



10-423/10-623 Generative AI

Machine Learning Department
School of Computer Science
Carnegie Mellon University

State Space Models

Matt Gormley & Pat Virtue

Lecture 22

Apr. 9, 2025

Reminders

- **Project Midway Report**
 - Due: Thu, Apr-17 at 11:59pm
- **HW623**
 - Only for students registered in 10-623
 - Due: Mon, Apr-21 at 11:59pm

Motivation

- Transformers are slow at test time: they require a KV cache that grows ~~quadratically~~ ^{linearly} in size with the sequence length
- State space models (SSMs) are fast at test time: they only hold a fixed size hidden state in memory (like RNNs)
- But we'll see that SSMs can also be trained efficiently with the right tricks
- As well, they elegantly transition between different granularities of representation for the input (e.g. sound at 16Hz vs. 8Hz)

Motivation

- <https://www.isattentionallyouneed.com/>

Is Attention All You Need?



Current Status: Yes

Time Remaining: 631d 22h 55m 11s

Proposition:

On January 1, 2027, a Transformer-like model will continue to hold the state-of-the-art position in most benchmarked tasks in natural language processing.

Motivation

- <https://www.isattentionallyyouneed.com/>

Is Attention All You Need?

For the Motion

Jonathan Frankle
@jefrankle
Harvard Professor
Chief Scientist Mosaic ML



Against the Motion

Sasha Rush
@srush_nlp
Cornell Professor
Research Scientist Hugging Face 🤗



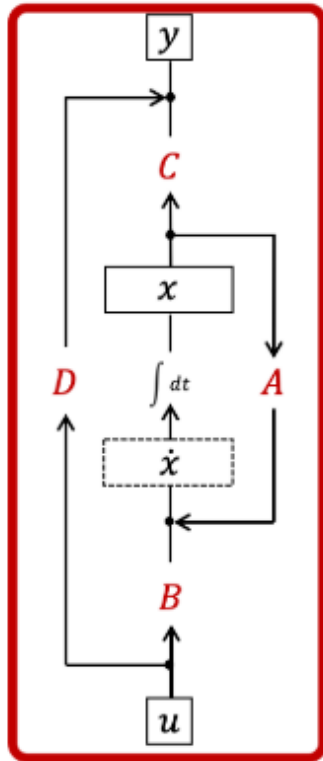
Wager

The wager is for donation of equity in Mosaic ML or Hugging Face to a charity of the winner's choice. Details to come.

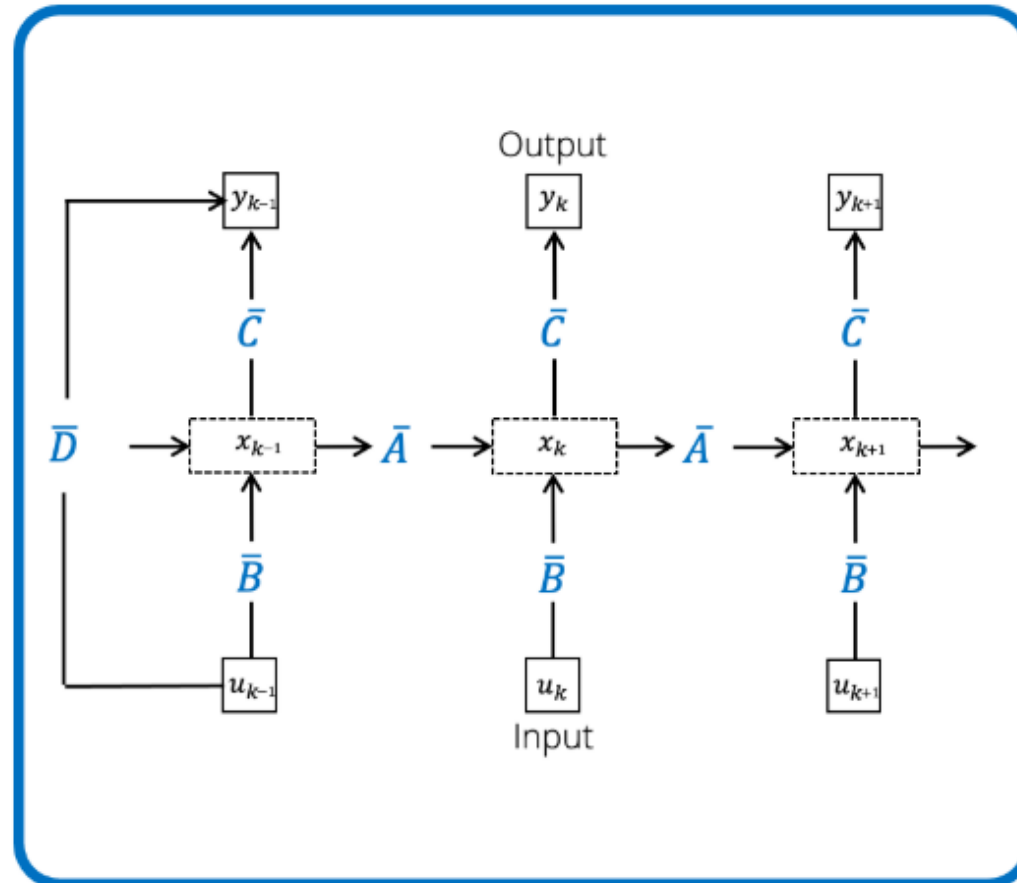
STATE SPACE MODEL (SSM)

