# Mamba: Linear-Time Sequence Modeling with Selective State Spaces
*S4 (easy mode)*

# Why S4? Not a basic SSM?



$$\dot{x} = Ax + Bu$$
$$y = Cx + Du$$

**Continuous State Space**

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 2 & 0 \\ 1 & 3 & 3 \end{bmatrix}$$

**Long-Range Dependencies**

$$x = \bar{A}x + \bar{B}u$$
$$y = \bar{C}x + \bar{D}u$$

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

**Fast Discrete Representations**

(1) SSM can model long sequence
(2) However, a basic SSM has prohibitive computation and memory requirements!
(3) S4 computes efficiently and require less memory (30x faster + 400x less usage than conventional SSM )

S4 introduces a novel parameterization that efficiently swaps between these representations,
allowing it to handle a wide range of tasks, be efficient at both training and inference, and excel at long sequences.
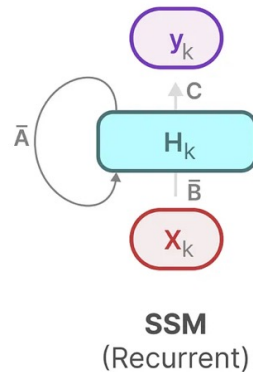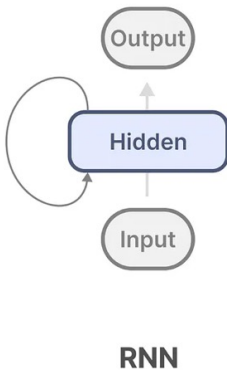
# Three representations: Recurrent

$$x = \bar{A}x + \bar{B}u$$
$$y = \bar{C}x + \bar{D}u$$

Recurrent Representation

Bilinear method:

$$\bar{A} = \left(I - \frac{\Delta}{2} * A\right)^{-1} \left(I + \frac{\Delta}{2} * A\right)$$
$$\bar{B} = \left(I - \frac{\Delta}{2} * A\right)^{-1} \Delta B$$
$$\bar{C} = C$$



**Timestep 0**

$$h_0 = \bar{B}x_0$$
$$y_0 = Ch_0$$

Timestep -1
does not exist so
$Ah_{-1}$
can be ignored

**Timestep 1**

$$h_1 = \bar{A}h_0 + \bar{B}x_1$$
$$y_1 = Ch_1$$

State of
**previous** timestep
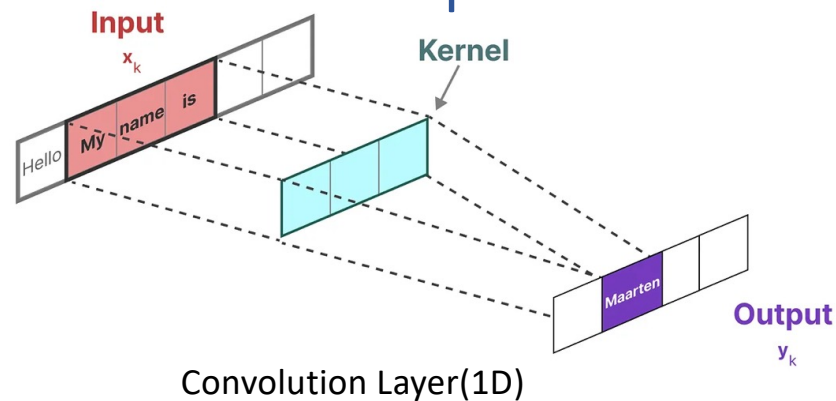
State of
**current** timestep

**Timestep 2**

$$h_2 = \bar{A}h_1 + \bar{B}x_2$$
$$y_2 = Ch_2$$

State of
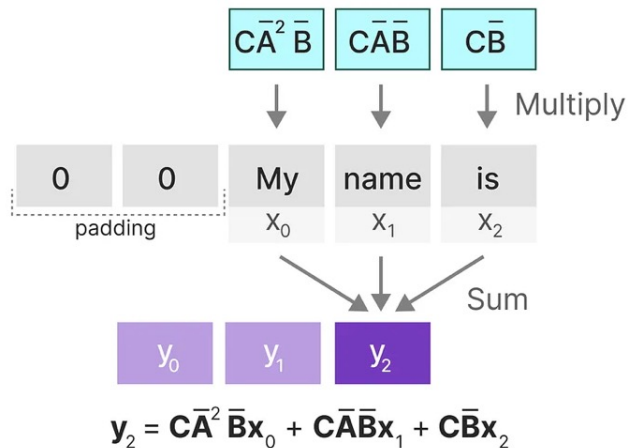**previous** timestep

State of
**current** timestep

RNN

SSM
(Recurrent)

