

# Efficiently Modeling Long Sequences with Structured State Spaces

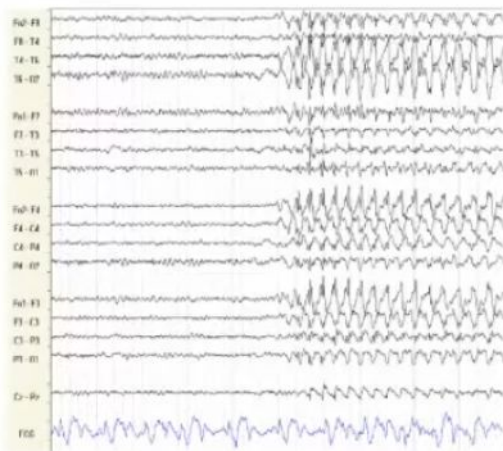
By Albert Gu, Karan Goel, and Christopher Ré

CPSC 670  
Spring 2023  
David Peng, Yujie Qiao

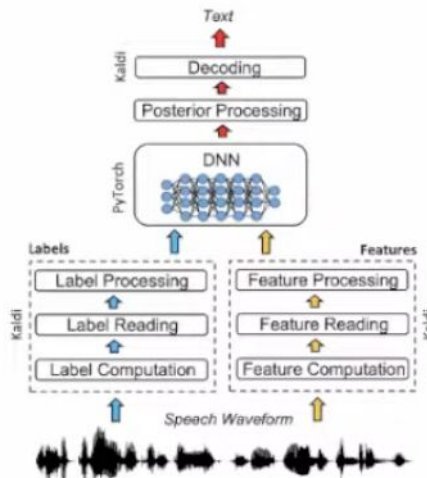
# Outline

- Introduction
- Background
- Methods
- Experiments
- Conclusions
- Discussion

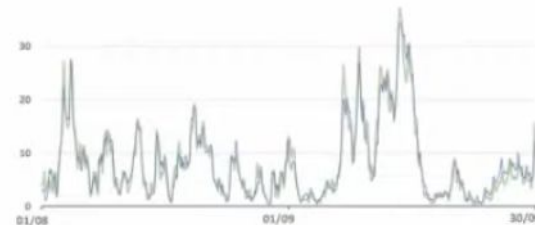
# Data/task: Long “Continuous” Time Series



EEG/ECG



Speech



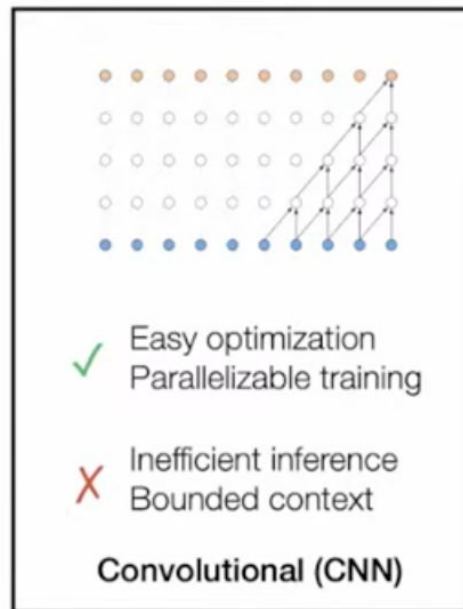
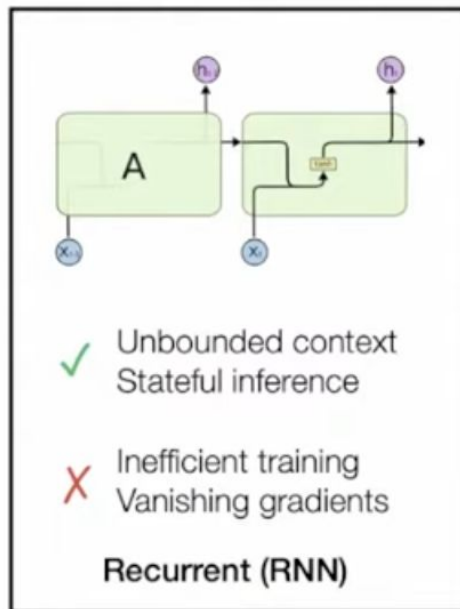
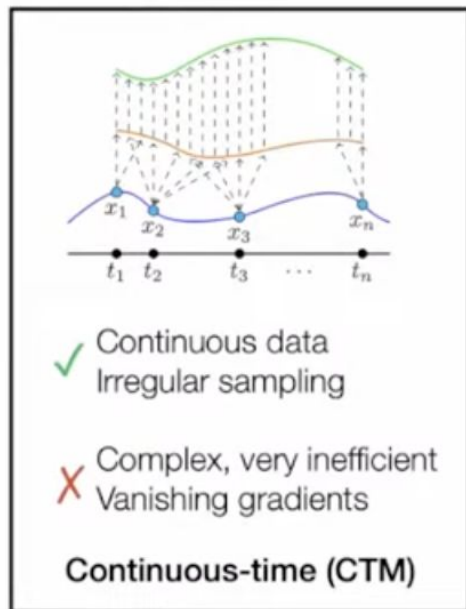
Energy Forecasting