Three Representations of a State Space Model

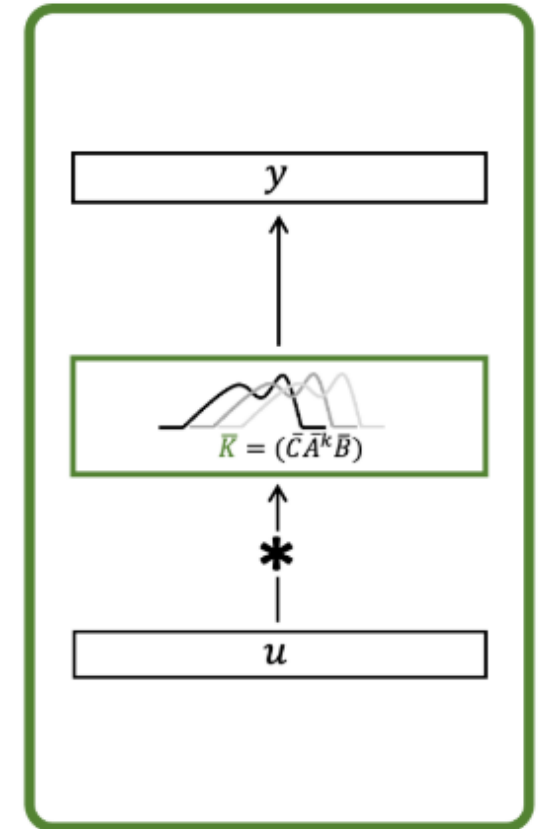
Continuous



Recurrent



Convolutional

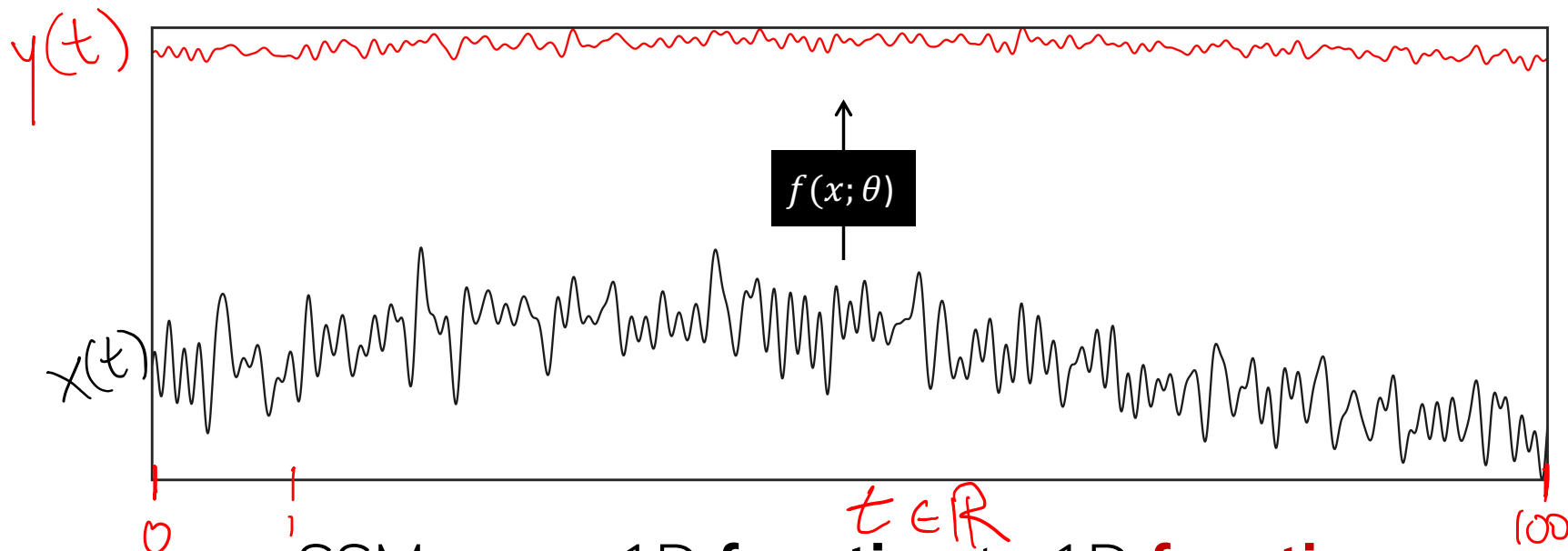


SSM: 1D Continuous Representation

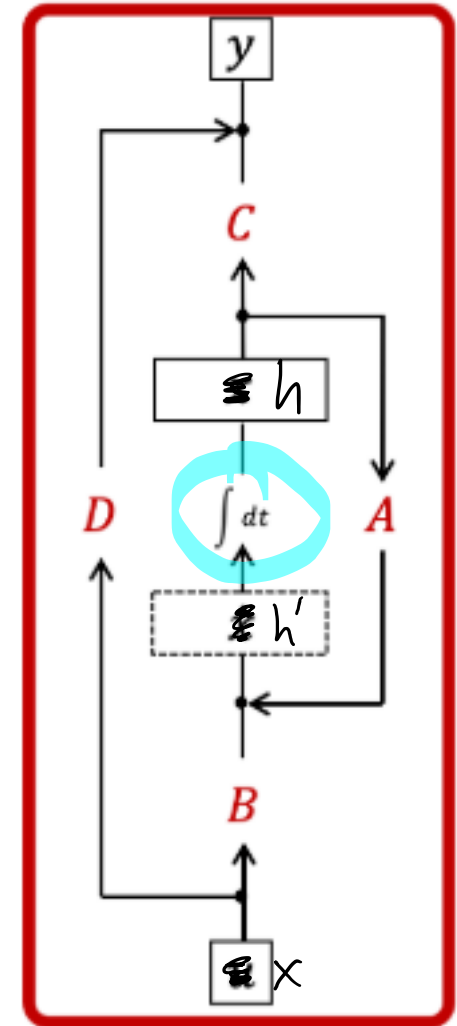
$x(t) \in \mathbb{R}$
 $y(t) \in \mathbb{R}$
 $h(t) \in \mathbb{R}^N$
 $A \in \mathbb{R}^{N \times N}$
 $B \in \mathbb{R}^{N \times 1}$

$C \in \mathbb{R}^{1 \times N}$
 $D \in \mathbb{R}^{1 \times 1}$

$$\begin{aligned}
 h'(t) &= \mathbf{A}h(t) + \mathbf{B}x(t) = \frac{d}{dt} h(t) \\
 y(t) &= \mathbf{C}h(t) + \mathbf{D}x(t)
 \end{aligned}$$



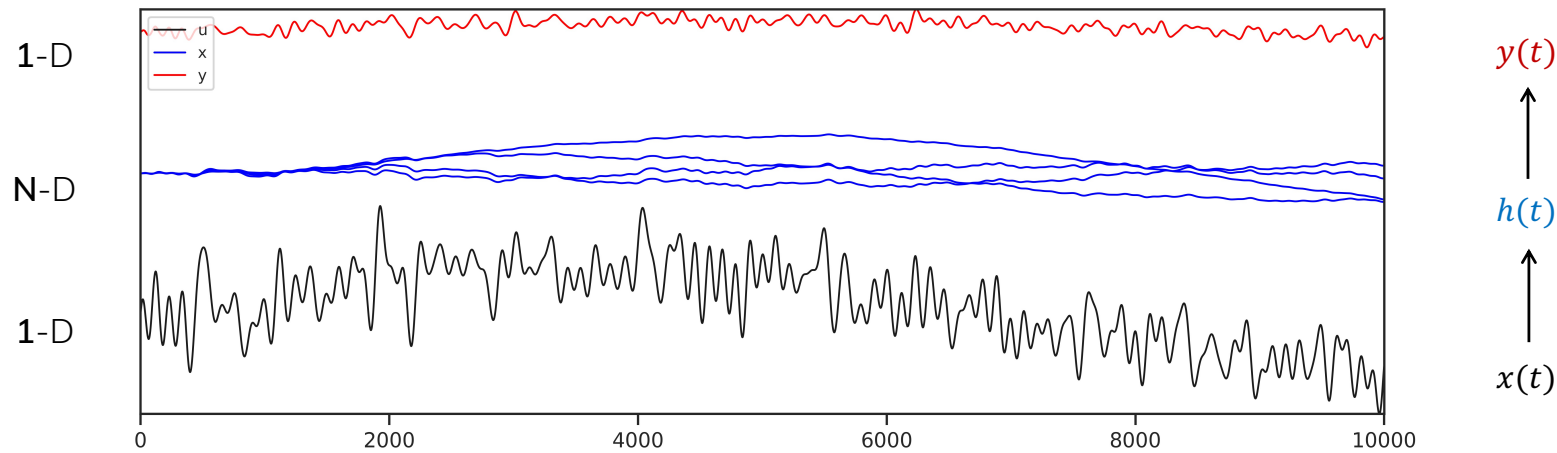
SSMs map 1D **function** to 1D **function**



SSM: 1D Continuous Representation

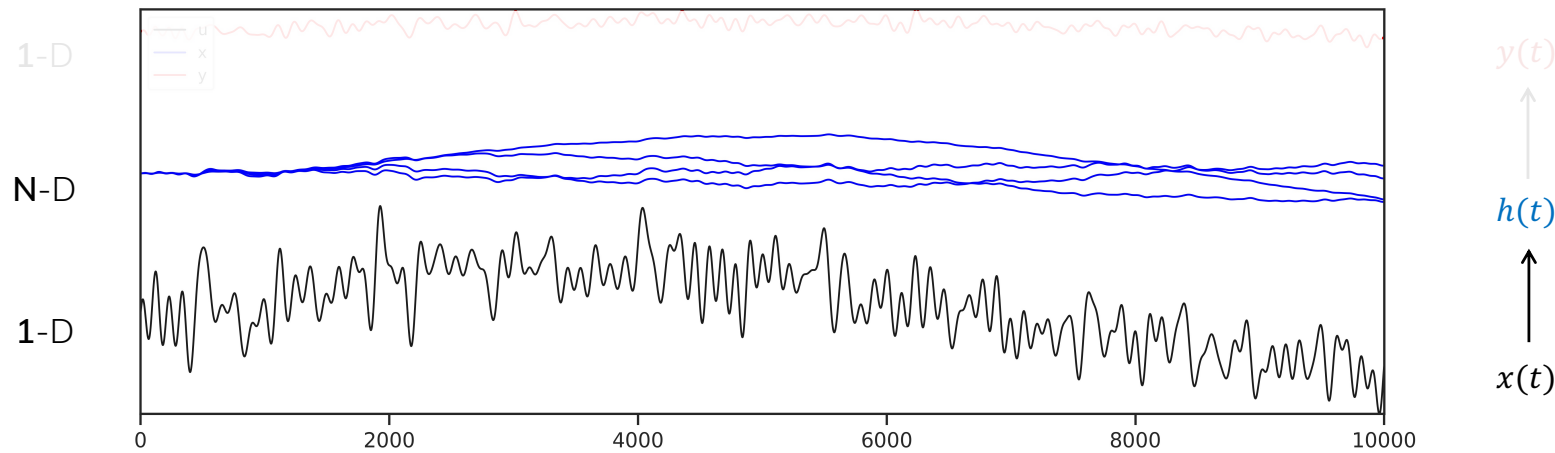
$$h'(t) = \mathbf{A}h(t) + \mathbf{B}x(t)$$

