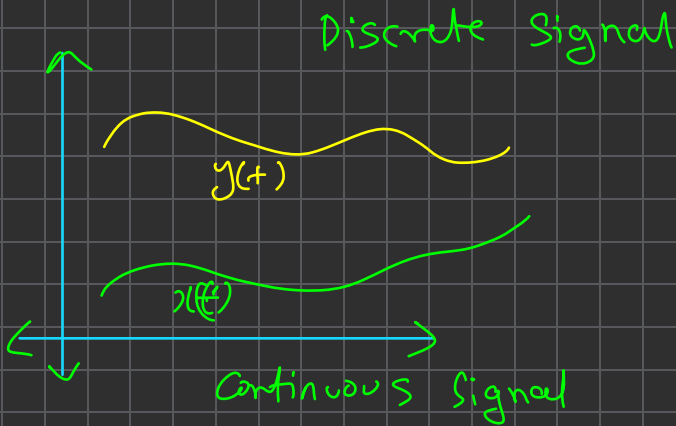
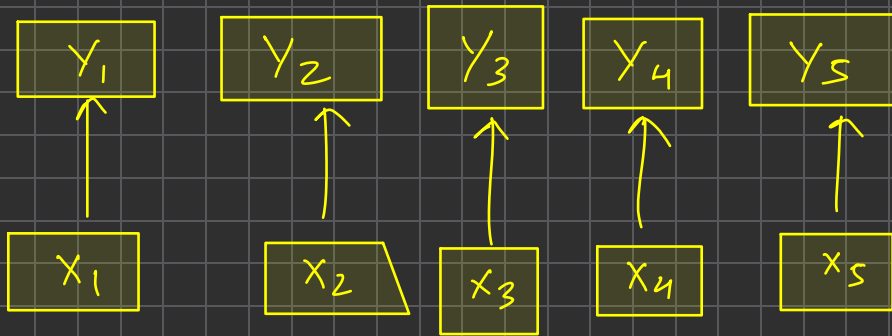


# Sequence Modeling



What is Mamba??



State Space  
Model  
(SSM)

+

Linear  
MLP  
(gated)

+

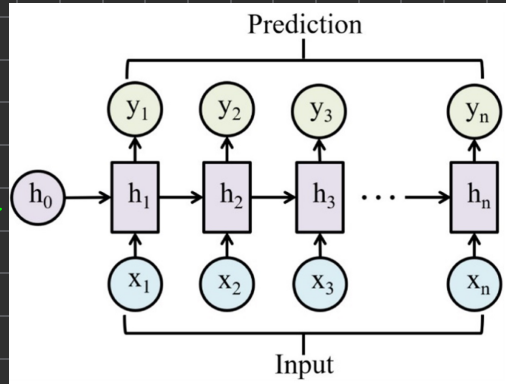
Some  
selection  
mechanism

① RNNs  $\rightarrow O(LD)$

Simple RNN

$L \rightarrow$  context length  
(Sequence length)

$D \rightarrow$  model parameter  
or hidden size  
( $d_{model}$ )



① Memory

Efficient

(Space complexity)

① vanishing & Exploding gradients

② Limited context size  
(can't retain info)

③ cannot be parallelized  
(Making training slow)

② CNN  $\rightarrow O(L^2D)$

$\rightarrow$  kernels

$\rightarrow$  Convolution Filters

① context window is small

② Requires stacking many Conv layers

(Ex: UNET)  
Increasing  
computation  
cost  
 $\nearrow$

### ③ Transformers

- ① quadratic complexity  $\rightarrow O(L^2 D)$
- ② Easily parallelizable
- ③ Inference is slow (unless KV cache is used)
- ④ Limited vocab size / sequence length

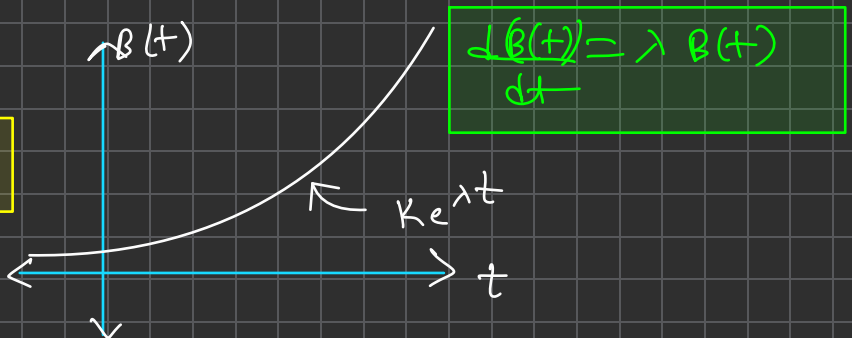
### \* Differential Equation

$$\frac{db}{dt} = \lambda b(t)$$

$$\text{let } b(t) = ke^{\lambda t}$$

$$\frac{d(b(t))}{dt} = ke^{\lambda t} (\lambda) = \lambda ke^{\lambda t} = \lambda b(t)$$