# Data/task: Long “Continuous” Time Series

## Characteristics

- Sequence data that's sampled from underlying continuous signal
- Tens of thousands of samples
  - e.g. raw recordings of speech are 16,000 Hertz
  - 1-2 second recording → input sequence with tens of thousands of samples
- Often requires feature engineering

# Paradigms for Long Time Series



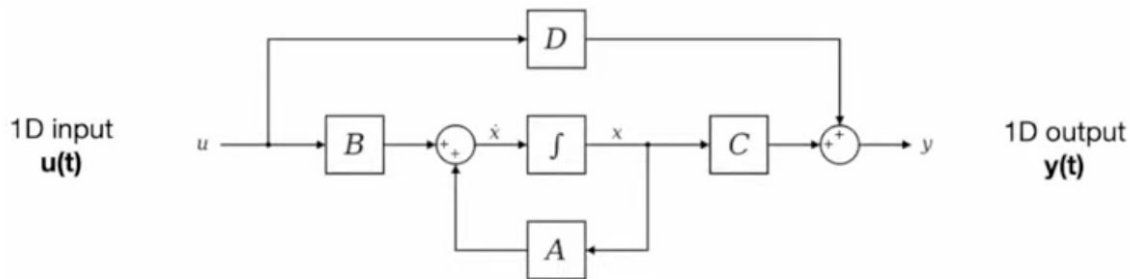
State space models (SSMs) can unify the strengths of CTM, RNN, CNN models at handling long-range dependencies

# Background: State Space Model (SSM)

- Linear time-invariant system
- Ordinary differential equation (ODE) used commonly in electrical engineering
- Given input signal (sequence)  $\mathbf{u}(t)$ , produce output signal (sequence)  $\mathbf{y}(t)$
- A larger dimensional hidden state  $\mathbf{x}(t)$  retains complex patterns
  - drop-in replacement for self-attention

$$\mathbf{x}'(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$$

$$\mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t)$$

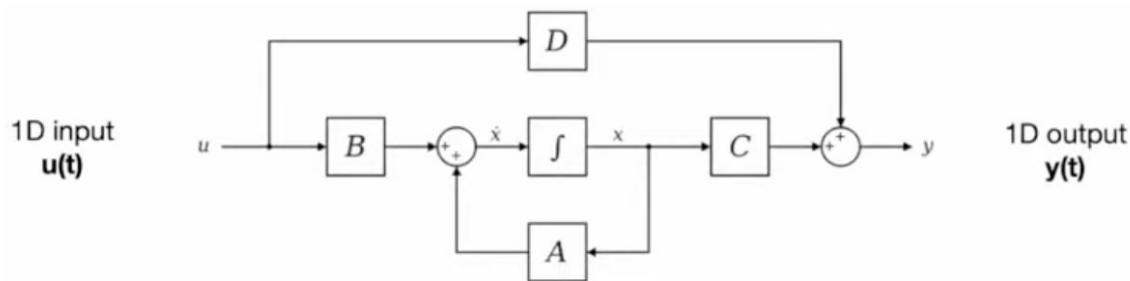


# Background: State Space Model (SSM)

- **A, B, C, D** are matrices
- Treat as parameters to be learned in a deep learning model

$$x'(t) = \mathbf{A}x(t) + \mathbf{B}u(t)$$

$$y(t) = \mathbf{C}x(t) + \mathbf{D}u(t)$$



$$u(t) \in \mathbb{R}^p \quad A \in \mathbb{R}^{d \times d}$$

$$x(t) \in \mathbb{R}^d \quad B \in \mathbb{R}^{d \times p}$$

$$y(t) \in \mathbb{R}^p \quad C \in \mathbb{R}^{1 \times d}$$

$$D \in \mathbb{R}^{d \times 1}$$



Solving the ODEs, we can derive...