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



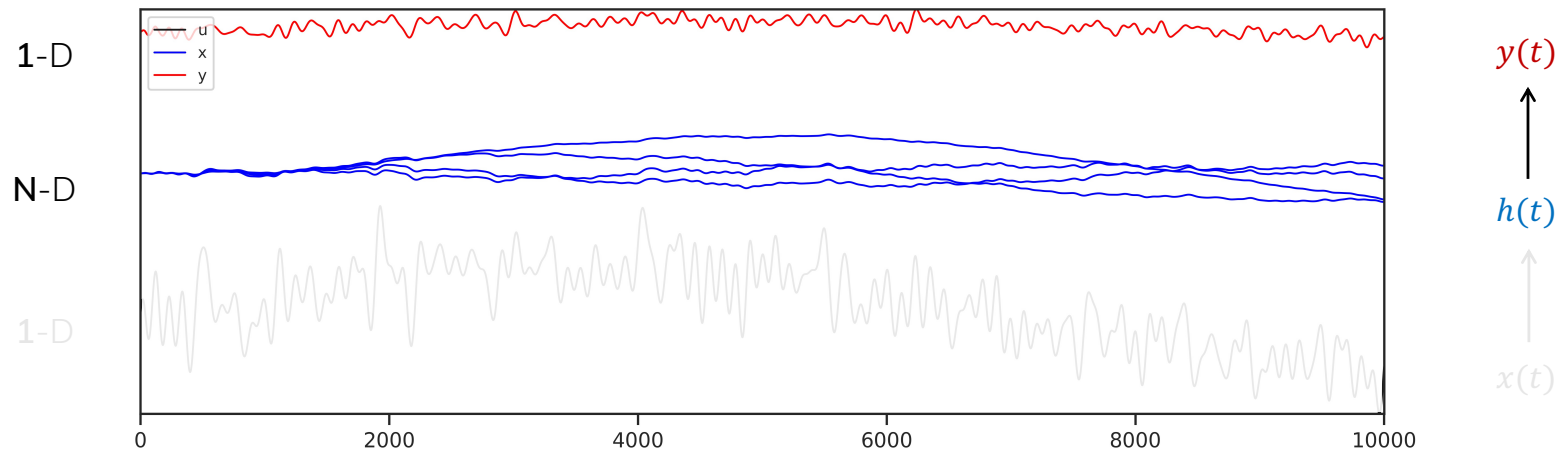
SSM: 1D Continuous Representation

$$h'(t) = \mathbf{A}h(t) + \mathbf{B}x(t)$$
$$y(t) = \mathbf{C}h(t) + \mathbf{D}x(t)$$



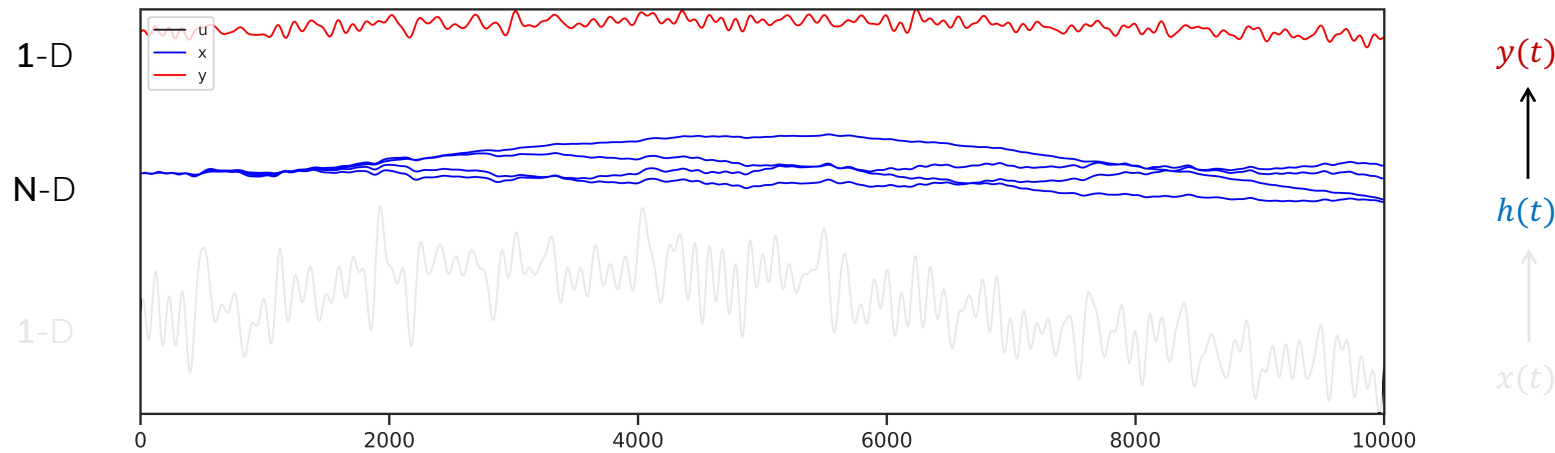
SSM: 1D Continuous Representation

$$h'(t) = Ah(t) + Bx(t)$$
$$y(t) = Ch(t) + Dx(t)$$



SSM: 1D Continuous Representation

$$h'(t) = Ah(t) + Bx(t)$$
$$y(t) = Ch(t) + Dx(t)$$



SSM: 1D Continuous Representation

$$h(t) = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$h_1(t) = 0$$

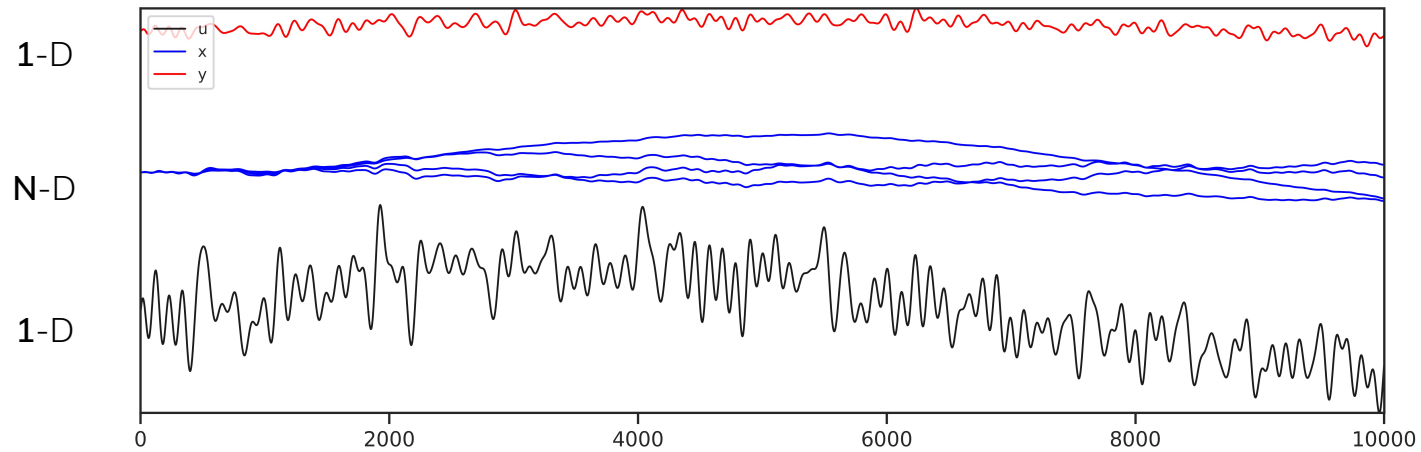
$$h_2(t) = 0$$

$$\vdots$$

$$h_N(t) = 0$$

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

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



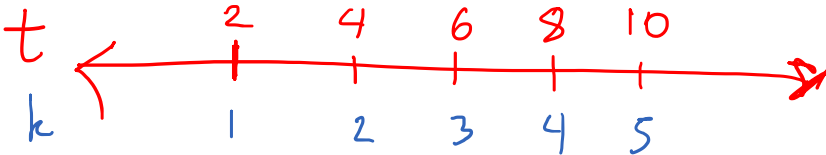
$y(t)$
↑
 $h(t)$
↑
 $x(t)$



SSM: 1D Discrete Recurrent Representation

Continuous Representation

- Uses parameters we will actually work with in the end
- Seamlessly represents any continuous 1D \rightarrow 1D function
- Impractical for real data



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

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

$$t \in \mathbb{R}^+$$

Discrete Recurrent Representation

- A discrete approximation using different parameters which are **functions** of the original parameters A,B,C,D
- Allows us to work with real data