# Three representations: Convolutional



**Input** $x_k$

**Kernel**

Hello | My | name | is

Maarten

**Output** $y_k$

Convolution Layer(1D)

**Kernel**

| $C\bar{A}^2\bar{B}$ | $C\bar{A}\bar{B}$ | $C\bar{B}$ |

↓ ↓ ↓ Multiply

**Input** $(x_k)$

| 0 | 0 | My | name | is |

padding | | $x_0$ | $x_1$ | $x_2$

Sum

**Output** $(y_k)$

| $y_0$ | $y_1$ | $y_2$ |

$$y_2 = C\bar{A}^2\bar{B}x_0 + C\bar{A}\bar{B}x_1 + C\bar{B}x_2$$

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

Convolutional Representation

| time | hidden | prediction |
|------|--------|------------|
| 0 | $\bar{B}x_0$ | $C\bar{B}x_0$ |
| 1 | $\bar{A}h_0 + \bar{B}x_1$ | $C\bar{A}\bar{B}x_0 + C\bar{B}x_1$ |
| 2 | $\bar{A}(\bar{A}h_0 + \bar{B}x_1) + \bar{B}x_2$ | $C\bar{A}^2\bar{B}x_0 + C\bar{A}\bar{B}x_1 + C\bar{B}x_2$ |
| ... | ... | ... |

If conv size is fixed, like conv size = 3,
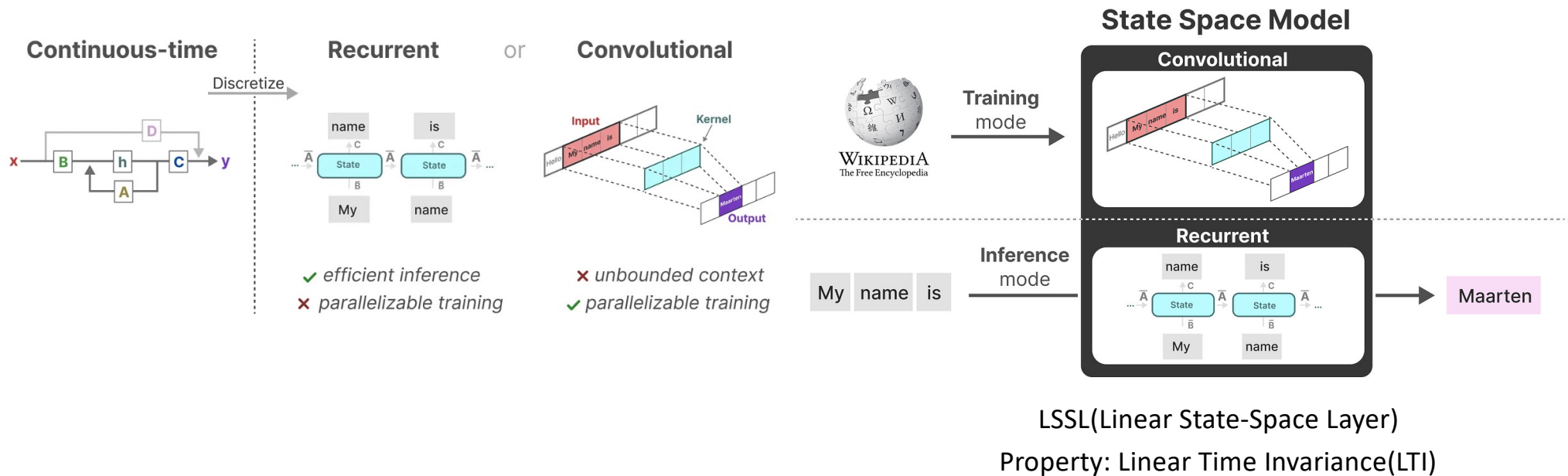Conv kernel should be: $\bar{K} = (C\bar{A}^{L-1}B, \dots, C\bar{A}^2\bar{B}, C\bar{A}\bar{B}, C\bar{B})$