$$x(s) = e^{As}x(0) + \int_0^s e^{A(s-t)}Bu(t) dt$$

- Given input data and a value of  $\mathbf{A}$ , we could numerically integrate to get the hidden state value at any time  $\mathbf{s}$
- Hidden state is a convolution over the input data
- Next step: turn integral into discrete approximation

# Turn continuous input into discrete input

- Use step size  $\Delta$
- $u[k] = u(k\Delta)$  for  $k = 1, 2, \dots, n$
- Convert state matrices into discrete approximations based on
  - step size  $\Delta$
  - state matrix  $A$
- Recurrence relation—like an RNN!

$$\begin{aligned}x_k &= \bar{A}x_{k-1} + \bar{B}u_k & \bar{A} &= (I - \Delta/2 \cdot A)^{-1}(I + \Delta/2 \cdot A) \\y_k &= \bar{C}x_k & \bar{B} &= (I - \Delta/2 \cdot A)^{-1}\Delta B & \bar{C} &= C\end{aligned}$$

- $D$  is set to identity matrix (acts like skip connection)

# Refine recurrence relation as convolution

- CNNs are easier to train than RNNs
- Convolution is easier to train than recurrence relation
- Calculate  $\overline{\mathbf{K}}$ , the SSM convolution kernel or filter

$$x_k = \overline{\mathbf{A}}x_{k-1} + \overline{\mathbf{B}}u_k$$

$$y_k = \overline{\mathbf{C}}x_k$$

$$\begin{array}{lll} x_0 = \overline{\mathbf{B}}u_0 & x_1 = \overline{\mathbf{A}}\overline{\mathbf{B}}u_0 + \overline{\mathbf{B}}u_1 & x_2 = \overline{\mathbf{A}}^2\overline{\mathbf{B}}u_0 + \overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{B}}u_2 \quad \dots \\ y_0 = \overline{\mathbf{C}}\overline{\mathbf{B}}u_0 & y_1 = \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{B}}u_1 & y_2 = \overline{\mathbf{C}}\overline{\mathbf{A}}^2\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_1 + \overline{\mathbf{C}}\overline{\mathbf{B}}u_2 \quad \dots \end{array}$$

$$y_k = \overline{\mathbf{C}}\overline{\mathbf{A}}^k\overline{\mathbf{B}}u_0 + \overline{\mathbf{C}}\overline{\mathbf{A}}^{k-1}\overline{\mathbf{B}}u_1 + \dots + \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}u_{k-1} + \overline{\mathbf{C}}\overline{\mathbf{B}}u_k$$

$$y = \overline{\mathbf{K}} * u.$$

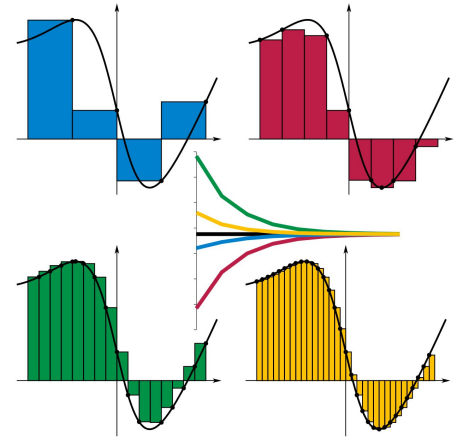
$$\overline{\mathbf{K}} \in \mathbb{R}^L := \mathcal{K}_L(\overline{\mathbf{A}}, \overline{\mathbf{B}}, \overline{\mathbf{C}}) := \left( \overline{\mathbf{C}}\overline{\mathbf{A}}^i\overline{\mathbf{B}} \right)_{i \in [L]} = (\overline{\mathbf{C}}\overline{\mathbf{B}}, \overline{\mathbf{C}}\overline{\mathbf{A}}\overline{\mathbf{B}}, \dots, \overline{\mathbf{C}}\overline{\mathbf{A}}^{L-1}\overline{\mathbf{B}})$$

