

Text Classification with Mamba - Project

TA Minh-Duc Bui
STA Khai-Xuan Trinh

Outline

1. Motivation

2. State Space Models

3. Mamba

4. Coding

Outline

1. Motivation

2. State Space Models

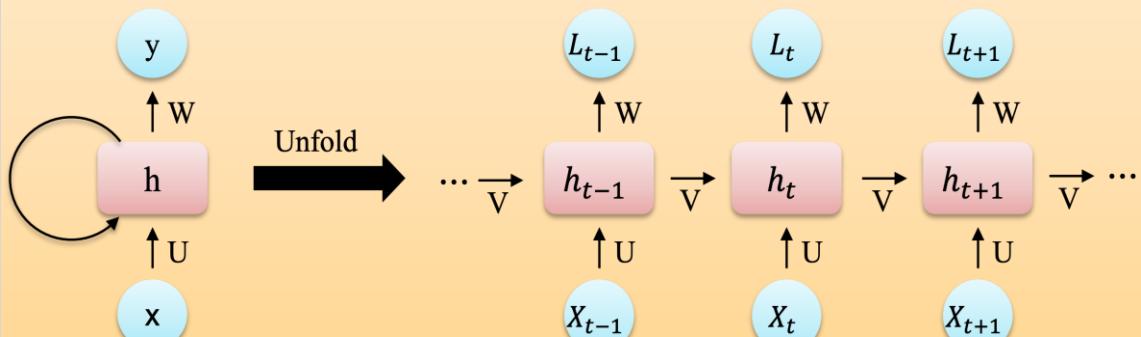
3. Mamba

4. Coding

RNNs vs. Transformers

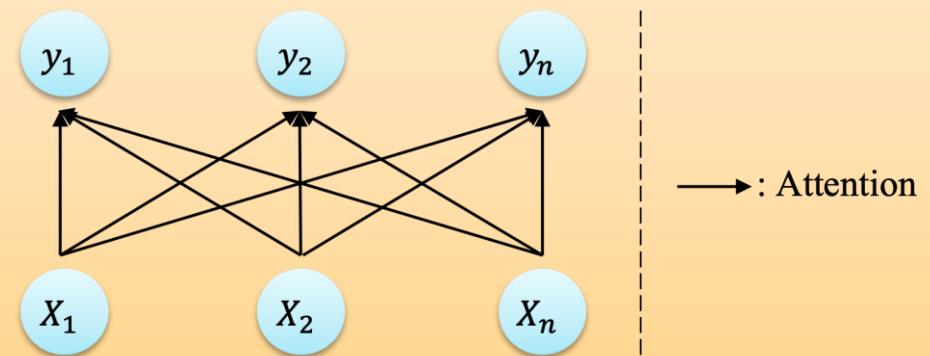
RNNs

- train stage: $O(N)$, **non-parallelable**
- inference stage: $O(N)$
- $O(1)$ per token



Transformers

- train stage: $O(N^2)$, **parallelable**
- inference stage: $O(N^2)$
- $O(N)$ per token

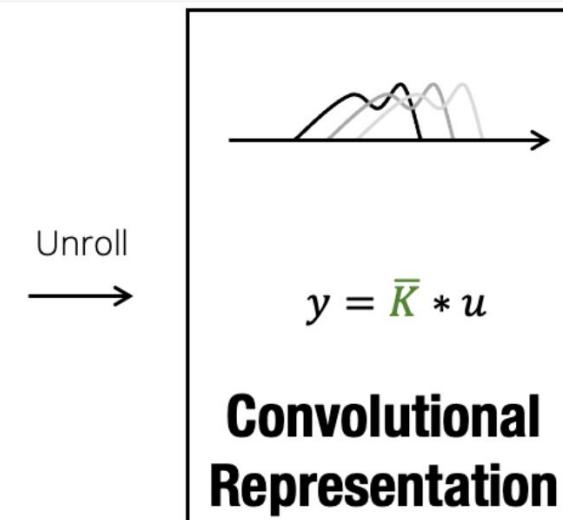
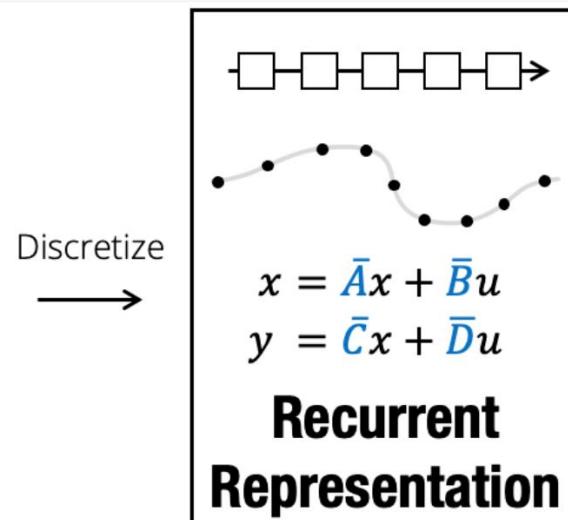
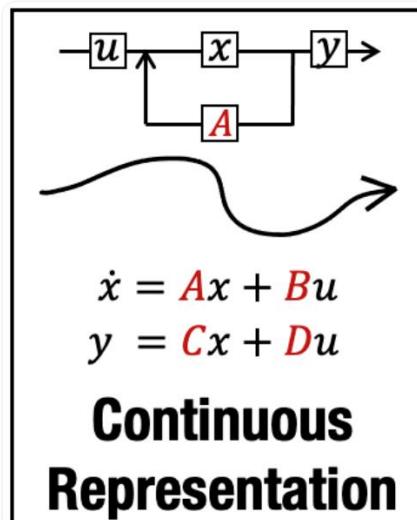


Ideal Models

- train stage: $O(N)$, **parallelable**
- inference stage: $O(N)$
- $O(1)$ per token

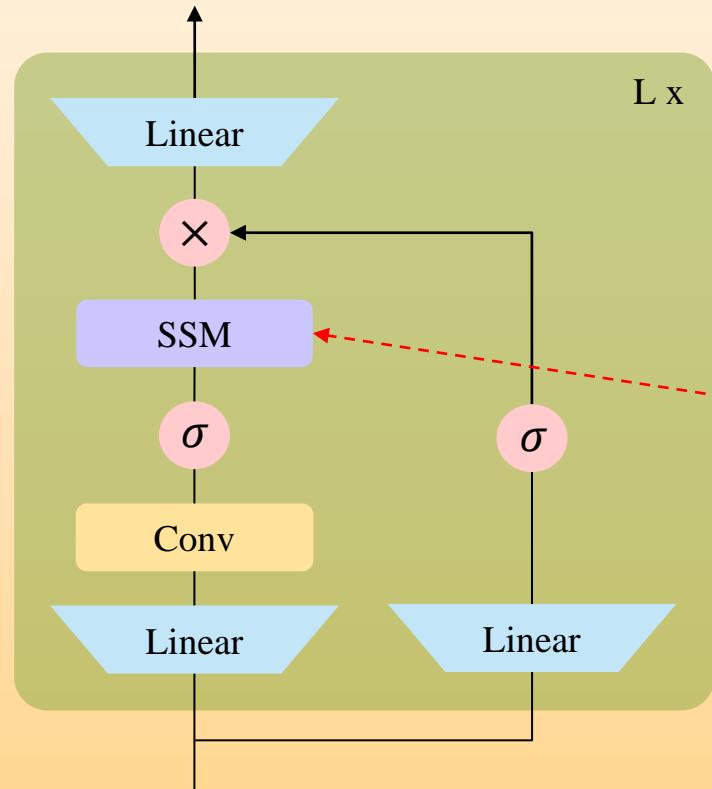
RNNs vs. Transformers

	RNN	Transformer	Ideal Model
Parallelizable	no	yes	yes
Training Stage	$O(n)$	$O(n^2)$	$O(n)$
Inference Stage	$O(n)$	$O(n^2)$	$O(n)$
Inference per Token	$O(1)$	$O(n)$	$O(1)$