Difficulties: L times for computing $y = \bar{K} * x$
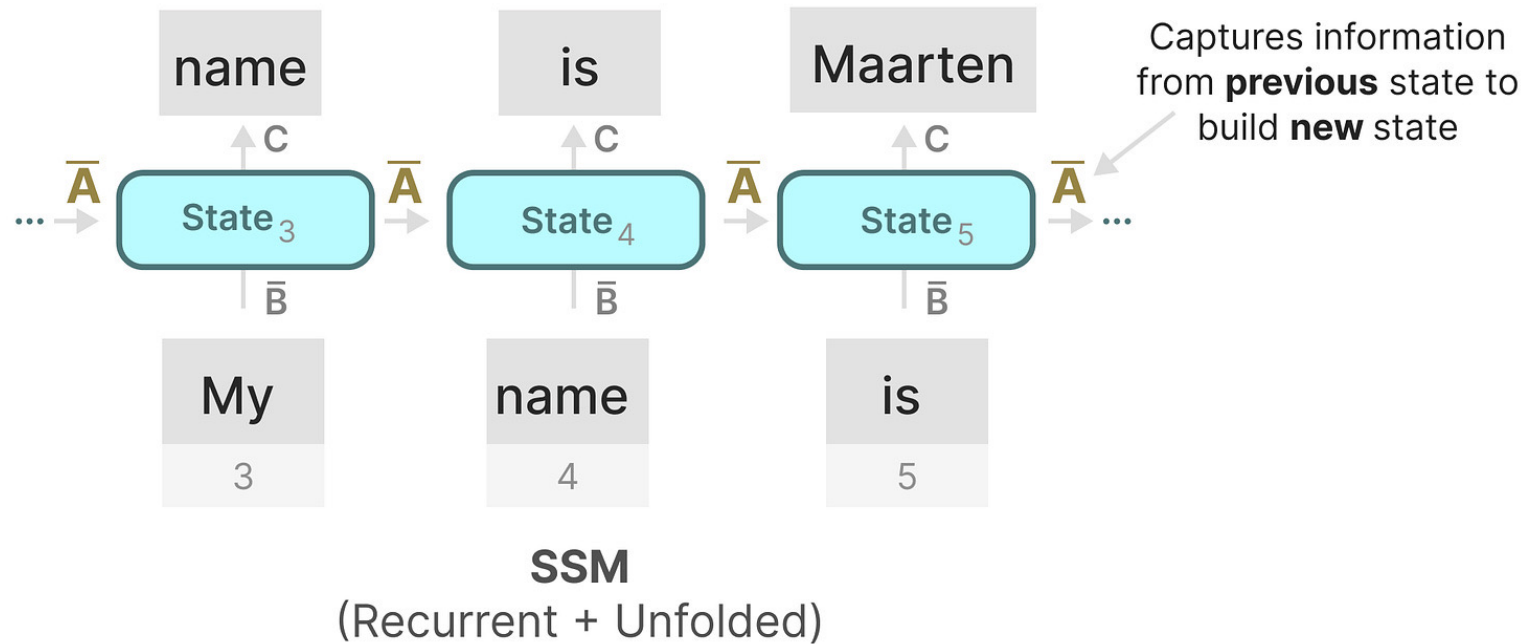Solve: truncated generating function ( convert the power of matrix to the inverse of matrix, like
$$y = fun_1(A) * fun_2(A^{-1}) * x$$

# Summary of three representations



**Continuous-time** — Discretize → **Recurrent** or **Convolutional**

✔ efficient inference
✘ parallelizable training

✘ unbounded context
✔ parallelizable training

**State Space Model**

Training mode → Convolutional

Inference mode → Recurrent → Maarten

LSSL(Linear State-Space Layer)
Property: Linear Time Invariance(LTI)

# The importance of Matrix A



name    is    Maarten

Captures information from **previous** state to build **new** state

$\overline{A}$ ... → State$_3$ → $\overline{A}$ → State$_4$ → $\overline{A}$ → State$_5$ → $\overline{A}$ ...