# Turning continuous function into discrete function

- Input of the model is actually discrete inputs
- Sample frequency  $T$
- $u[k] = u(kT)$  for  $k = 1, 2, \dots, n$
- How can we find  $x[k] = x(kT)$  ?
- We can derive...

$$x[n + 1] = e^{AT} x[n] + Te^{AT} Bu[n]$$

$$\int_0^{(n+1)T} f(u(t))dt \approx T \sum_{k=0}^n f(u[k])$$



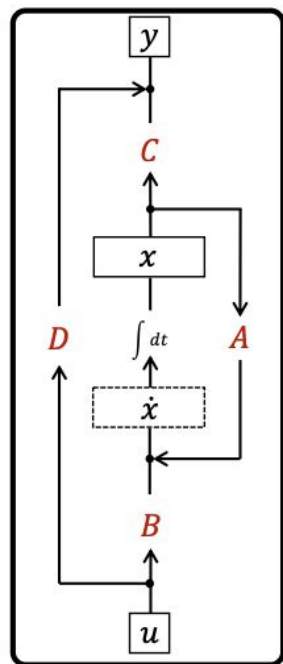
## Redefine recurrence relation as convolution

$$\mathbf{x}[n + 1] = e^{AT} \mathbf{x}[n] + T e^{AT} B u[n]$$

- Convolution is easier to calculate than recurrence
- We can derive...

$$\mathbf{x}[n + 1] = e^{AT(n+1)} \mathbf{x}[0] + T \sum_{j=0}^k e^{A(k-j)T} B u[j]$$

# Three forms of SSM

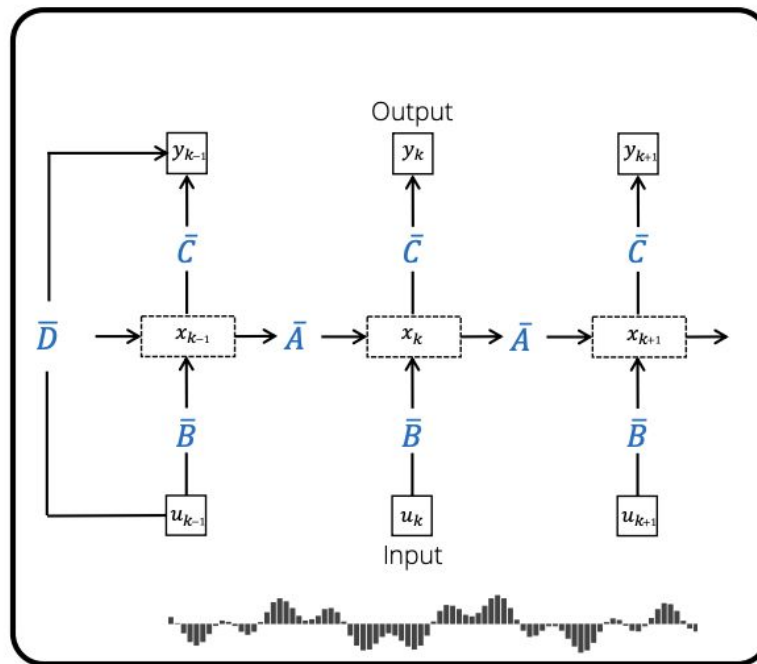


Continuous-time

- ✓ continuous data
- ✓ irregular sampling

Discretize

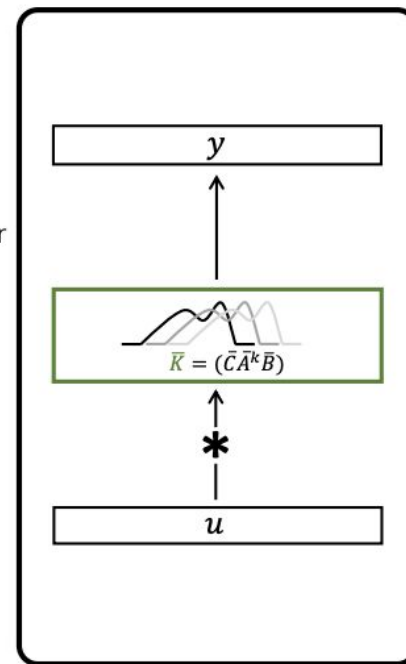
$\Delta t$



Recurrent

- ✓ unbounded context
- ✓ efficient inference

or



Convolutional

- ✓ local information
- ✓ parallelizable training

# Applying SSM to deep learning

- SSM layer acts as a drop-in replacement for attention
- Model learns matrices ***A***, ***B***, ***C***, ***D***
- We want the training model to be
  - Stable: training leads to reasonable output
  - Efficient: minimize number of computations
  - Expressive: can learn non-linear, complex functions