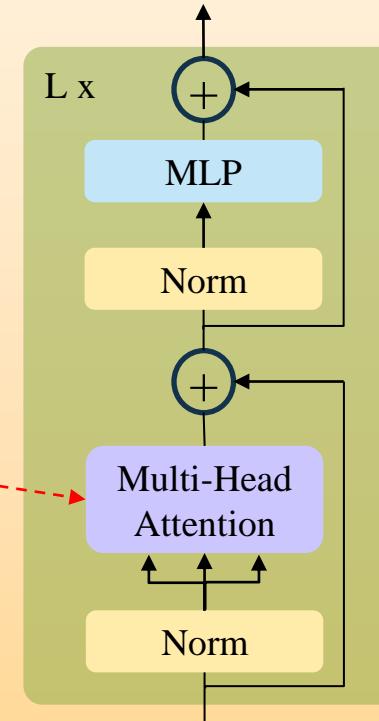


Mamba Block

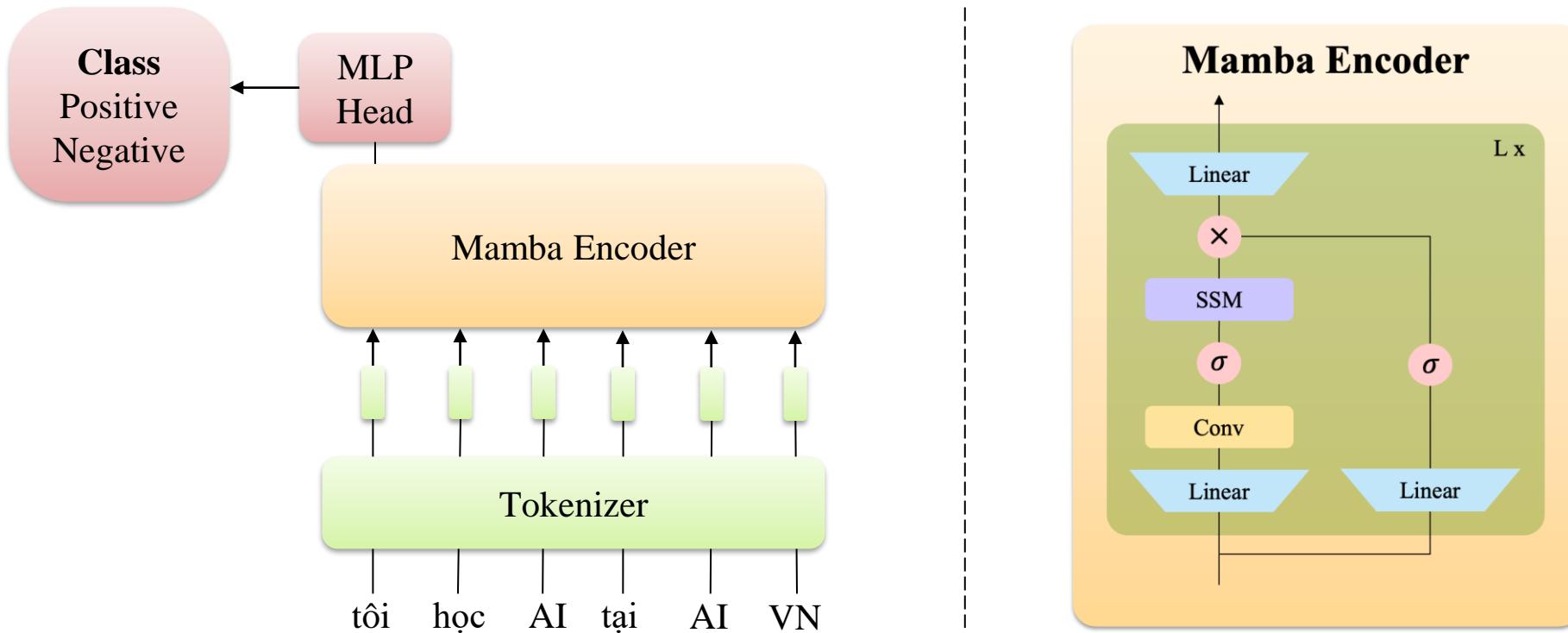
Mamba Encoder



Transformer Encoder



Mamba Architecture



Outline

1. Motivation

2. State Space Models

3. Mamba

4. Coding

State Space Models (SSMs)

What

The SSMs are traditionally used in control theory to model a dynamic system via state variables.

How

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

Why

- train stage: $\mathcal{O}(N)$, parallelable
- inference stage: $\mathcal{O}(N)$
- $\mathcal{O}(1)$ per token

Benefits

- suitable for long-range tasks
- faster training and inferencing
- low computational cost
- less memory

Discretize SSMs

Continuous 😞

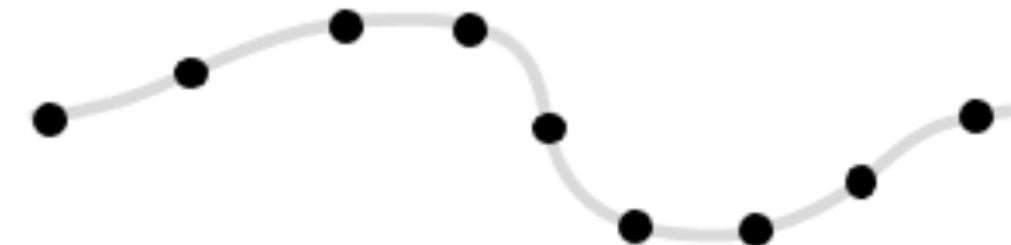
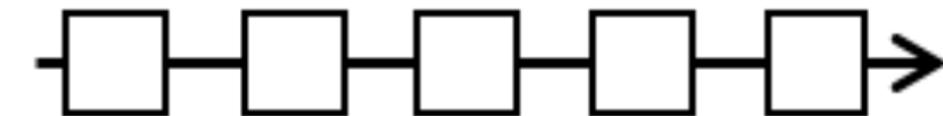
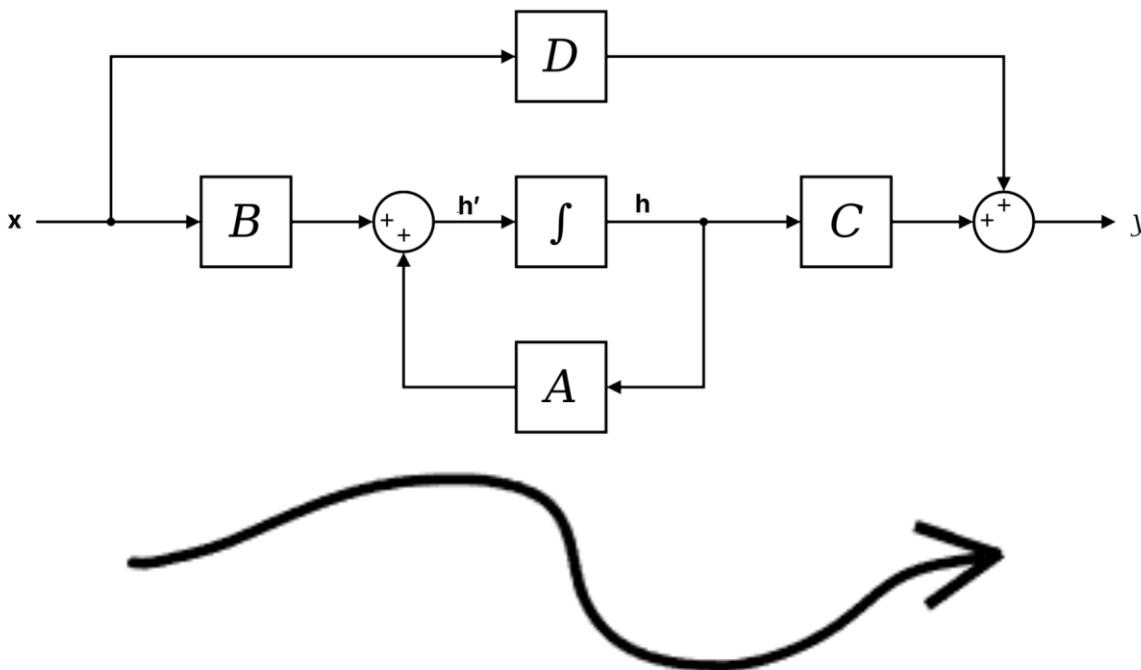
$$\begin{aligned} h'(t) &= \mathbf{A}h(t) + \mathbf{B}x(t) \\ y(t) &= \mathbf{C}h(t) + \mathbf{D}x(t) \end{aligned}$$