$$\dot{b} = \lambda b$$



# \* State Space Model (linear & Time-Invariant)

objective: Map  $x(t)$  to an output signal  $y(t)$

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

$$y(t) = C \cdot h(t) + D \cdot x(t)$$

parameter matrices  $\rightarrow A, B, C, D$

(do not depend on Time)

$h(t) \rightarrow$  state representation

$y(t) \rightarrow$  output signal

$\rightarrow$  Discretization

$h(0), h(1), h(2), h(3), \text{etc}$

So finding  $h(t)$  we want  $h(t_k) = h(k\Delta)$

$\Delta \rightarrow$  step size

$$\therefore b'(t) = \lambda b(t)$$

$$\lim_{\Delta \rightarrow 0} \frac{b(t+\Delta) - b(t)}{\Delta} = b'(t)$$

$\Delta \rightarrow$  very small

$$\therefore \frac{b(t+\Delta) - b(t)}{\Delta} \cong b'(t)$$

$$b(t+\Delta) \cong b'(t)\Delta + b(t)$$

$$\therefore \boxed{b(t+\Delta) \cong \lambda b(t)\Delta + b(t)}$$

recurrent Formulation

let  $\lambda = 2, \Delta = 1$

① at  $t=0 \quad b=5$

$$b(0+1) \cong 2(b(0)1 + b(0))$$

$$\cong 2 \times 5 \times 1 + 5$$

$$b(1) \cong 15$$

$\therefore$  at  $t=1, b=15$

$\therefore$  Discretization of our State Space Model

$$h(t+\Delta) \cong \Delta h'(t) + h(t)$$

$$\therefore \overbrace{h'(t) = A \cdot h(t) + B \cdot x(t)}$$

$$h(t+\Delta) \cong (A \cdot h(t) + B \cdot x(t))\Delta + h(t)$$

$$\cong \Delta A h(t) + \Delta B x(t) + h(t)$$

$$= h(t) (I + \Delta A) + \Delta B x(t)$$

$$\boxed{h(t+\Delta) \cong \bar{A} h(t) + \bar{B} x(t)}$$

$$h(t+\Delta) \approx \bar{A}h(t) + \bar{B} \cdot x(t)$$

$$\bar{A} = (\mathbf{I} + \Delta A)$$

$\Delta \rightarrow$  Time step

$$\bar{B} = \Delta B$$

$\mathbf{I} \rightarrow$  Identity matrix

In paper :-

$$\left. \begin{array}{l} h'(t) = Ah(t) + B \cdot x(t) \\ y(t) = Ch(t) \end{array} \right\} \begin{array}{l} \text{upon solving} \\ h_t = \bar{A}h_{t-1} + \bar{B}x_t \\ y_t = Ch_t \end{array}$$

using Zero-Order hold (ZOH)

$$\bar{A} = f_A(\Delta, A) = e^{(\Delta A)}$$

$$\bar{B} = f_B(\Delta, A, B) = \frac{(e^{(\Delta A)} - \mathbf{I}) \cdot \Delta B}{\Delta A}$$

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

$$y_t = Ch_t$$

parameters  $\rightarrow \Delta, A, B, C$   
(learnable)

# Recurrent Computation

Initial condition  $\rightarrow$  at  $t=0$ ,  $h_{t-1} = 0$   
(Initial state is zero)

$$\therefore h_t = \bar{A}h_{t-1} + \bar{B}x_t \quad y_t = Ch_t$$

$$\therefore t=0 \quad h_{t-1} = 0$$

$$h_0 = 0 + \bar{B}x_0$$

$$y_0 = Ch_0$$

$$\therefore t=1$$

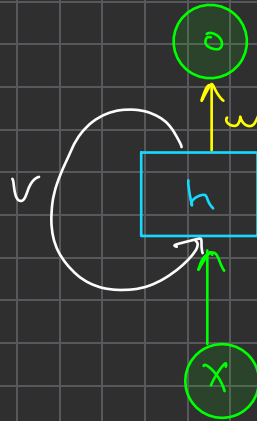
$$h_1 = \bar{A}h_0 + \bar{B}x_1$$

$$y_1 = Ch_1$$

$$\therefore t=2$$

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

$$y_2 = Ch_2$$



RC  $\rightarrow$  Memory efficient (during Inference)

$\hookrightarrow$  not good during Training

parallelizable!