$$x'(t) = \mathbf{A}x(t) + \mathbf{B}u(t)$$

$$y(t) = \mathbf{C}x(t) + \mathbf{D}u(t)$$

# Previous work with SSM: Linear State Space Layer (LSSL)

- Paper published a year earlier by Gu et al.
- Proof of concept for SSM at handling long-range dependencies
- Problem: computation is slow and memory inefficient
- $O(N^2 L)$  operations for hidden state dimension  $N$ , sequence length  $L$ 
  - Much more than comparably-sized RNN, CNN

$$y_k = \overline{CA}^k \overline{B}u_0 + \overline{CA}^{k-1} \overline{B}u_1 + \cdots + \overline{CAB}u_{k-1} + \overline{CB}u_k$$
$$y = \overline{K} * u.$$

$$\overline{K} \in \mathbb{R}^L := \mathcal{K}_L(\overline{A}, \overline{B}, \overline{C}) := \left( \overline{CA}^i \overline{B} \right)_{i \in [L]} = (\overline{CB}, \overline{CAB}, \dots, \overline{CA}^{L-1} \overline{B})$$

- How can we calculate  $\overline{A}^L$  efficiently? Bottleneck problem



# LSSL: non-random initialization of matrix $\mathbf{A}$

- Random initialization leads to poor performance
- Use HiPPO Matrix
  - Paper published earlier by Gu et al.
  - Allows  $\mathbf{x}(\mathbf{t})$  to memorize history of input  $\mathbf{u}(\mathbf{t})$
  - “compresses” past history to be able to reconstruct most of it
  - Only needs to be calculated once

$$\text{(HiPPO Matrix)} \quad \mathbf{A}_{nk} = - \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$

# New work: Structured State Space sequence model (S4)

- Further exploits HiPPO matrices to efficiently compute  $\overline{\mathbf{K}}$
- Main contribution: some matrix math tricks...
  - Argue  $\mathbf{A}$  is a type of matrix called a **Normal Plus Low-Rank (NPLR)**
  - Rewrite  $\overline{\mathbf{K}}$  in the form of its truncated generating function
  - Generating function uses matrix inverses instead of matrix powers. Apply an identity and reduce equation to a simpler case
  - Show that reduced case is equivalent to **Cauchy matrices**, which have been shown to be fast and stable

**Theorem 3** (S4 Convolution). *Given any step size  $\Delta$ , computing the SSM convolution filter  $\overline{\mathbf{K}}$  can be reduced to 4 Cauchy multiplies, requiring only  $\tilde{O}(N + L)$  operations and  $O(N + L)$  space.*

We reiterate that Theorem 3 is our core technical contribution

# Experiments - Efficiency Benchmarks

S4 is theoretically much more efficient than the LSSL

- 30x faster, 400x less memory usage

Table 2: Deep SSMS: The S4 parameterization with Algorithm 1 is asymptotically more efficient than the LSSL.

Dim.	TRAINING STEP (MS)			MEMORY ALLOC. (MB)		
	128	256	512	128	256	512
LSSL	9.32	20.6	140.7	222.1	1685	13140
<b>S4</b>	<b>4.77</b>	<b>3.07</b>	<b>4.75</b>	<b>5.3</b>	<b>12.6</b>	<b>33.5</b>
Ratio	1.9×	6.7×	<b>29.6×</b>	42.0×	133×	<b>392×</b>

Table 3: Benchmarks vs. efficient Transformers

	LENGTH 1024		LENGTH 4096	
	Speed	Mem.	Speed	Mem.
Transformer	1×	1×	1×	1×
Performer	1.23×	<u>0.43</u> ×	3.79×	<u>0.086</u> ×
Linear Trans.	<b>1.58</b> ×	<b>0.37</b> ×	<b>5.35</b> ×	<b>0.067</b> ×
<b>S4</b>	<b>1.58</b> ×	<u>0.43</u> ×	<u>5.19</u> ×	0.091×