Discretize

Recurrent 😊

$$\begin{aligned} h_t &= \bar{\mathbf{A}}h_{t-1} + \bar{\mathbf{B}}x_t \\ y_t &= \mathbf{C}h_t + \mathbf{D}x_t \end{aligned}$$

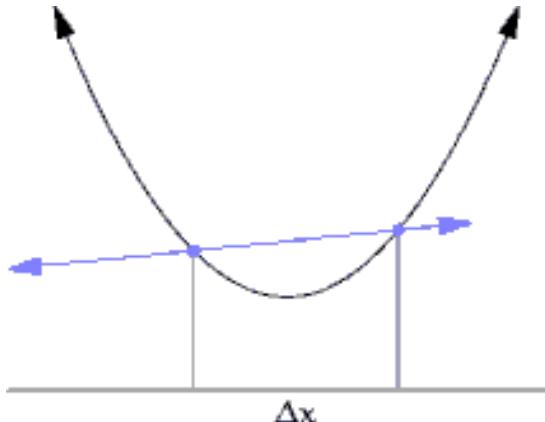
$$\bar{\mathbf{A}}, \bar{\mathbf{B}} = \text{Bilinear}(\mathbf{A}, \mathbf{B}, \Delta) \quad \text{hoặc} \quad \bar{\mathbf{A}}, \bar{\mathbf{B}} = \text{ZOH}(\mathbf{A}, \mathbf{B}, \Delta)$$



Discretize SSMs

1

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



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

when Δ is small

2

$$\begin{aligned} h'(t) &\approx \frac{h(t + \Delta) - h(t)}{\Delta} \\ h'(t)\Delta &\approx h(t + \Delta) - h(t) \\ h(t + \Delta) &\approx h'(t)\Delta + h(t). \end{aligned}$$

3

$$\begin{aligned} h_{t+\Delta} &\approx (\mathbf{A}h_t + \mathbf{B}x_t)\Delta + h_t \\ h_{t+\Delta} &\approx \Delta\mathbf{A}h_t + \Delta\mathbf{B}x_t + h_t \\ h_{t+\Delta} &\approx (\Delta\mathbf{A} + \mathbf{I})h_t + \Delta\mathbf{B}x_t \\ h_{t+\Delta} &\approx \bar{\mathbf{A}}h_t + \bar{\mathbf{B}}x_t, \end{aligned}$$

Example of Recurrent Representation

Recurrent Representation

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

$$y_t = \mathbf{C}h_t + \mathbf{D}x_t$$

$$\bar{\mathbf{A}}, \bar{\mathbf{B}} = \text{Bilinear}(\mathbf{A}, \mathbf{B}, \Delta) \quad \text{hoặc} \quad \bar{\mathbf{A}}, \bar{\mathbf{B}} = \text{ZOH}(\mathbf{A}, \mathbf{B}, \Delta)$$

$$(\text{ZOH}) \quad \bar{\mathbf{A}} = \exp(\Delta \mathbf{A})$$

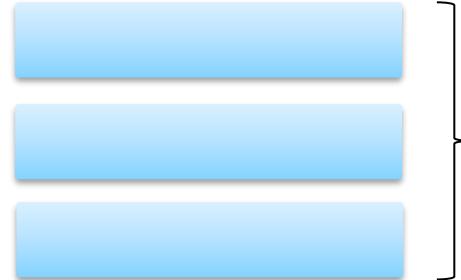
$$\bar{\mathbf{B}} = (\Delta \mathbf{A})^{-1} (\exp(\Delta \cdot \mathbf{A}) - \mathbf{I}) \cdot \Delta \mathbf{B}.$$

$$(\text{Bilinear}) \quad \bar{\mathbf{A}} = (\mathbf{I} - \Delta/2\mathbf{A})^{-1} (\mathbf{I} + \Delta/2\mathbf{A})$$

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

Example of Recurrent Representation

tôi
đi
học



input: $L \times D$

L : sequence len
 D : hidden state
 N : SSM dimension

$$h_t = \bar{\mathbf{A}}h_{t-1} + \bar{\mathbf{B}}x_t$$
$$y_t = \mathbf{C}h_t + \mathbf{D}x_t$$

$$h_0 = \underbrace{\bar{\mathbf{A}}}_{N \times N} \times \underbrace{h_{-1}}_{N \times D} + \underbrace{\bar{\mathbf{B}}}_{N \times 1} \times \underbrace{x}_{1 \times D} = \underbrace{h_0}_{N \times D}$$
$$y_0 = \underbrace{\mathbf{C}}_{1 \times N} \times \underbrace{h_0}_{N \times D} = \underbrace{y_0}_{1 \times D}$$

Example of Recurrent Representation

tôi	1	2	1	4
đi	5	1	7	3
học	4	3	2	4

} input: 3×4

L : sequence len
 D : hidden state
 N : SSM dimension

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

$$y_t = \mathbf{C}h_t + \mathbf{D}x_t$$

A	0.5	1	6
	1	2	3
	2	5	4

3×3

B	1	Δ	1
	5		
	3		

3×1

C	1	3	5
-----	---	---	---

1×3

(ZOH) $\bar{\mathbf{A}} = \exp(\Delta \mathbf{A})$
 $\bar{\mathbf{B}} = (\Delta \mathbf{A})^{-1}(\exp(\Delta \cdot \mathbf{A}) - \mathbf{I}) \cdot \Delta \mathbf{B}$.

(Bilinear) $\bar{\mathbf{A}} = (\mathbf{I} - \Delta/2\mathbf{A})^{-1}(\mathbf{I} + \Delta/2\mathbf{A})$
 $\bar{\mathbf{B}} = (\mathbf{I} - \Delta/2\mathbf{A})^{-1} \cdot \Delta \mathbf{B}$

$\bar{\mathbf{A}}$	1	2	1
	2	3	2
	3	4	1

3×3

$\bar{\mathbf{B}}$	2		
	3		
	1		

3×1

\mathbf{C}	1	3	5
--------------	---	---	---

1×3

Example of Recurrent Representation

tôi

1	2	1	4
---	---	---	---

đi

5	1	7	3
---	---	---	---

học

4	3	2	4
---	---	---	---

} input: 3×4

L : sequence len
 D : hidden state
 N : SSM dimension

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

$$y_t = \mathbf{C} h_t + \mathbf{D} x_t$$

$$h_0 = \underbrace{\begin{matrix} \bar{\mathbf{A}} \\ \hline 1 & 2 & 1 \\ 2 & 3 & 2 \\ 3 & 4 & 1 \end{matrix}}_{3 \times 3} \times \underbrace{\begin{matrix} h_{-1} \\ \hline 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{matrix}}_{3 \times 4}$$

$$+ \underbrace{\begin{matrix} \bar{\mathbf{B}} \\ \hline 2 \\ 3 \\ 1 \end{matrix}}_{3 \times 1} \times \underbrace{\begin{matrix} x_0 (\text{tôi}) \\ \hline 1 & 2 & 1 & 4 \end{matrix}}_{1 \times 4} = \underbrace{\begin{matrix} 2 & 4 & 2 & 8 \\ 3 & 6 & 3 & 12 \\ 1 & 2 & 1 & 4 \end{matrix}}_{3 \times 4}$$

$$y_0 = \underbrace{\begin{matrix} \mathbf{C} \\ \hline 1 & 3 & 5 \end{matrix}}_{1 \times 3} \times \underbrace{\begin{matrix} h_0 \\ \hline 2 & 4 & 2 & 8 \\ 3 & 6 & 3 & 12 \\ 1 & 2 & 1 & 4 \end{matrix}}_{3 \times 4} = \underbrace{\begin{matrix} 16 & 32 & 16 & 64 \end{matrix}}_{1 \times 4}$$