$$h_{k+1} = \overline{\mathbf{A}}h_k + \overline{\mathbf{B}}x_k$$

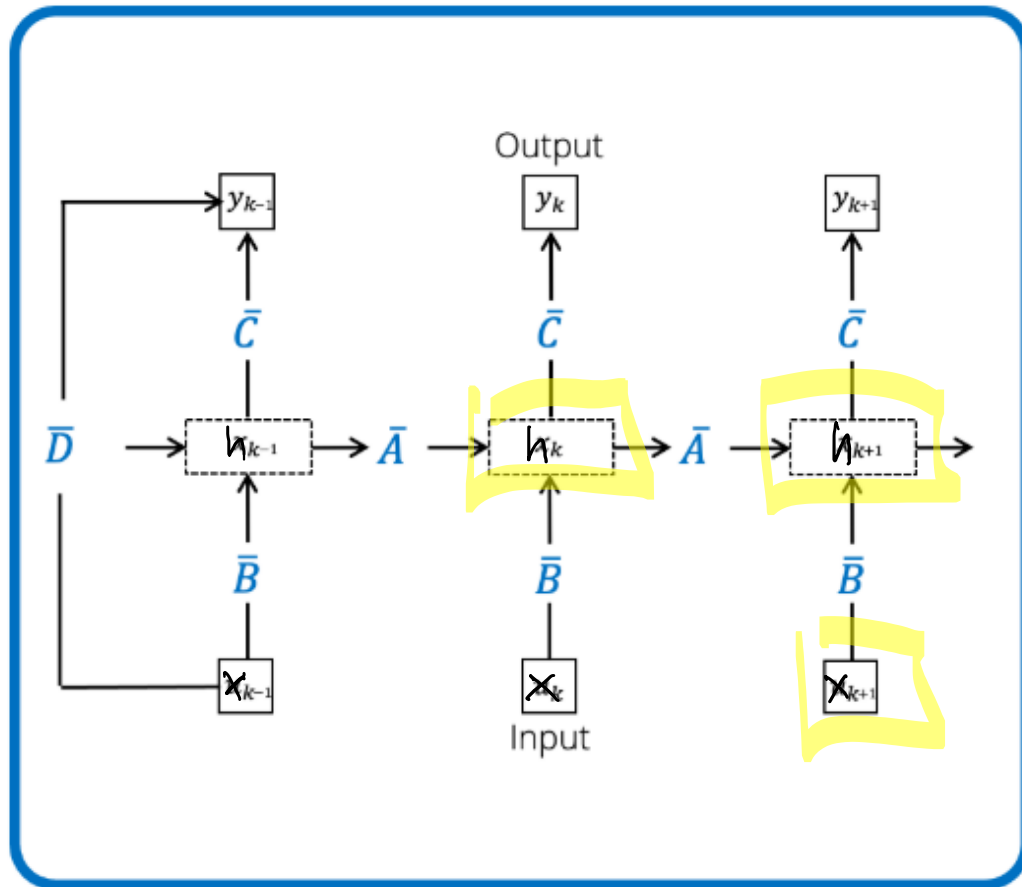
$$y_k = \overline{\mathbf{C}}h_k + \overline{\mathbf{D}}x_k$$

$$k \in \mathbb{N}^+$$

SSM: 1D Discrete Recurrent Representation

Q: How is this different from an RNN?

Question: How can we depict this recurrent computation?



Discrete Recurrent Representation

- A discrete approximation using different parameters which are **functions** of the original parameters A, B, C, D
- Allows us to work with real data

$$h_{k+1} = \bar{A}h_k + \bar{B}x_k$$

$$y_k = \bar{C}h_k + \bar{D}x_k$$

residual connection

no bias term

no activation function

are 1D

How to discretize a continuous SSM?

S4 uses a bilinear transformation to discretize the continuous SSM

$$\bar{\mathbf{A}} = e^{\Delta A}$$

$$\bar{\mathbf{B}} = \mathbf{A}^{-1}(e^{\Delta A} - \mathbf{I})\mathbf{B}$$

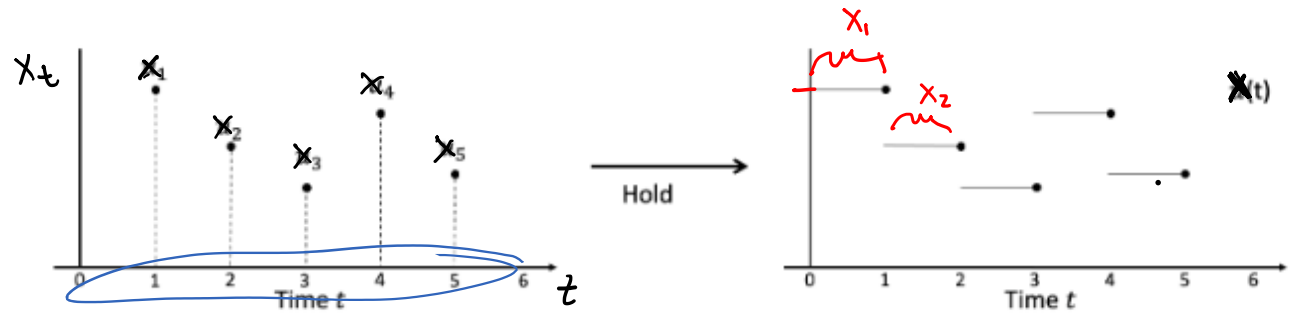
$$\bar{\mathbf{C}} = \mathbf{C}$$

$$\bar{\mathbf{D}} = \mathbf{D}$$

$\Delta = \text{step size}$

The bilinear transformation uses a first order Pade approximation:

$$e^x \approx \frac{1+x/2}{1-x/2}$$

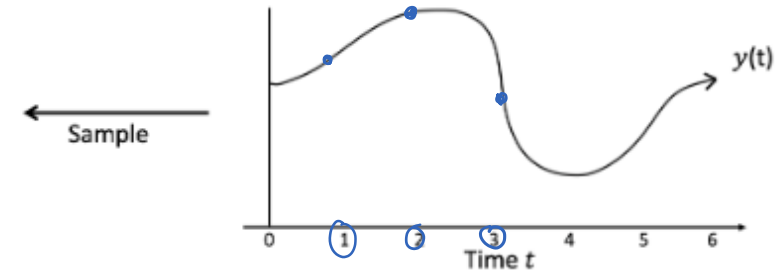
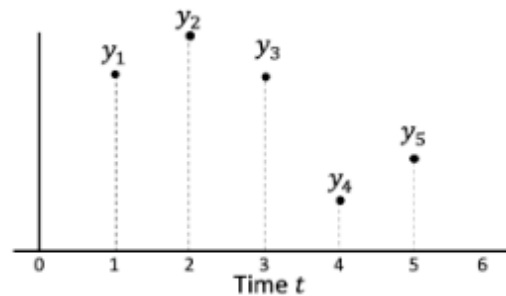


$$\begin{aligned} h_{k+1} &= \bar{\mathbf{A}}h_k + \bar{\mathbf{B}}x_k \\ y_k &= \bar{\mathbf{C}}h_k + \bar{\mathbf{D}}x_k \end{aligned}$$

Discrete SSM

$$\begin{aligned} h'(t) &= \mathbf{A}h(t) + \mathbf{B}x(t) \\ y(t) &= \mathbf{C}h(t) + \mathbf{D}x(t) \end{aligned}$$

Continuous SSM



The discrete-time SSM (left), a sequence-to-sequence map, is exactly equivalent to applying the bilinear transformation to the continuous-time SSM (right), a function-to-function map, on the held signal.

How to discretize a continuous SSM?

S4 uses a bilinear transformation to discretize the continuous SSM

$$\bar{\mathbf{A}} = e^{\Delta \mathbf{A}}$$

$$\bar{\mathbf{B}} = \mathbf{A}^{-1}(e^{\Delta \mathbf{A}} - \mathbf{I})\mathbf{B}$$

$$\bar{\mathbf{C}} = \mathbf{C}$$

$$\bar{\mathbf{D}} = \mathbf{D}$$

The bilinear transformation uses a first order Pade approximation:

$$e^x \approx \frac{1+x/2}{1-x/2}$$

$$\bar{\mathbf{A}} = \left(\mathbf{I} - \frac{\Delta}{2} \cdot \mathbf{A}\right)^{-1} \left(\mathbf{I} + \frac{\Delta}{2} \cdot \mathbf{A}\right)$$

$$\bar{\mathbf{B}} = \left(\mathbf{I} - \frac{\Delta}{2} \cdot \mathbf{A}\right)^{-1} \Delta \mathbf{B}$$

$$\bar{\mathbf{C}} = \mathbf{C}$$

$$\bar{\mathbf{D}} = \mathbf{D}$$

think of this as
a PyTorch layer.
b/c we will learn $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$
and not $\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{C}}, \bar{\mathbf{D}}$

SSM: 1D Convolutional Representation

We unroll the recurrent computation as:
Assume a zero initial state: $h_{-1} = 0$.

$$\begin{aligned}
 h_0 &= \bar{\mathbf{B}}x_0 \\
 \rightarrow h_1 &= \bar{\mathbf{A}}\bar{\mathbf{B}}x_0 + \bar{\mathbf{B}}x_1 \\
 \rightarrow h_2 &= \bar{\mathbf{A}}^2\bar{\mathbf{B}}x_0 + \bar{\mathbf{A}}\bar{\mathbf{B}}x_1 + \bar{\mathbf{B}}x_2 \\
 &\vdots \bar{\mathbf{A}} h_1 \\
 y_0 &= \bar{\mathbf{C}}\bar{\mathbf{B}}x_0 \\
 \rightarrow y_1 &= \bar{\mathbf{C}}\bar{\mathbf{A}}\bar{\mathbf{B}}x_0 + \bar{\mathbf{C}}\bar{\mathbf{B}}x_1 \\
 \rightarrow y_2 &= \bar{\mathbf{C}}\bar{\mathbf{A}}^2\bar{\mathbf{B}}x_0 + \bar{\mathbf{C}}\bar{\mathbf{A}}\bar{\mathbf{B}}x_1 + \bar{\mathbf{C}}\bar{\mathbf{B}}x_2 \\
 &\vdots \text{each } x_t \text{ is multiplied by} \\
 &\quad \bar{\mathbf{C}}\bar{\mathbf{A}}^{(t-k)}\bar{\mathbf{B}} \\
 y_k &= \bar{\mathbf{C}}\bar{\mathbf{A}}^k\bar{\mathbf{B}}x_0 + \bar{\mathbf{C}}\bar{\mathbf{A}}^{k-1}\bar{\mathbf{B}}x_1 + \dots + \bar{\mathbf{C}}\bar{\mathbf{A}}^1\bar{\mathbf{B}}x_{k-1} + \bar{\mathbf{C}}\bar{\mathbf{A}}^0\bar{\mathbf{B}}x_k
 \end{aligned}$$