# Convolutional Computation

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

$$y_t = c h_t$$

at  $t=0$ ,  $h=0$   
given

①  $t=0$   
 $\therefore h_0 = \bar{A} \overset{0}{h} + \bar{B} x_0$

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

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

②  $t=1$

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

$$h_1 = \bar{A} \bar{B} x_0 + \bar{B} x_1$$

$$y_1 = c h_1$$

$$y_1 = c (\bar{A} \bar{B} x_0 + \bar{B} x_1)$$

$$y_1 = c \bar{A} \bar{B} x_0 + c \bar{B} x_1$$

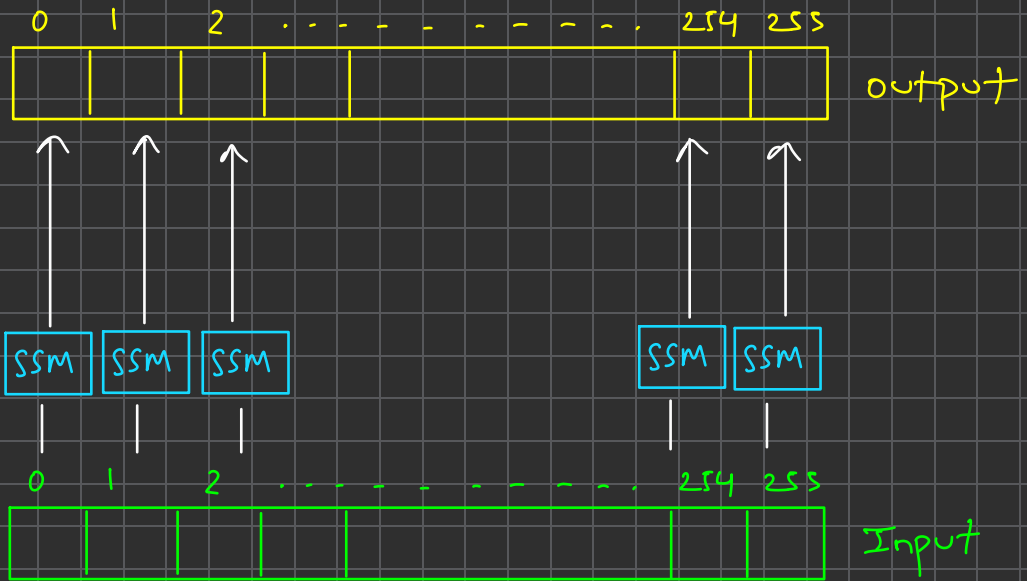
$$\therefore y_k = [c(\bar{A})^k \bar{B} x_0 + c(\bar{A})^{k-1} \bar{B} x_1 + \dots + c \bar{A} \bar{B} x_{k-1} + c \bar{B} x_k]$$

Kernel  $\bar{K} = (c \bar{B}, c \bar{A} \bar{B}, \dots, c(\bar{A})^k \bar{B}, \dots)$  [3ca]

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

- Can be parallelized : because output  $y_k$  doesn't depend on  $y_{k-1}$ , but only on the kernel & input

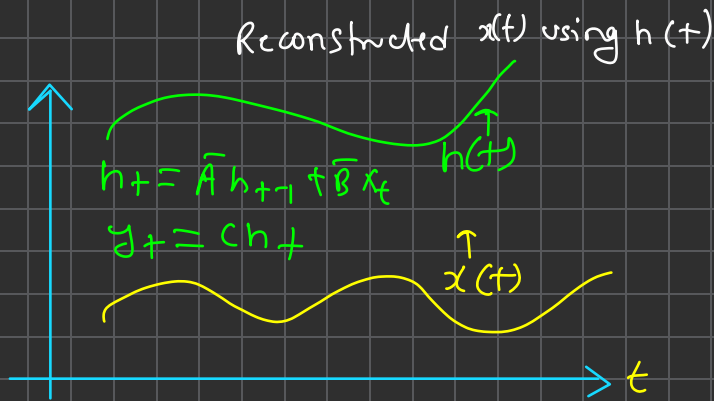
# Multi-dimension SSM



In this way we can parallelize all these operations by working on batches of input  
 $A, B, C, D, x(t), y(t)$  become Tensors

# HIPPO THEORY

With hippo theory, we Build the A matrix in such a way that it approximates all the Input signal seen so far into a vector of coefficient (representing Legendre polynomials).