Example of Recurrent Representation

tôi $\begin{matrix} 1 & 2 & 1 & 4 \end{matrix}$
đi $\begin{matrix} 5 & 1 & 7 & 3 \end{matrix}$
học $\begin{matrix} 4 & 3 & 2 & 4 \end{matrix}$

} input: 3×4

L : sequence len
 D : hidden state
 N : SSM dimension

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

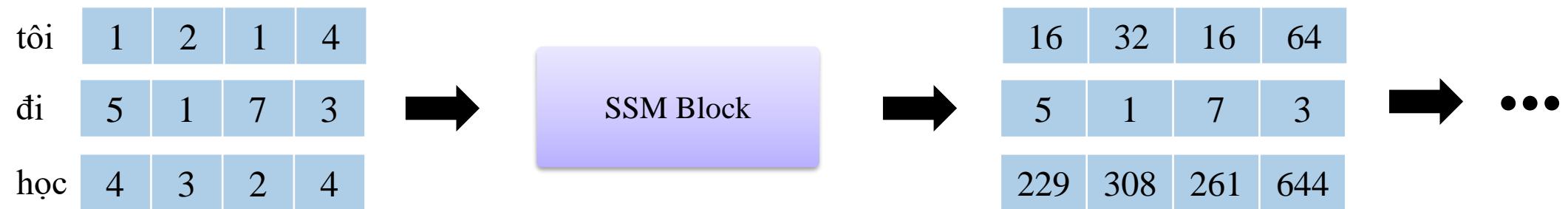
$$y_t = \mathbf{C}h_t + \mathbf{D}x_t$$

$$h_1 = \underbrace{\begin{matrix} 1 & 2 & 1 \\ 2 & 3 & 2 \\ 3 & 4 & 1 \end{matrix}}_{3 \times 3} \times \underbrace{\begin{matrix} 2 & 4 & 2 & 8 \\ 3 & 6 & 3 & 12 \\ 1 & 2 & 1 & 4 \end{matrix}}_{3 \times 4} \quad h_0$$

$$+ \underbrace{\begin{matrix} 2 \\ 3 \\ 1 \end{matrix}}_{3 \times 1} \times \underbrace{\begin{matrix} 5 & 1 & 7 & 3 \end{matrix}}_{1 \times 4} = \underbrace{\begin{matrix} 19 & 20 & 23 & 42 \\ 30 & 31 & 36 & 69 \\ 24 & 39 & 26 & 79 \end{matrix}}_{3 \times 4}$$

$$y_1 = \underbrace{\begin{matrix} 1 & 3 & 5 \end{matrix}}_{1 \times 3} \times \underbrace{\begin{matrix} 19 & 20 & 23 & 42 \\ 30 & 31 & 36 & 69 \\ 24 & 39 & 26 & 79 \end{matrix}}_{3 \times 4} = \underbrace{\begin{matrix} 229 & 308 & 261 & 644 \end{matrix}}_{1 \times 4}$$

Example of Recurrent Representation



Discretize SSMs

Recurrent Representation

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

$$y_t = \mathbf{C}h_t + \mathbf{D}x_t$$

$\bar{\mathbf{A}}, \bar{\mathbf{B}}$ = Bilinear($\mathbf{A}, \mathbf{B}, \Delta$) hoặc $\bar{\mathbf{A}}, \bar{\mathbf{B}}$ = ZOH($\mathbf{A}, \mathbf{B}, \Delta$)

(ZOH) $\bar{\mathbf{A}} = \exp(\Delta \mathbf{A})$

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

(Bilinear) $\bar{\mathbf{A}} = (\mathbf{I} - \Delta/2\mathbf{A})^{-1}(\mathbf{I} + \Delta/2\mathbf{A})$

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

Learnable Params

- A, B, C, D, Δ

Ideal Model

Parallelizable yes 😕

Training Stage $O(n)$ 😊

Inference Stage $O(n)$ 😊

Inference per Token $O(1)$ 😊

Convolutional Representation of SSMs

t=0

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

$$y_0 = \mathbf{C}h_0 = \mathbf{C}\bar{\mathbf{B}}x_0$$

t=1

$$h_1 = \bar{\mathbf{A}}h_0 + \bar{\mathbf{B}}x_1 = \bar{\mathbf{A}}\bar{\mathbf{B}}x_0 + \bar{\mathbf{B}}x_1$$

$$y_1 = \mathbf{C}h_1 = \mathbf{C}(\bar{\mathbf{A}}\bar{\mathbf{B}}x_0 + \bar{\mathbf{B}}x_1) = \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}x_0 + \mathbf{C}\bar{\mathbf{B}}x_1$$

t=2

$$h_2 = \bar{\mathbf{A}}h_1 + \bar{\mathbf{B}}x_2 = \bar{\mathbf{A}}(\bar{\mathbf{A}}\bar{\mathbf{B}}x_0 + \bar{\mathbf{B}}x_1) + \bar{\mathbf{B}}x_2 = \bar{\mathbf{A}}^2\bar{\mathbf{B}}x_0 + \bar{\mathbf{A}}\bar{\mathbf{B}}x_1 + \bar{\mathbf{B}}x_2$$

$$y_2 = \mathbf{C}h_2 = \mathbf{C}(\bar{\mathbf{A}}^2\bar{\mathbf{B}}x_0 + \bar{\mathbf{A}}\bar{\mathbf{B}}x_1 + \bar{\mathbf{B}}x_2) = \mathbf{C}\bar{\mathbf{A}}^2\bar{\mathbf{B}}x_0 + \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}x_1 + \mathbf{C}\bar{\mathbf{B}}x_2$$

...

...

t=k

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

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

$$y_t = \mathbf{C}h_t + \mathbf{D}x_t$$

Convolutional Representation of SSMs

t=k

$$y_k = \mathbf{C}\overline{\mathbf{A}}^k\overline{\mathbf{B}}x_0 + \mathbf{C}\overline{\mathbf{A}}^{k-1}\overline{\mathbf{B}}x_1 + \dots + \mathbf{C}\overline{\mathbf{A}}\overline{\mathbf{B}}x_{k-1} + \mathbf{C}\overline{\mathbf{B}}x_k.$$

$$\overline{\mathbf{K}} = (\mathbf{C}\overline{\mathbf{B}}, \mathbf{C}\overline{\mathbf{A}}\overline{\mathbf{B}}, \dots, \mathbf{C}\overline{\mathbf{A}}^{L-1}\overline{\mathbf{B}}),$$

$$x = (x_0, x_1, \dots, x_{L-1}),$$

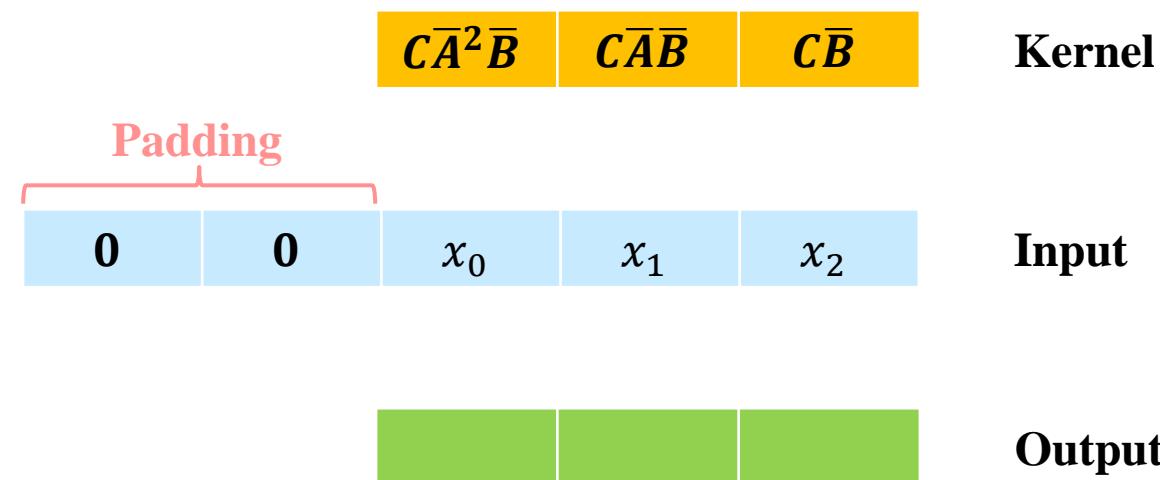
$$y = \overline{\mathbf{K}} * x.$$

Example of Convolutional Representation

tôi 
đi 
học 

} input: $L \times D$

L : sequence len
 D : hidden state
 N : SSM dimension

$$\bar{\mathbf{K}} = (\mathbf{C}\bar{\mathbf{B}}, \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}, \dots, \mathbf{C}\bar{\mathbf{A}}^{L-1}\bar{\mathbf{B}}),$$
$$x = (x_0, x_1, \dots, x_{L-1}),$$
$$y = \bar{\mathbf{K}} * x.$$