We can represent this as a *global* convolution computation:

$$\mathbf{y} = \bar{\mathbf{K}} * \mathbf{x}$$

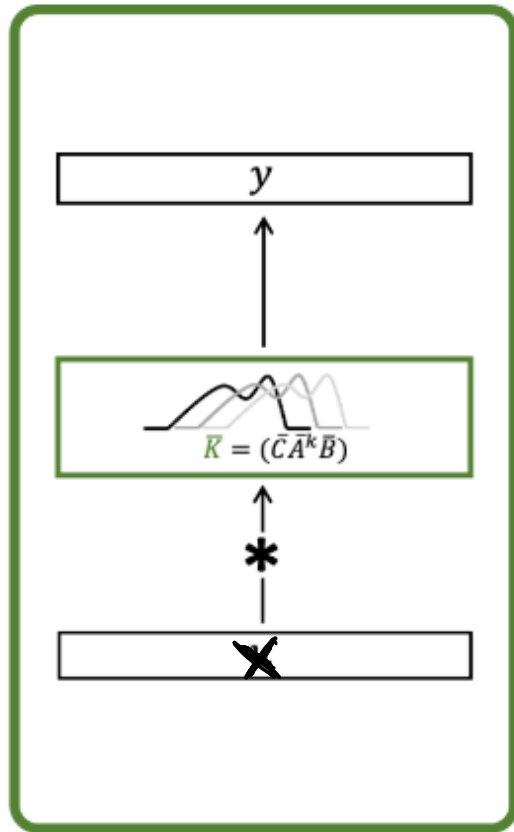
$L = \text{sequence length}$

where the SSM convolution kernel is:

$$\bar{\mathbf{K}} \in \mathbb{R}^L = (\underbrace{\bar{\mathbf{C}}\bar{\mathbf{A}}^0\bar{\mathbf{B}}}_{\in \mathbb{R}}, \underbrace{\bar{\mathbf{C}}\bar{\mathbf{A}}^1\bar{\mathbf{B}}}_{\in \mathbb{R}}, \dots, \bar{\mathbf{C}}\bar{\mathbf{A}}^{L-2}\bar{\mathbf{B}}, \bar{\mathbf{C}}\bar{\mathbf{A}}^{L-1}\bar{\mathbf{B}})$$

$$\begin{aligned}
 h_{k+1} &= \bar{\mathbf{A}}h_k + \bar{\mathbf{B}}x_k \\
 y_{k+1} &= \bar{\mathbf{C}}h_{k+1} \quad \square \quad \text{leave out } \bar{\mathbf{D}}x_k
 \end{aligned}$$

SSM: 1D Convolutional Representation

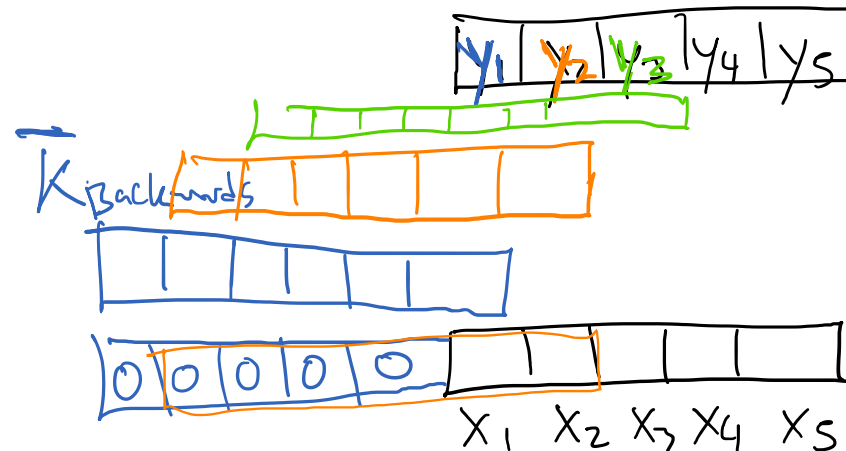


We can represent this as a *global* convolution computation:

$$y = \bar{K} * x$$

where the SSM convolution kernel is:

$$\bar{K} \in \mathbb{R}^L = (\bar{C}\bar{A}^0\bar{B}, \bar{C}\bar{A}^1\bar{B}, \dots, \bar{C}\bar{A}^{L-2}\bar{B}, \bar{C}\bar{A}^{L-1}\bar{B})$$



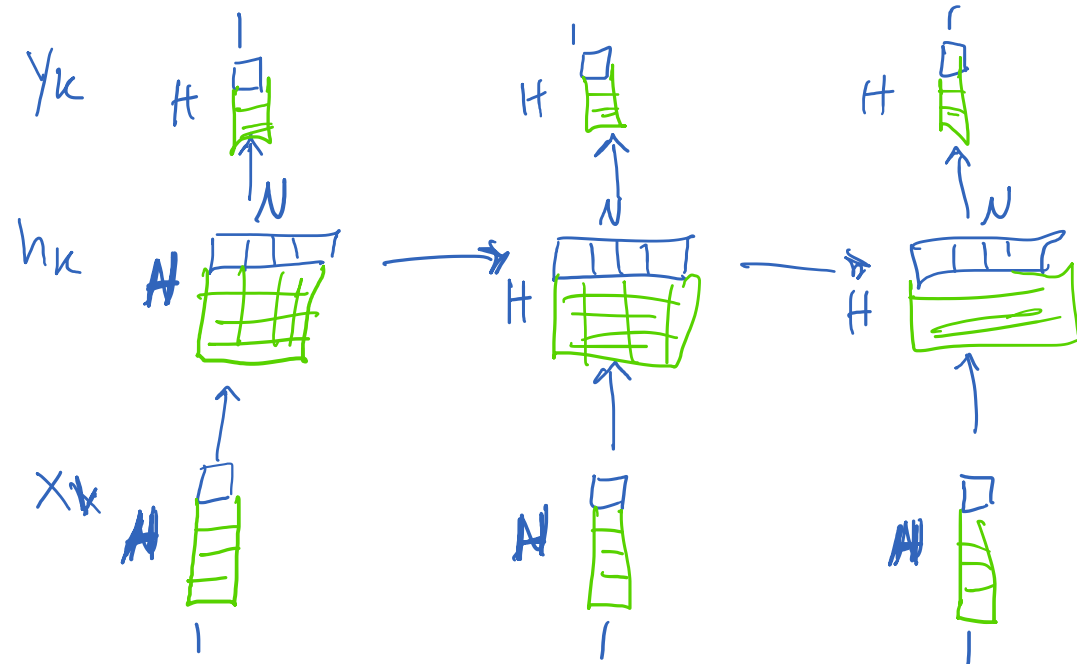
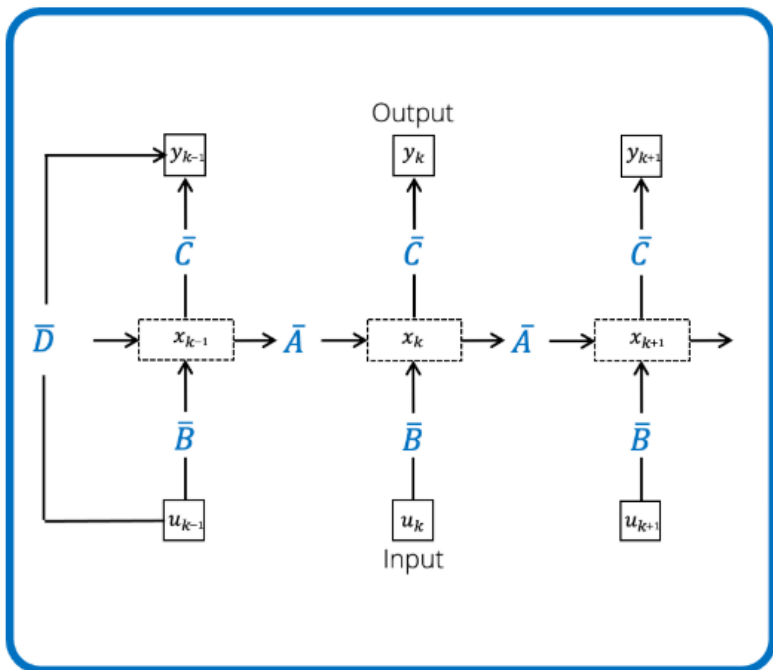
★ can compute all $y_k \forall k \in \{1, \dots, L\}$ in parallel

$$y_k = \bar{C}\bar{A}^k\bar{B}x_0 + \bar{C}\bar{A}^{k-1}\bar{B}x_1 + \dots + \bar{C}\bar{A}^1\bar{B}x_{k-1} + \bar{C}\bar{A}^0\bar{B}x_k$$