# Experiments - Long Range Arena (LRA) and raw speech classification

## LRA

- Contains 6 tasks with lengths 1K-16K steps
- S4 solves the Path-X task that involves reasoning about LRDs over sequences of length  $128 \times 128 = 163840$

Table 4: (**Long Range Arena**) (*Top*) Original Transformer variants in LRA. Full results in Appendix [D.2](#). (*Bottom*) Other models reported in the literature. *Please read Appendix [D.5](#) before citing this table.*

MODEL	LISTOPS	TEXT	RETRIEVAL	IMAGE	PATHFINDER	PATH-X	AVG
Transformer	36.37	64.27	57.46	42.44	71.40	✗	53.66
Reformer	<u>37.27</u>	56.10	53.40	38.07	68.50	✗	50.56
BigBird	36.05	64.02	59.29	40.83	74.87	✗	54.17
Linear Trans.	16.13	<u>65.90</u>	53.09	42.34	75.30	✗	50.46
Performer	18.01	65.40	53.82	42.77	77.05	✗	51.18
FNet	35.33	65.11	59.61	38.67	<u>77.80</u>	✗	54.42
Nyströmformer	37.15	65.52	<u>79.56</u>	41.58	70.94	✗	57.46
Luna-256	37.25	64.57	79.29	<u>47.38</u>	77.72	✗	<u>59.37</u>
<b>S4</b>	<b>59.60</b>	<b>86.82</b>	<b>90.90</b>	<b>88.65</b>	<b>94.20</b>	<b>96.35</b>	<b>86.09</b>

# Experiments - Long Range Arena (LRA) and raw speech classification

## Raw Speech Classification

- Used the SC10 subset of the Speech Commands dataset
- Classified raw speech (length-16000)
- S4 achieves 98.3% accuracy, higher than all baselines that use the 100× shorter MFCC features

# Experiments - First-ever non-random performance on Path-X task

- Pathfinder task but with longer sequence lengths
- 88% accuracy (50% random guessing for all prior models)

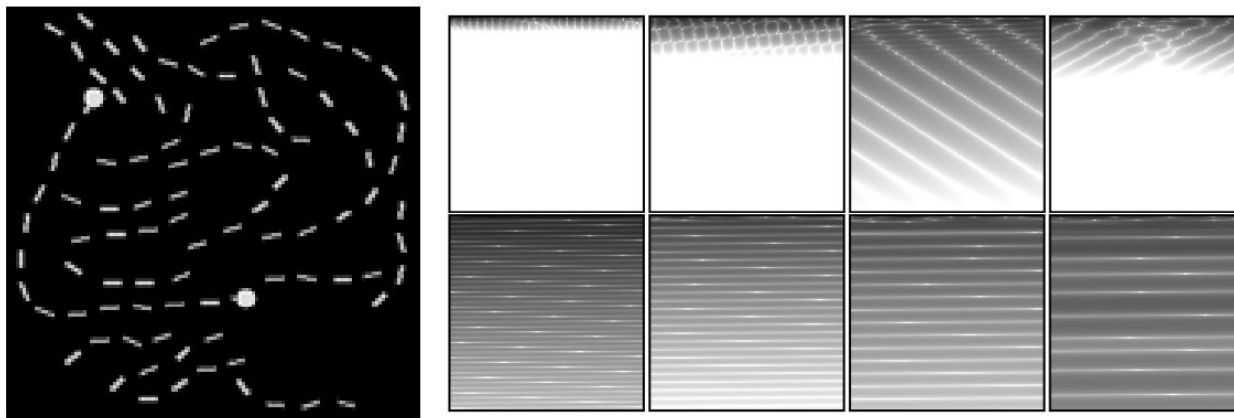


Figure 2: Visualizations of a trained S4 model on LRA Path-X. SSM convolution kernels  $\overline{\mathbf{K}} \in \mathbb{R}^{16384}$  are reshaped into a  $128 \times 128$  image. (*Left*) Example from the Path-X task, which involves deducing if the markers are connected by a path (*Top*) Filters from the first layer (*Bottom*) Filters from the last layer.