Example of Convolutional Representation

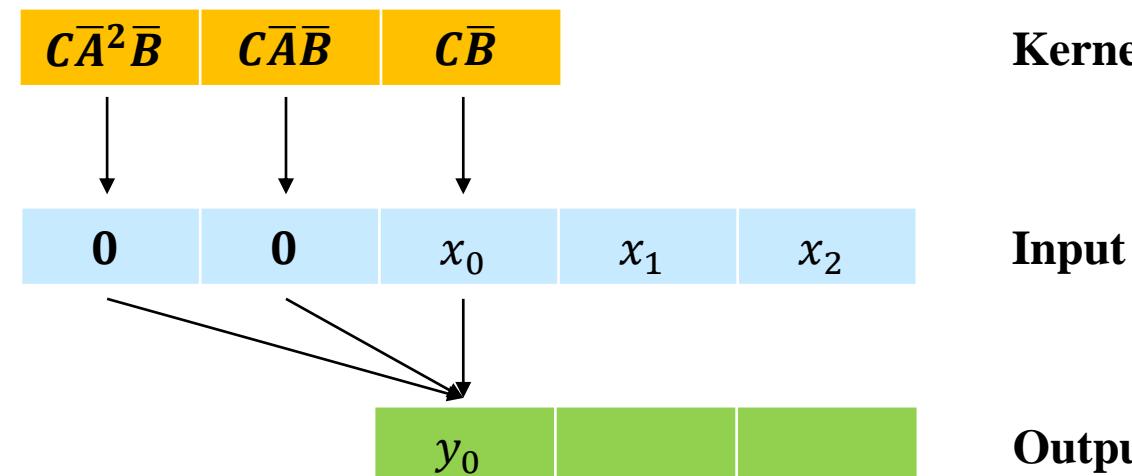
tôi 
đi 
học 

} input: $L \times D$

L : sequence len
 D : hidden state
 N : SSM dimension

$\bar{\mathbf{K}} = (\mathbf{C}\bar{\mathbf{B}}, \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}, \dots, \mathbf{C}\bar{\mathbf{A}}^{L-1}\bar{\mathbf{B}})$,
 $x = (x_0, x_1, \dots, x_{L-1})$,
 $y = \bar{\mathbf{K}} * x$.

Step 1



$$y_0 = \mathbf{C}\bar{\mathbf{B}}x_0$$

Example of Convolutional Representation

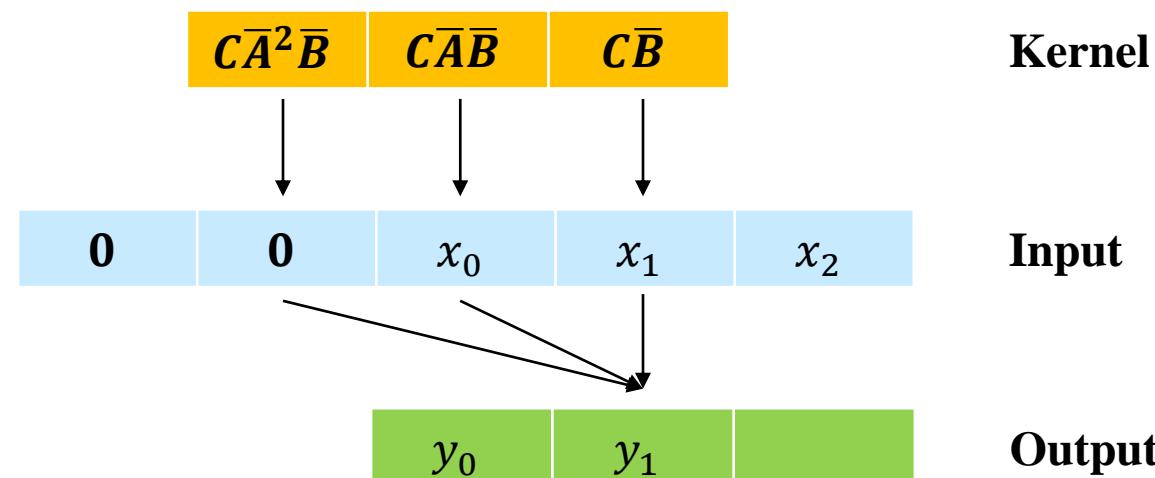
tôi 
đi 
học 

} input: $L \times D$

L : sequence len
 D : hidden state
 N : SSM dimension

$\bar{\mathbf{K}} = (\mathbf{C}\bar{\mathbf{B}}, \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}, \dots, \mathbf{C}\bar{\mathbf{A}}^{L-1}\bar{\mathbf{B}}),$
 $x = (x_0, x_1, \dots, x_{L-1}),$
 $y = \bar{\mathbf{K}} * x.$

Step 2



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

Example of Convolutional Representation

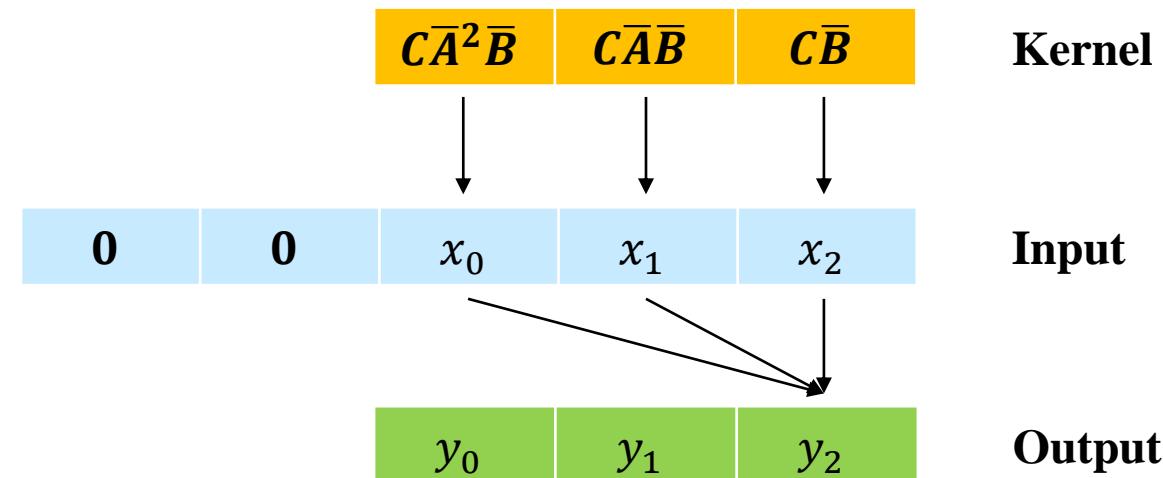
tôi 
đi 
học 

} input: $L \times D$

L : sequence len
 D : hidden state
 N : SSM dimension

$\bar{\mathbf{K}} = (\mathbf{C}\bar{\mathbf{B}}, \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}, \dots, \mathbf{C}\bar{\mathbf{A}}^{L-1}\bar{\mathbf{B}}),$
 $x = (x_0, x_1, \dots, x_{L-1}),$
 $y = \bar{\mathbf{K}} * x.$