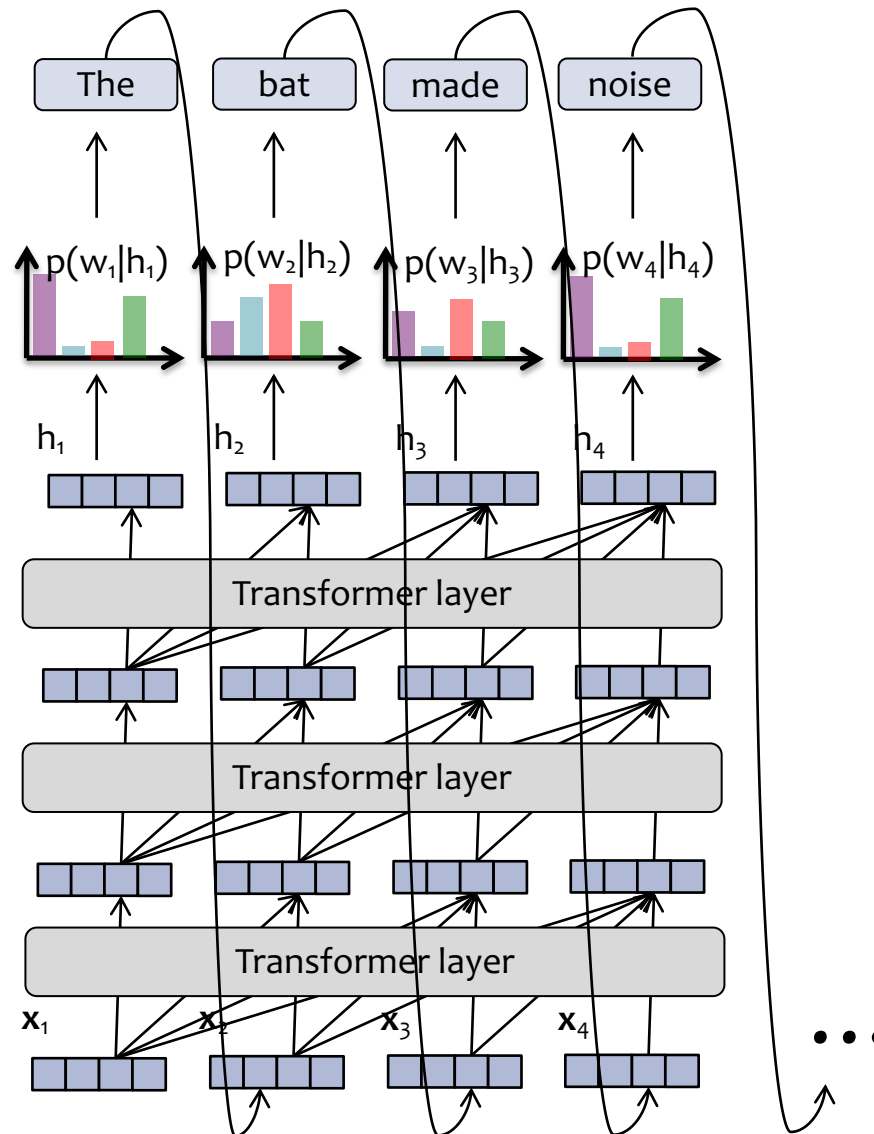
THE STRUCTURED STATE SPACE SEQUENCE MODEL (S4)

SSM as a Neural Network Layer

- We can take H copies of the 1D recurrent representation
- Let each copy have its own parameters
- This is just like multiple (indep.) heads in Attention
- And just like multiple (indep.) channels in Convolution
- So we get...



Transformer Language Model



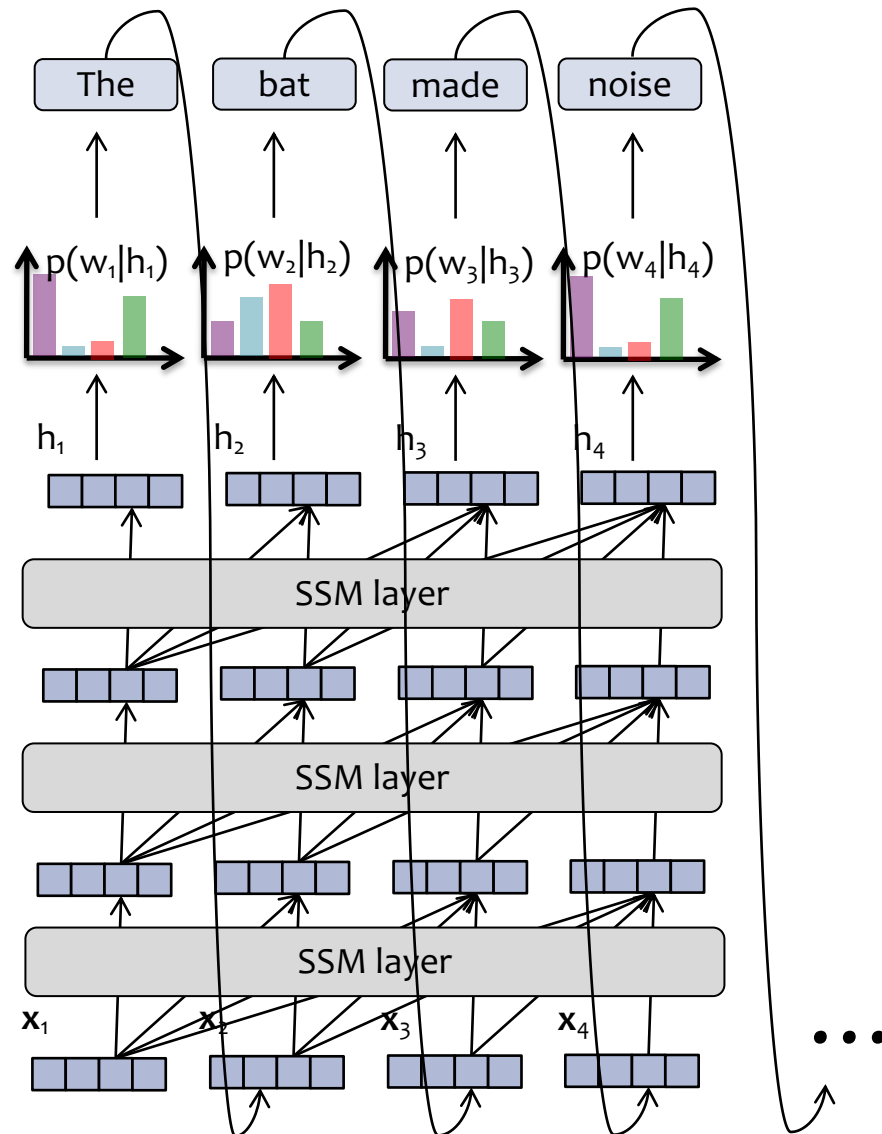
Each layer of a Transformer LM consists of several **sublayers**:

1. attention
2. feed-forward neural network
3. layer normalization
4. residual connections

Each hidden vector looks back at the hidden vectors of the **current and previous timesteps in the previous layer**.

The language model part is just like an RNN-LM.

SSM inside a Deep Language Model



Each layer of an S4 LM consists of several **sublayers as well** including an SSM, nonlinearity, etc.

