	24	0.09449	0.24259	<b>0.07560</b>	<b>0.20989</b>	0.09309	0.24015	<b>0.08626</b>	<b>0.22559</b>	0.11136	0.26315	<b>0.09345</b>	<b>0.25515</b>
	48	0.15016	0.30996	<b>0.13229</b>	<b>0.29278</b>	0.15483	0.31445	<b>0.13684</b>	<b>0.28936</b>	0.15137	0.30316	<b>0.14945</b>	<b>0.30129</b>
	168	0.25197	0.41087	<b>0.21046</b>	<b>0.37453</b>	<b>0.23193</b>	<b>0.38947</b>	0.30071	0.43671	0.20403	0.35646	<b>0.18370</b>	<b>0.33714</b>
	336	0.22258	0.38170	<b>0.20867</b>	<b>0.37298</b>	0.26321	0.41659	<b>0.24875</b>	<b>0.40827</b>	0.22188	0.37417	<b>0.21195</b>	<b>0.36425</b>
	720	0.21932	0.38844	<b>0.18445</b>	<b>0.35793</b>	0.27722	0.43063	<b>0.23646</b>	<b>0.39648</b>	0.25612	0.40089	<b>0.25604</b>	<b>0.40085</b>
ETT <sub>m1</sub>	24	0.01279	0.08410	<b>0.01210</b>	<b>0.08312</b>	0.03016	0.13717	<b>0.01815</b>	<b>0.09147</b>	0.02317	0.11778	<b>0.02300</b>	<b>0.10107</b>
	48	0.08974	0.25869	<b>0.02872</b>	<b>0.12820</b>	0.06944	0.20255	<b>0.05848</b>	<b>0.19686</b>	0.04130	0.15783	<b>0.03931</b>	<b>0.15601</b>
	96	0.05341	0.17696	<b>0.05182</b>	<b>0.15017</b>	0.19414	0.37236	<b>0.13336</b>	<b>0.30091</b>	0.05432	0.18033	<b>0.05258</b>	<b>0.17605</b>
	288	0.22354	0.40455	<b>0.13780</b>	<b>0.29825</b>	0.40140	0.55355	<b>0.30266</b>	<b>0.46864</b>	0.11893	0.27181	<b>0.07521</b>	<b>0.21728</b>
	672	<b>0.40726</b>	<b>0.55824</b>	<b>0.40726</b>	0.55826	0.51164	0.64390	<b>0.27543</b>	<b>0.45377</b>	<b>0.09156</b>	0.23690	0.09280	<b>0.23621</b>
Weather	24	0.00223	0.03468	<b>0.00154</b>	<b>0.02497</b>	0.11676	0.25142	<b>0.11256</b>	<b>0.23844</b>	0.00740	0.06422	<b>0.00736</b>	<b>0.06329</b>
	48	0.00422	0.04106	<b>0.00292</b>	<b>0.03026</b>	<b>0.17822</b>	<b>0.31846</b>	0.19134	0.32408	0.01002	<b>0.07648</b>	<b>0.00978</b>	0.07727
	168	0.00537	0.05975	<b>0.00319</b>	<b>0.04464</b>	0.26585	0.39764	<b>0.25138</b>	<b>0.37400</b>	0.01038	0.07082	<b>0.00528</b>	<b>0.05638</b>
	336	0.00524	0.05772	<b>0.00417</b>	<b>0.03673</b>	0.29713	0.41571	<b>0.24748</b>	<b>0.37725</b>	0.00729	0.06492	<b>0.00566</b>	<b>0.05888</b>
	720	0.00933	0.07630	<b>0.00272</b>	<b>0.03823</b>	0.35875	0.46647	<b>0.26479</b>	<b>0.39214</b>	0.00960	0.08758	<b>0.00925</b>	<b>0.07136</b>
Count	2		39		4		36		4		36		

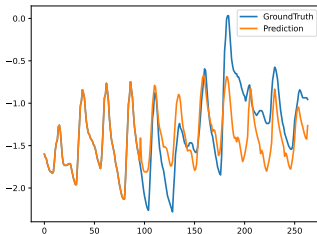
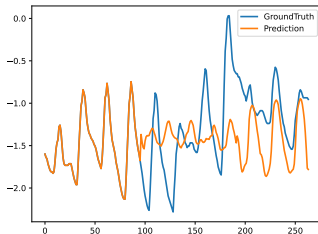
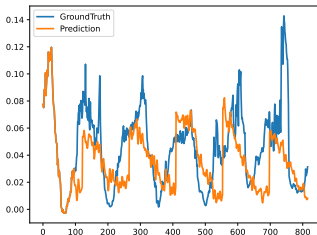
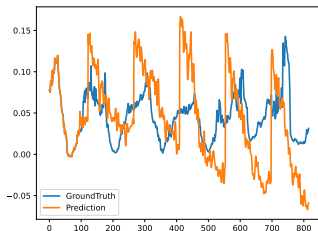


- ▶ Left: Change of  $L_f$  during training (ETTh<sub>2</sub>, 24h)
  - ▶ LipReg effectively altered input continuity
- ▶ Right: MSE with different  $\lambda$  (ETTh<sub>2</sub>)
  - ▶ Transformers prefer low continuity





Left: original; Right: with LipReg





- ▶ More Experiments:
  - ▶ Fine-tune Vision Transformers
  - ▶ Speech Transformers
  - ▶ Neural ODE



- ▶ In the Frequency Domain, the Lipschitz Regularizer is:

$$\begin{aligned}\sum_{i=0}^{n-1} (x_{i+1} - x_i)^2 &\approx \int_{\mathbb{R}} \left( \frac{df(t)}{dt} \right)^2 dt \\ &= \int_{\mathbb{R}} (2\pi i\xi)^2 \hat{f}^2(\xi) (-d\xi) \\ &= 4\pi^2 C \mathbb{E}_{p(\xi)}[\xi^2]\end{aligned}$$

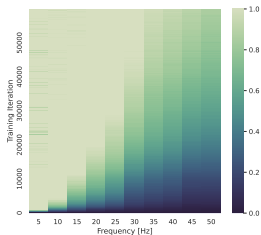
- ▶ An expectation over the frequency of the function
- ▶ Use it to penalize the frequency of the neural network



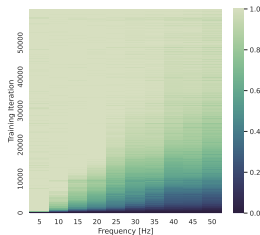
- ▶ Spectral Bias: low-frequency part is learned first
- ▶ Use LipReg to penalize the low-frequency part of NN

$$\mathcal{L}(y, \hat{y}) = \mathcal{L}_{\text{MSE}}(y, \hat{y}) - \lambda e^{-\epsilon t} \mathcal{L}_{\text{Lip}}(\hat{y})$$

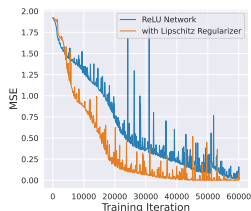
- ▶ Experiment: Curve fitting with ReLU Network
  - ▶ Training Iteration & Error in Frequency
  - ▶ Significantly reduce Spectral Bias



(a) Without LipReg



(b) With LipReg

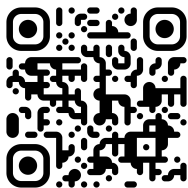


(c) Training Loss

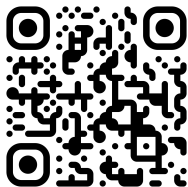


# Thank you for your attention!

Link to Paper



Link to Code



Poster No. 66