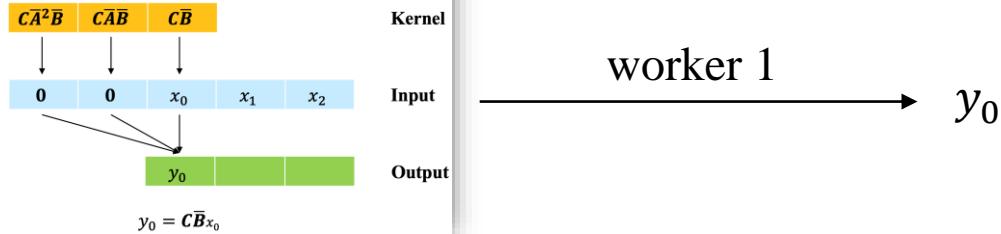
Step 3



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

Example of Convolutional Representation

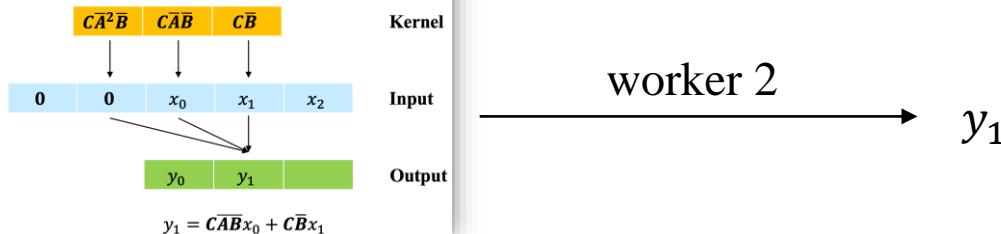
Step 1



worker 1

 y_0

Step 2

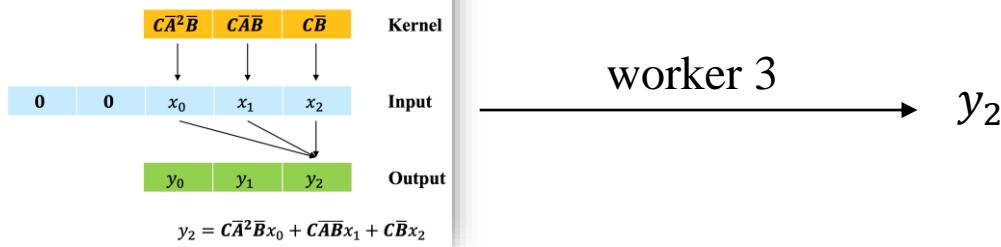


worker 2

 y_1

Can be computed in parallel!

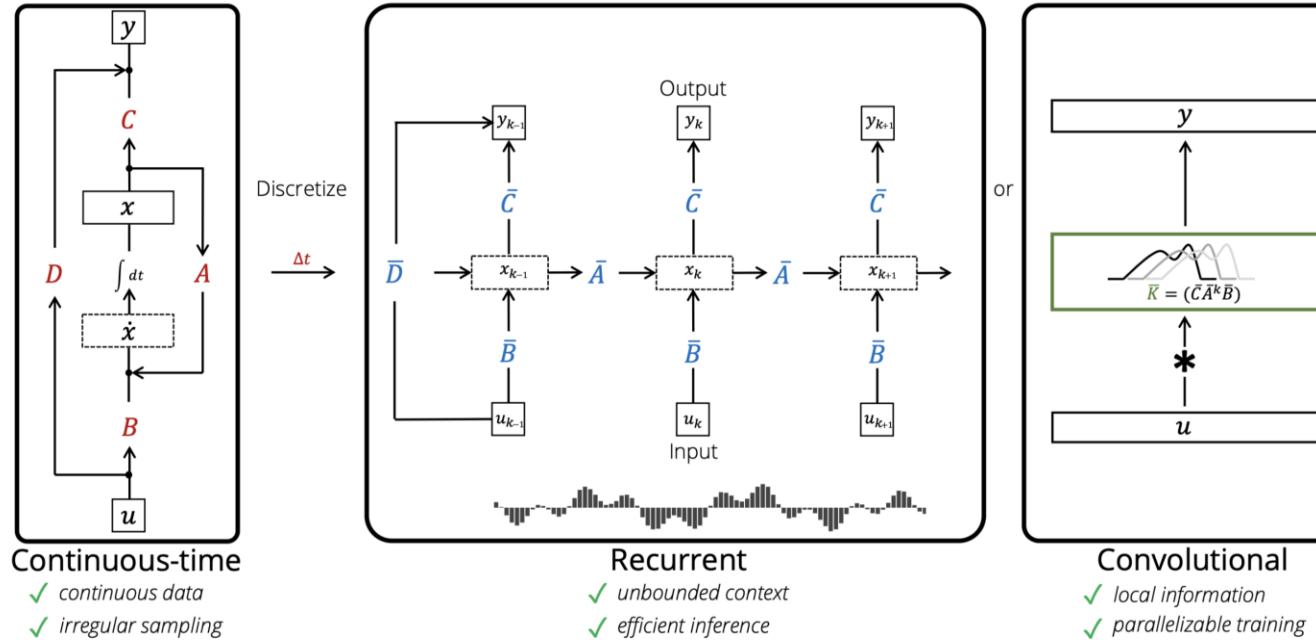
Step 3



worker 3

 y_2

SSMs is Ideal Models



Recurrent

$$\begin{aligned} h_t &= \bar{\mathbf{A}}h_{t-1} + \bar{\mathbf{B}}x_t \\ y_t &= \mathbf{C}h_t + \mathbf{D}x_t \end{aligned}$$

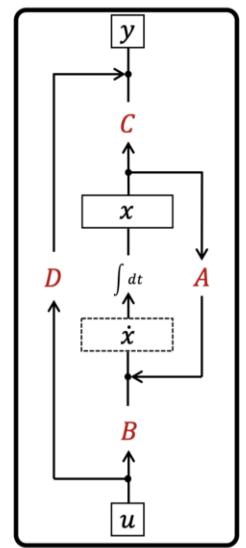
Convolutional

$$\begin{aligned} \bar{\mathbf{K}} &= (\mathbf{CB}, \mathbf{CAB}, \dots, \mathbf{CA}^{L-1}\mathbf{B}), \\ x &= (x_0, x_1, \dots, x_{L-1}), \\ y &= \bar{\mathbf{K}} * x. \end{aligned}$$

Ideal Model

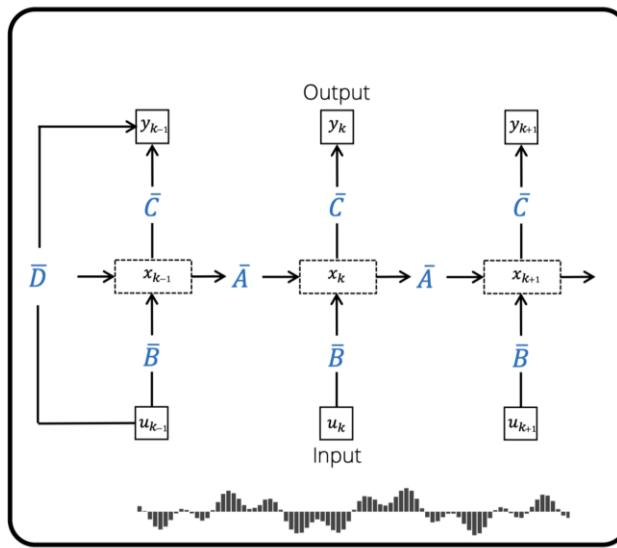
Parallelizable	yes	😊
Training Stage	$O(n)$	😊
Inference Stage	$O(n)$	😊
Inference per Token	$O(1)$	😊