# Experiments - General Sequence Model

Key goal of sequence modeling research is to develop a single model that can be applied in many domains (e.g. images, audio, text, time-series)

S4 can be used as a general sequence model to perform effectively and efficiently in a wide variety of settings including image classification, image and text generation, and time series forecasting

Table 5: **(SC10 classification)** Transformer, CTM, RNN, CNN, and SSM models. (MFCC) Standard pre-processed MFCC features (length 161). (Raw) Unprocessed signals (length 16000). ( $0.5\times$ ) Frequency change at test time. ✗ denotes not applicable or computationally infeasible on single GPU. Please read Appendix D.5 before citing this table.

	MFCC	RAW	$0.5\times$
Transformer	90.75	✗	✗
Performer	80.85	30.77	30.68
ODE-RNN	65.9	✗	✗
NRDE	89.8	16.49	15.12
ExpRNN	82.13	11.6	10.8
LipschitzRNN	88.38	✗	✗
CKConv	<b>95.3</b>	71.66	<u>65.96</u>
WaveGAN-D	✗	<u>96.25</u>	✗
LSSL	93.58	✗	✗
<b>S4</b>	<u>93.96</u>	<b>98.32</b>	<b>96.30</b>

Table 6: **(Pixel-level 1-D image classification)** Comparison against reported test accuracies from prior works (Transformer, RNN, CNN, and SSM models). Extended results and citations in Appendix D.

	sMNIST	PMNIST	sCIFAR
Transformer	98.9	97.9	62.2
LSTM	98.9	95.11	63.01
r-LSTM	98.4	95.2	72.2
UR-LSTM	99.28	96.96	71.00
UR-GRU	99.27	96.51	74.4
HiPPO-RNN	98.9	98.3	61.1
LMU-FFT	-	98.49	-
LipschitzRNN	99.4	96.3	64.2
TCN	99.0	97.2	-
TrellisNet	99.20	98.13	73.42
CKConv	99.32	98.54	63.74
LSSL	<u>99.53</u>	<b>98.76</b>	<u>84.65</u>
<b>S4</b>	<b>99.63</b>	<u>98.70</u>	<b>91.13</b>

# Experiments - General Sequence Model

Table 7: (CIFAR-10 density estimation) As a generic sequence model, S4 is competitive with previous autoregressive models (in bits per dim.) while incorporating no 2D inductive bias, and has fast generation through its recurrence mode.

Model	bpd	2D bias	Images / sec
Transformer	3.47	<b>None</b>	0.32 (1×)
Linear Transf.	3.40	<b>None</b>	17.85 (56×)
PixelCNN	3.14	2D conv.	-
Row PixelRNN	3.00	2D BiLSTM	-
PixelCNN++	2.92	2D conv.	<u>19.19</u> (59.97×)
Image Transf.	2.90	2D local attn.	0.54 (1.7×)
PixelSNAIL	<u>2.85</u>	2D conv. + attn.	0.13 (0.4×)
Sparse Transf.	<b>2.80</b>	2D sparse attn.	-
<b>S4 (base)</b>	2.92	<b>None</b>	<b>20.84 (65.1×)</b>
<b>S4 (large)</b>	<u>2.85</u>	<b>None</b>	3.36 (10.5×)

Table 8: (WikiText-103 language modeling) S4 ap-sequence model, S4 is competitive with previous autoregressive models the performance of Transformers with much faster generation. (Top) Transformer baseline which our implementation is based on, with attention replaced by S4. (Bottom) Attention-free models (RNNs and CNNs).

Model	Params	Test ppl.	Tokens / sec
Transformer	247M	<b>20.51</b>	0.8K (1×)
GLU CNN	229M	37.2	-
AWD-QRNN	151M	33.0	-
LSTM + Hebb.	-	29.2	-
TrellisNet	180M	29.19	-
Dynamic Conv.	255M	25.0	-
TaLK Conv.	240M	23.3	-
<b>S4</b>	249M	<b>20.95</b>	<b>48K (60×)</b>



# SSM Ablations: the Importance of HiPPO

The HiPPO matrices were critical to initialize an SSM for S4

Investigation on generic SSMs with various initializations

A random Gaussian initialization (with variance scaled down until it did not NaN) & the HiPPO initialization