Each hidden vector looks back at the hidden vectors of the **current and previous timesteps in the previous layer.**

The language model part is just like an RNN-LM or Transformer-LM

Efficiency of SSM, RNN, & Transformer

For SSMs:

1. At test time, generation does NOT need a KV-cache in our **Recurrent representation**, so we can effortlessly generate truly long sequences (unlike Transformers, but just like RNNs)
2. At train time, we can use the **Convolution representation** to do fast parallel training (just like Transformers, but unlike RNNs)

| | Train | Test |
|------------|-------|------|
| Recurrence | Slow | Fast |
| Attention | Fast | Slow |
| SSM | Fast | Fast |

S4 Model

We need several additional tricks to get training to work well:

- HiPPO Matrix

- we initialize the matrix A very carefully

$$h_{k+1} = \bar{A} h_k + \bar{B} x_k$$

- Efficient computation

- we decompose A so that we can compute the kernel K very efficiently and in a numerically stable way

$$h'(t) = A h(t) + B x(t)$$

$$K = (\dots, \underbrace{\bar{C} \bar{A}^{-L} \bar{B}}_{\text{unstable}})$$

Selective State Space Model

with Hardware-aware State Expansion

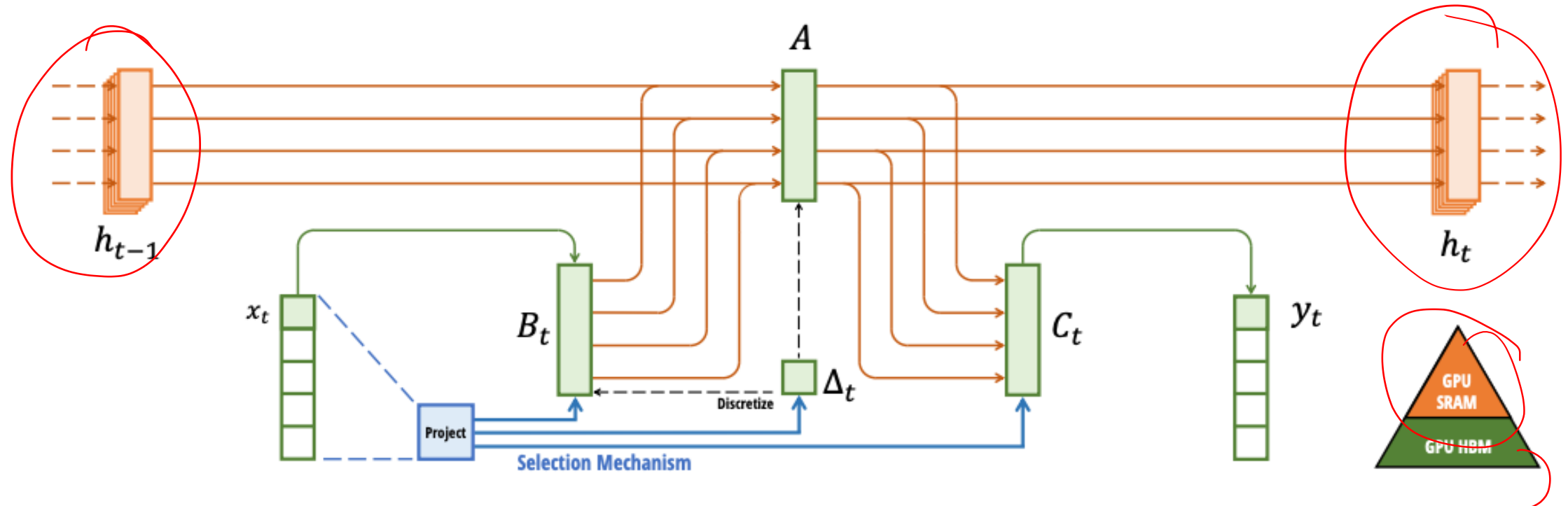
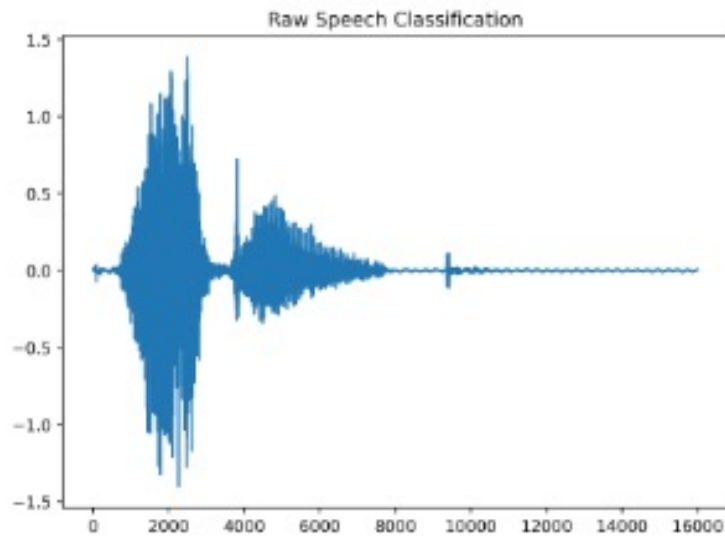


Figure 1: (**Overview.**) Structured SSMs independently map each channel (e.g. $D = 5$) of an input x to output y through a higher dimensional latent state h (e.g. $N = 4$). Prior SSMs avoid materializing this large effective state (DN , times batch size B and sequence length L) through clever alternate computation paths requiring time-invariance: the (Δ, A, B, C) parameters are constant across time. Our selection mechanism adds back input-dependent dynamics, which also requires a careful hardware-aware algorithm to only materialize the expanded states in more efficient levels of the GPU memory hierarchy.

S4 Results: Train and test on different input granularities



| | | ↓ Train: 16K Hz | ↓ Test: 8K Hz | |
|--------------|--------------|-----------------------|---------------------|-------|
| | MFCC | RAW | 0.5× | |
| Transformer | 90.75 | ✗ | ✗ | |
| Performer | 80.85 | 30.77 | 30.68 | ← bad |
| ODE-RNN | 65.9 | ✗ | ✗ | |
| NRDE | 89.8 | 16.49 | 15.12 | |
| ExpRNN | 82.13 | 11.6 | 10.8 | |
| LipschitzRNN | 88.38 | ✗ | ✗ | |
| CKConv | 95.3 | 71.66 | 65.96 | ← |
| WaveGAN-D | ✗ | 96.25 | ✗ | |
| LSSL | 93.58 | ✗ | ✗ | |
| S4 | 93.96 | 98.32 | 96.30 | |

MAMBA

Selective State Space Models

Algorithm 1 SSM (S4)

Input: $x : (B, L, D)$

Output: $y : (B, L, D)$

1: $A : (D, N) \leftarrow$ Parameter

▸ Represents structured $N \times N$ matrix

2: $B : (D, N) \leftarrow$ Parameter

3: $C : (D, N) \leftarrow$ Parameter

4: $\Delta : (D) \leftarrow \tau_{\Delta}(\text{Parameter})$

5: $\bar{A}, \bar{B} : (D, N) \leftarrow \text{discretize}(\Delta, A, B)$

6: $y \leftarrow \text{SSM}(\bar{A}, \bar{B}, C)(x)$

▸ Time-invariant: recurrence or convolution

7: **return** y

Algorithm 2 SSM + Selection (S6)

Input: $x : (B, L, D)$

Output: $y : (B, L, D)$

1: $A : (D, N) \leftarrow$ Parameter

▸ Represents structured $N \times N$ matrix

2: $B : (B, L, N) \leftarrow s_B(x) \leftarrow = \text{linear}(W_B, x)$

3: $C : (B, L, N) \leftarrow s_C(x) \leftarrow = \text{linear}(W_C, x)$

4: $\Delta : (B, L, D) \leftarrow \tau_{\Delta}(\text{Parameter} + s_{\Delta}(x))$

5: $\bar{A}, \bar{B} : (B, L, D, N) \leftarrow \text{discretize}(\Delta, A, B)$

6: $y \leftarrow \text{SSM}(\bar{A}, \bar{B}, C)(x)$

▸ Time-varying: recurrence (*scan*) only

7: **return** y

- **Selective** state space models differ from S4 in that they let the parameters B and C vary at each timestep

Selective State Space Models

Algorithm 1 SSM (S4)

Input: $x : (B, L, D)$

Output: $y : (B, L, D)$

1: $A : (D, N) \leftarrow$ Parameter

▷ Represents structured $N \times N$ matrix

2: $B : (D, N) \leftarrow$ Parameter

3: $C : (D, N) \leftarrow$ Parameter

4: $\Delta : (D) \leftarrow \tau_{\Delta}(\text{Parameter})$

5: $\bar{A}, \bar{B} : (D, N) \leftarrow \text{discretize}(\Delta, A, B)$

6: $y \leftarrow \text{SSM}(\bar{A}, \bar{B}, C)(x)$

▷ Time-invariant recurrence or convolution

7: **return** y

$$h_t = \underline{A} h_{t-1} + \underline{B} x_t$$

$$y_t = \underline{C}^{\top} h_t$$

Algorithm 2 SSM + Selection (S6)

Input: $x : (B, L, D)$

Output: $y : (B, L, D)$

1: $A : (D, N) \leftarrow$ Parameter

▷ Represents structured $N \times N$ matrix

2: $B : (B, L, N) \leftarrow s_B(x)$

3: $C : (B, L, N) \leftarrow s_C(x)$

4: $\Delta : (B, L, D) \leftarrow \tau_{\Delta}(\text{Parameter} + s_{\Delta}(x))$