SSMs is Ideal Models

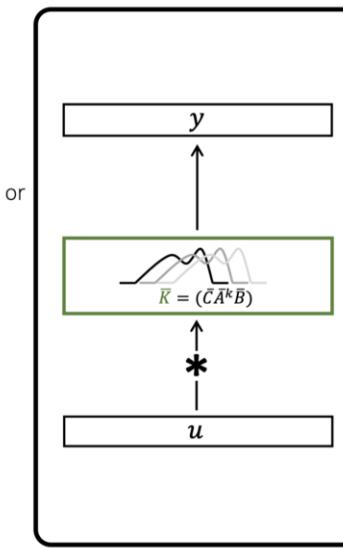


- ✓ continuous data
- ✓ irregular sampling

Discretize
 Δt



- ✓ unbounded context
- ✓ efficient inference



- ✓ local information
- ✓ parallelizable training

Recurrent

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

$$y_t = \mathbf{C}h_t + \mathbf{D}x_t$$

Convolutional

$$\bar{\mathbf{K}} = (\mathbf{CB}, \mathbf{CAB}, \dots, \mathbf{CA}^{L-1}\bar{\mathbf{B}}),$$

$$x = (x_0, x_1, \dots, x_{L-1}),$$

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

The importance of the A matrix



SSMs still have some of the same issues as RNNs, like vanishing/exploding gradients.

↓ solution

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

↓



The performance in the Sequential MNIST increased from 60% to 98%.

Recurrent

$$\begin{aligned} h_t &= \bar{\mathbf{A}}h_{t-1} + \bar{\mathbf{B}}x_t \\ y_t &= \mathbf{C}h_t + \mathbf{D}x_t \end{aligned}$$

Convolutional

$$\begin{aligned} \bar{\mathbf{K}} &= (\mathbf{CB}, \mathbf{CAB}, \dots, \mathbf{CA}^{L-1}\bar{\mathbf{B}}), \\ x &= (x_0, x_1, \dots, x_{L-1}), \end{aligned}$$

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

QUIZ!

TIME

Outline

1. Motivation

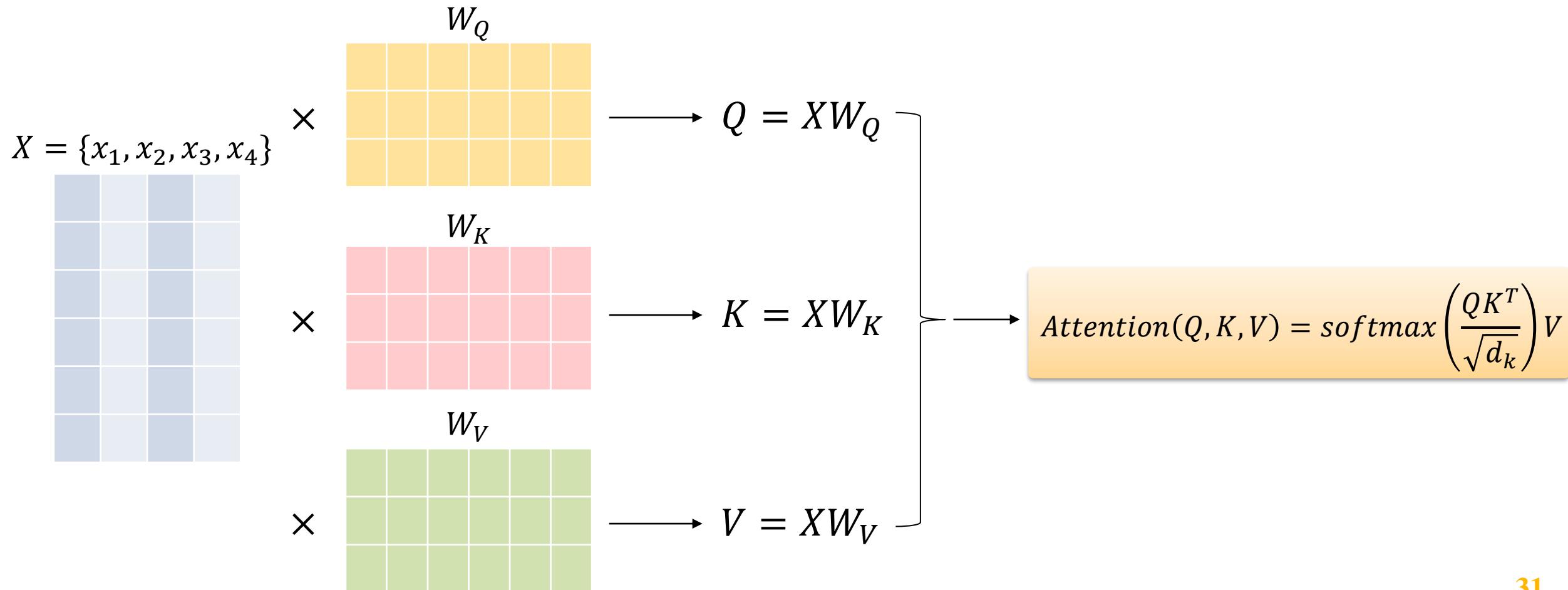
2. State Space Models

3. Mamba

4. Coding

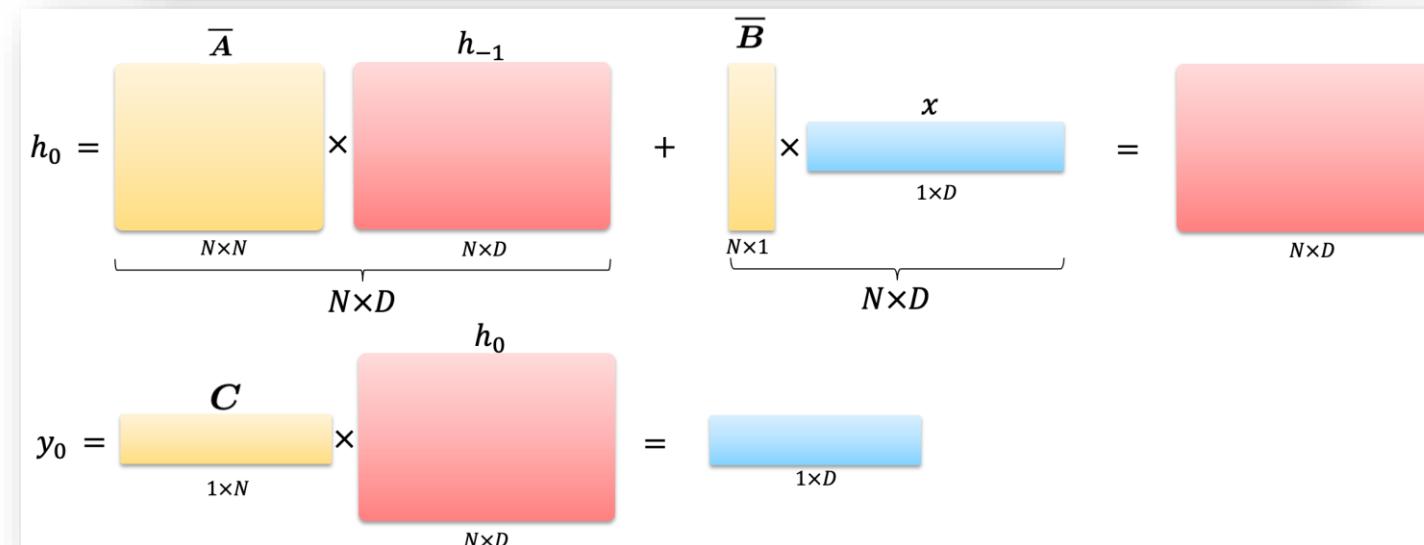
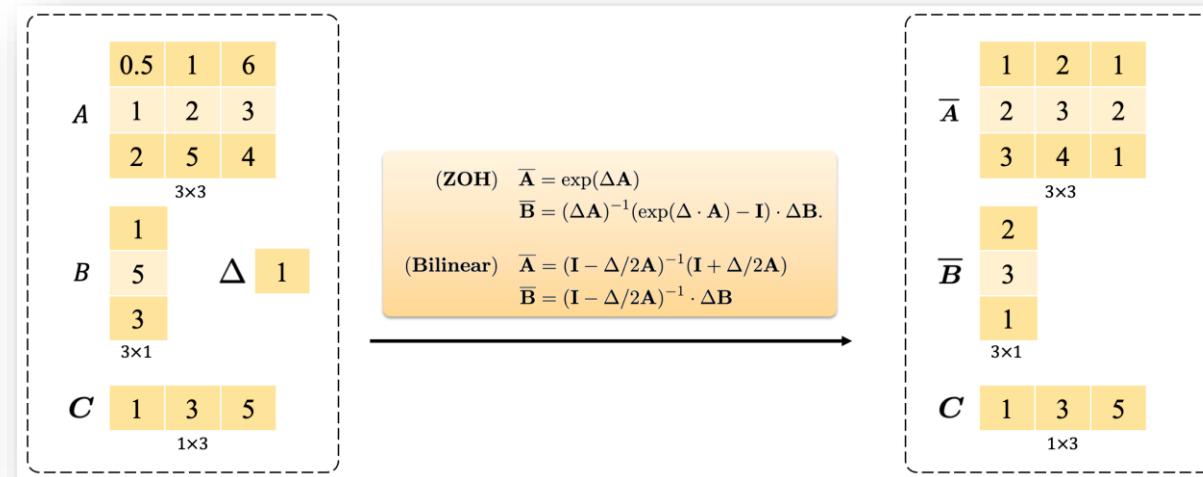
Mamba Motivation