(Hippo Matrix)  $\hookrightarrow (n \times k)$

$\downarrow$

$O(n \log n)$   
order

$A_{nk} = \begin{cases} (2n+1)^{1/2} (2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$

$\uparrow$

Hippo Legs (Scaled Legendre)

- H3 (High-order polynomial projection operator Theory) is a mathematical framework designed to efficiently store & retrieve sequential information by projecting data onto special basis functions

For Example  $N=3$  (State Space Model size)

$$A_{11} = n+1 = 3+1 = 4 \quad (k=n)$$

$$A_{12} = 0 \quad (k > n)$$

$$A_{13} = 0 \quad (k > n)$$

$$A_{21} = (2n+1)^{1/2} (2k+1)^{1/2} = (2(2)+1)^{1/2} (2(1)+1)^{1/2} \\ \therefore (n > k) \quad = \sqrt{5} \sqrt{3} = \sqrt{15}$$

$$A_{nk} = \begin{bmatrix} 4 & 0 & 0 \\ \sqrt{5} & - & - \\ - & - & - \end{bmatrix}^T$$

# Improving SSMs with Selection

## Selection Mechanism

$$SSM(\bar{A}, \bar{B}, C)(x)$$

$$h_t = \bar{A}(G \odot h_{t-1}) + \bar{B}(F \odot x(t))$$

①  $\rightarrow$  element wise multiplication

$G \rightarrow$  Selection gate controlling how much memory to retain

$F \rightarrow$  Feature gate, determines  $x(t)$  in interacts with memory

Algorithm ② : Dynamic Selection

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

$A \rightarrow (D, N)$  parameter (Hippo-H3)

$B \rightarrow (D, N)$  parameter

$C \rightarrow (D, N)$  parameter

$\Delta \rightarrow (D)$  Step Size (20H)

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

$$y \leftarrow SSM(\bar{A}, \bar{B}, C, F, G)(x)$$

(Time invariant) (recurrent + conv) can be used

$\uparrow$   
gated mechanism

## Algorithm ②

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

$A \in \mathbb{R}^{N \times N}$   $A \rightarrow (D, N)$  parameter (Hippo  $M_3$ )

$B \in \mathbb{R}^{N \times D}$   $B \rightarrow (B, L, N)$

$(B, L, N) \leftarrow S_B(x)$

transforms the input Sequence into a 3D Tensor

$C \in \mathbb{R}^{B \times L \times N}$   $C \rightarrow (B, L, N)$

$\Delta \rightarrow (B, L, D) \leftarrow \text{step size}$

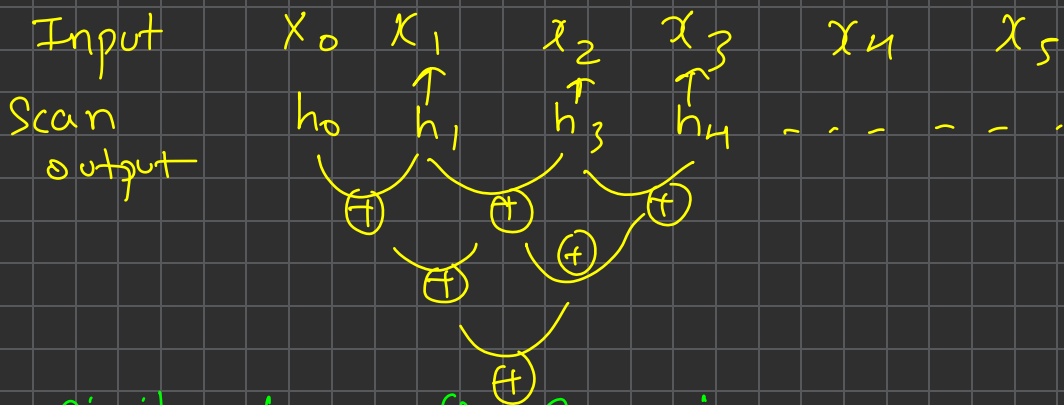
compute  $\bar{A}, \bar{B} : (B, L, D, N)$

$\uparrow \text{discretize}(\Delta, A, B)$

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

$\hookrightarrow$  Time varying  $\therefore$  Recurrent (scan) only

# The Scan operation



Similar to prefix Sum in an array

## Kernel Fusion :

When we perform a tensor operation, our deep Learning framework (eg pytorch) Loads the tensor in the fast memory (SRAM) of the GPU, performs the operations (eg. matrix multiplication), & then saves back the result in the high Bandwidth Memory of the GPU

(Use a single cuda kernel)