5: $\bar{A}, \bar{B} : (B, L, D, N) \leftarrow \text{discretize}(\Delta, A, B)$

6: $y \leftarrow \text{SSM}(\bar{A}, \bar{B}, C)(x)$

▷ Time-varying: recurrence (*scan*) only

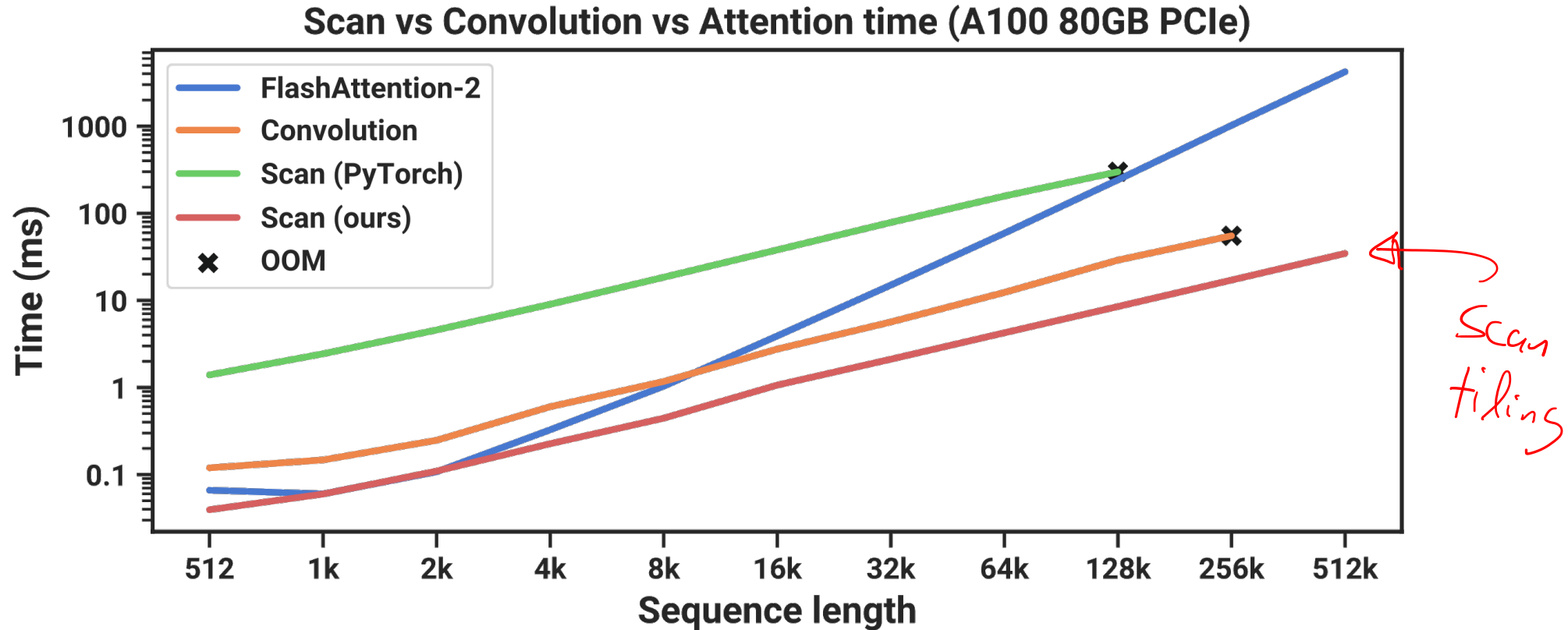
7: **return** y

$$h_t = \underline{A}_t h_{t-1} + \underline{B}_t x_t$$

$$y_t = \underline{C}_t^{\top} h_t$$

Mamba's Scan Implementation

- We can no longer compute the kernel K once up front
- Instead we perform an efficient scan implementation



Scan tiling

Mamba Results

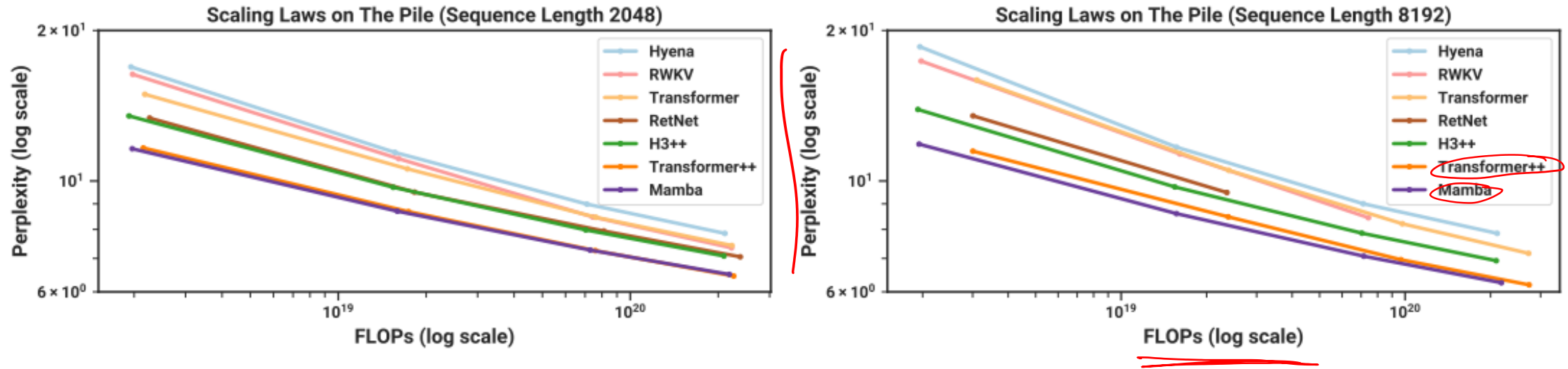
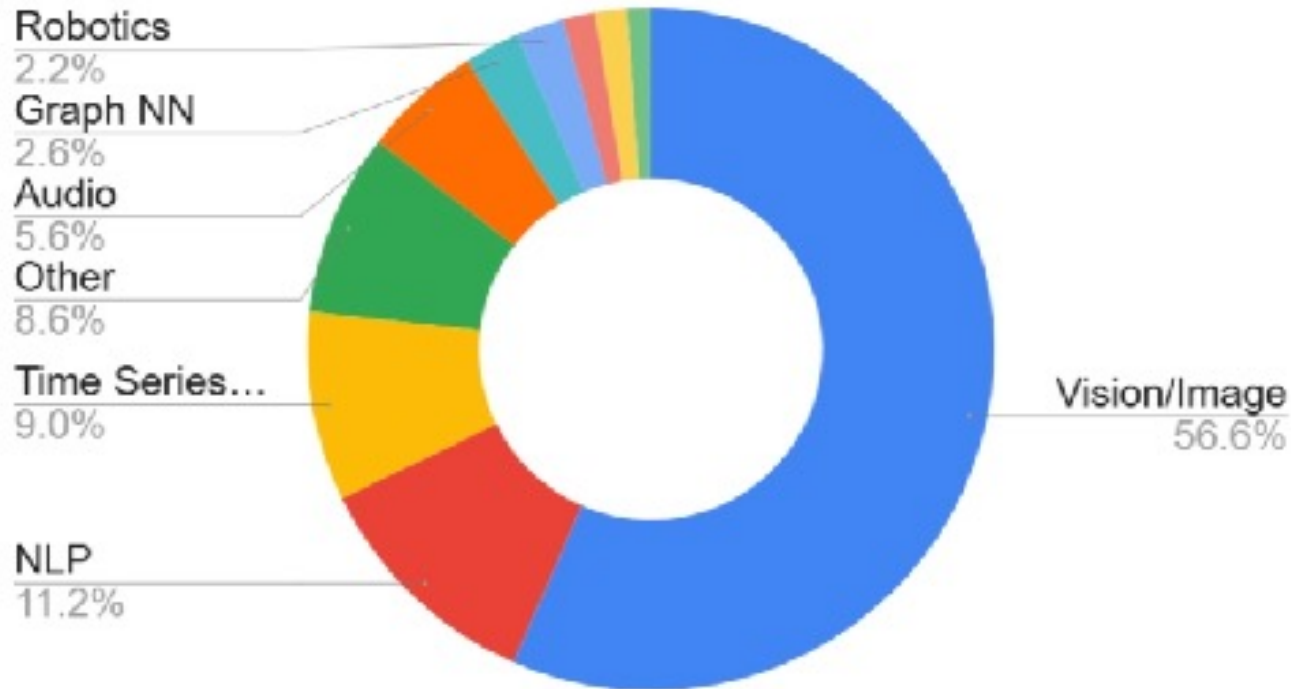


Figure 4: (**Scaling Laws.**) Models of size $\approx 125M$ to $\approx 1.3B$ parameters, trained on the Pile. Mamba scales better than all other attention-free models and is the first to match the performance of a very strong “Transformer++” recipe that has now become standard, particularly as the sequence length grows.

- Main takeaway: Mamba is the first (only?) non-attention based LM to challenge a Transformer

Mamba Use in the Real World

Mamba Paper Categories - 267 papers up till June 27th



Strong out-of-the-box
on **general modalities**
(not just language!)

Motivation

- <https://www.isattentionallyouneed.com/>

Is Attention All You Need?



Current Status: Yes

Time Remaining: 631d 22h 55m 11s

Proposition:

On January 1, 2027, a Transformer-like model will continue to hold the state-of-the-art position in most benchmarked tasks in natural language processing.