C ↑    C ↑    C ↑

$\overline{B}$ ↓    $\overline{B}$ ↓    $\overline{B}$ ↓
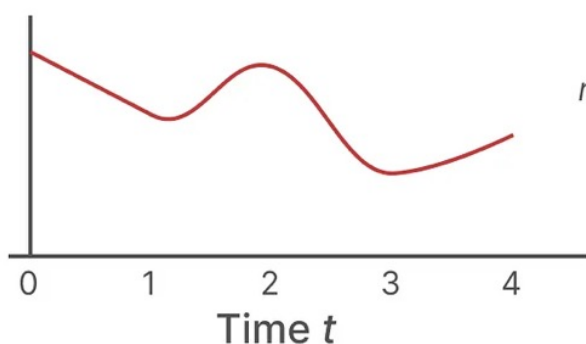
My    name    is

3    4    5

**SSM**
(Recurrent + Unfolded)

How can we get a Matrix A to retain the large memory (since it only look at previous states)

# HIPPO Matrix

**Input Signal**



**Reconstructed Signal**

HiPPO
(*compress* and
*reconstruct* signal
information)

small degration
of **newer** steps

large degration
of **older** steps

Mathematically, it does so by tracking the coefficients of a Legendre polynomial which allows it to approximate all of the previous history.

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

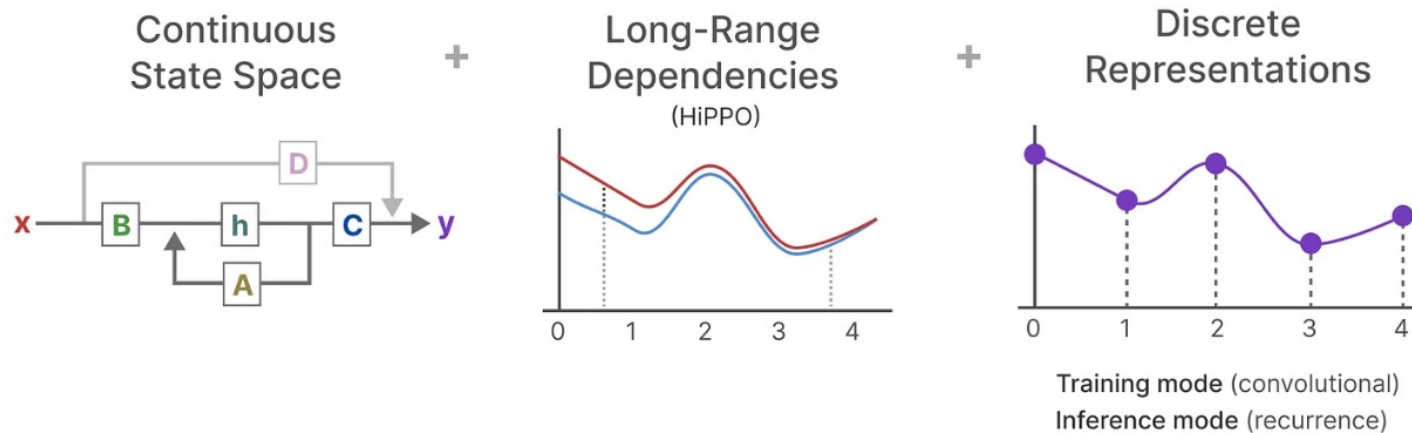https://hazyresearch.stanford.edu/blog/2020-12-05-hippo

https://proceedings.neurips.cc/paper/2019/hash/952285b9b7e7a1be5aa7849f32ffff05-Abstract.html

# Summary of S4

- State Space Model
- HiPPO for handling long-range dependencies
- Discretization for creating recurrent and convolution representations



**Structured State Spaces for Sequences** (S4)

Continuous State Space + Long-Range Dependencies (HiPPO) + Discrete Representations

Training mode (convolutional)
Inference mode (recurrence)

Differ from basic SSMs:
(1) Modeling challenge for LRD: using a special formula for the A matrix (HiPPO).
(2) Computational Challenge: introducing a special representation and algorithm to be able to work with this matrix (truncated generating function)