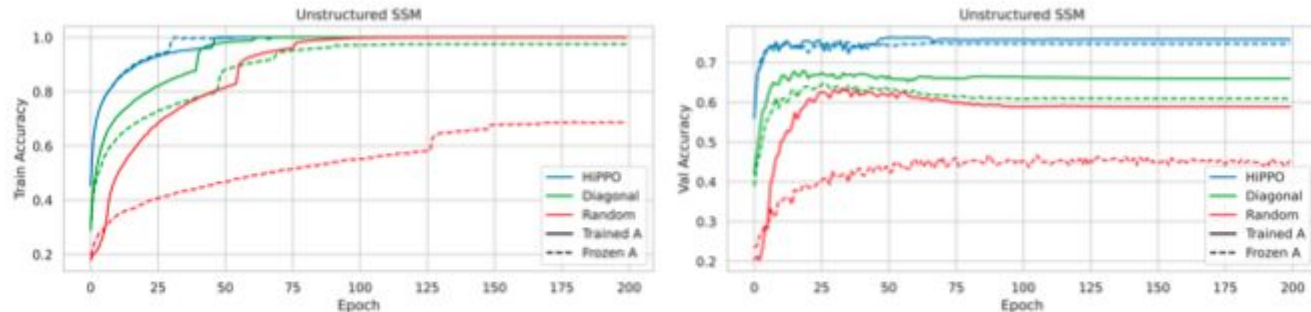


Figure 3: CIFAR-10 classification with unconstrained, real-valued SSMs with various initializations. (Left) Train accuracy. (Right) Validation accuracy.

# SSM Ablations: the Importance of HiPPO

S4's effectiveness primarily comes from the HiPPO initialization, not the NPLR parameterization

The HiPPO initialization (the full S4 method) achieves 84.27% test accuracy with just 100K parameters

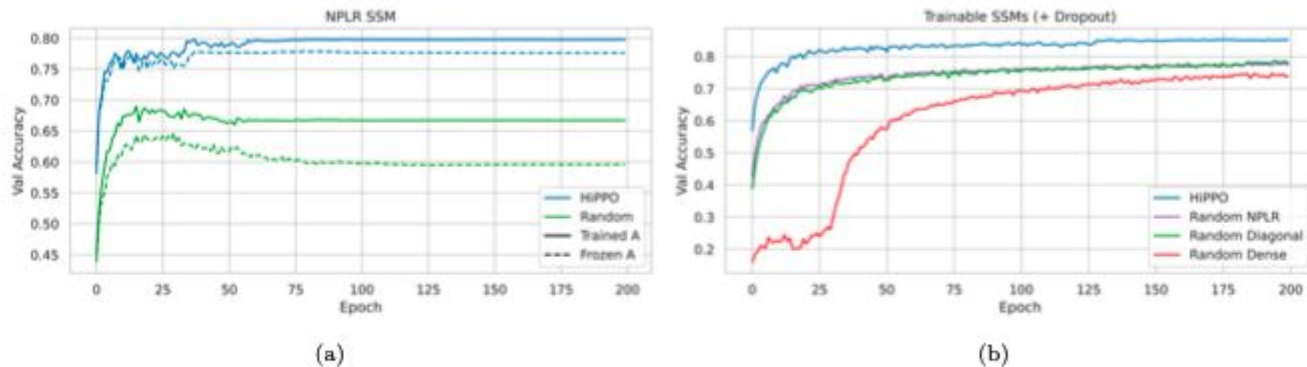


Figure 4: CIFAR-10 validation accuracy of SSMs with different initializations and parameterizations. (Left) NPLR parameterization with random versus HiPPO initialization. (Right) All methods considered in this section, including minor Dropout regularization. S4 achieves SotA accuracy on sequential CIFAR-10 with just 100K parameters.

# Conclusion

## Structured State Space sequence model (S4)

- A sequence model that uses a new parameterization for the SSM's continuous- time, recurrent, and convolutional views to efficiently model LRDs in a principled manner
- S4 achieves strong empirical results across a diverse range of established benchmarks, which suggest that S4 has the potential to be an effective general sequence modeling solution
  - 91% accuracy on sequential CIFAR-10 with no data augmentation or auxiliary losses
  - Comparable to Transformers on image and language modeling tasks, while performing generation 60× faster
  - SoTA on every task from the Long Range Arena benchmark

# Discussion

- What are some limitations of S4?
- Does the HiPPO matrices need to be stable?
- How do you think S4 will perform in traditional NLP classification tasks?
- Is there any other ways for memorizing the sequences other than the state space model (SSM)? How would one extend the work in this paper to other directions?