A step back: Attention in Transformer



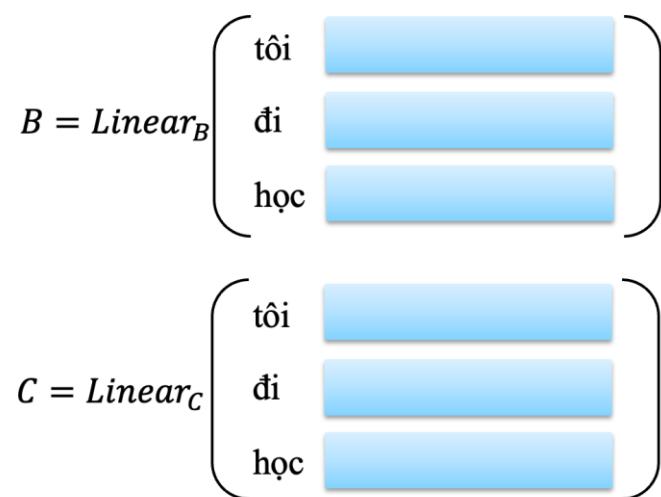
Mamba Motivation

A step back: S4



Mamba Contribution #1: Selection Mechanism

	0.5	1	6
A	1	2	3
	2	5	4



$$\Delta = \text{Softplus}(\Delta + \text{Linear}_\Delta(\Delta))$$

(ZOH) $\bar{\mathbf{A}} = \exp(\Delta \mathbf{A})$
 $\bar{\mathbf{B}} = (\Delta \mathbf{A})^{-1}(\exp(\Delta \cdot \mathbf{A}) - \mathbf{I}) \cdot \Delta \mathbf{B}$.

(Bilinear) $\bar{\mathbf{A}} = (\mathbf{I} - \Delta/2\mathbf{A})^{-1}(\mathbf{I} + \Delta/2\mathbf{A})$
 $\bar{\mathbf{B}} = (\mathbf{I} - \Delta/2\mathbf{A})^{-1} \cdot \Delta \mathbf{B}$

$\bar{\mathbf{A}}$	1	2	1
	2	3	2
	3	4	1

3×3

$\bar{\mathbf{B}}$	2		
	3		
	1		

3×1

$\bar{\mathbf{C}}$	1	3	5
--------------------	---	---	---

1×3

Mamba Contribution #1: Selection Mechanism

Compare S4 and Mamba (S6)

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

Mamba Contribution #1: Selection Mechanism

S4

A	
	3×3
B	
	3×1
C	
	1×3

(ZOH) $\bar{A} = \exp(\Delta A)$
 $\bar{B} = (\Delta A)^{-1}(\exp(\Delta \cdot A) - I) \cdot \Delta B.$

(Bilinear) $\bar{A} = (I - \Delta/2A)^{-1}(I + \Delta/2A)$
 $\bar{B} = (I - \Delta/2A)^{-1} \cdot \Delta B$

\bar{A}	
	3×3
\bar{B}	
	3×1
C	
	1×3

Mamba

A	
	3×3
$B = Linear_B$	

$C = Linear_C$	
----------------	--

$\Delta = Softplus(\Delta + Linear_{\Delta}(\Delta))$

(ZOH) $\bar{A} = \exp(\Delta A)$
 $\bar{B} = (\Delta A)^{-1}(\exp(\Delta \cdot A) - I) \cdot \Delta B.$

(Bilinear) $\bar{A} = (I - \Delta/2A)^{-1}(I + \Delta/2A)$
 $\bar{B} = (I - \Delta/2A)^{-1} \cdot \Delta B$

\bar{A}	
	3×3
\bar{B}	
	3×1
C	
	1×3

Mamba Contribution #1: Selection Mechanism

Mamba becomes time-varying

t=0

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

$$y_0 = \mathbf{C}h_0 = \mathbf{C}\bar{\mathbf{B}}x_0$$

t=1

$$h_1 = \bar{\mathbf{A}}h_0 + \bar{\mathbf{B}}x_1 = \bar{\mathbf{A}}\bar{\mathbf{B}}x_0 + \bar{\mathbf{B}}x_1$$

$$y_1 = \mathbf{C}h_1 = \mathbf{C}(\bar{\mathbf{A}}\bar{\mathbf{B}}x_0 + \bar{\mathbf{B}}x_1) = \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}x_0 + \mathbf{C}\bar{\mathbf{B}}x_1$$

t=2

$$h_2 = \bar{\mathbf{A}}h_1 + \bar{\mathbf{B}}x_2 = \bar{\mathbf{A}}(\bar{\mathbf{A}}\bar{\mathbf{B}}x_0 + \bar{\mathbf{B}}x_1) + \bar{\mathbf{B}}x_2 = \bar{\mathbf{A}}^2\bar{\mathbf{B}}x_0 + \bar{\mathbf{A}}\bar{\mathbf{B}}x_1 + \bar{\mathbf{B}}x_2$$

$$y_2 = \mathbf{C}h_2 = \mathbf{C}(\bar{\mathbf{A}}^2\bar{\mathbf{B}}x_0 + \bar{\mathbf{A}}\bar{\mathbf{B}}x_1 + \bar{\mathbf{B}}x_2) = \mathbf{C}\bar{\mathbf{A}}^2\bar{\mathbf{B}}x_0 + \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}x_1 + \mathbf{C}\bar{\mathbf{B}}x_2$$

...

...

t=k

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

Convolutional



$$\bar{\mathbf{K}} = (\mathbf{C}\bar{\mathbf{B}}, \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}, \dots, \mathbf{C}\bar{\mathbf{A}}^{L-1}\bar{\mathbf{B}}),$$

$$x = (x_0, x_1, \dots, x_{L-1}),$$

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

0.5	1	6
1	2	3
2	5	4

$$B = \text{Linear}_B \begin{cases} \text{tôi} \\ \text{di} \\ \text{hoc} \end{cases}$$

$$C = \text{Linear}_C \begin{cases} \text{tôi} \\ \text{di} \\ \text{hoc} \end{cases}$$

$$\Delta = \text{Softplus}(\Delta + \text{Linear}_{\Delta}(\Delta))$$

$$(\text{ZOH}) \quad \bar{\mathbf{A}} = \exp(\Delta \mathbf{A})$$

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

$$(\text{Bilinear}) \quad \bar{\mathbf{A}} = (\mathbf{I} - \Delta/2\mathbf{A})^{-1}(\mathbf{I} + \Delta/2\mathbf{A})$$

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

1	2	1
2	3	2
3	4	1

3x3

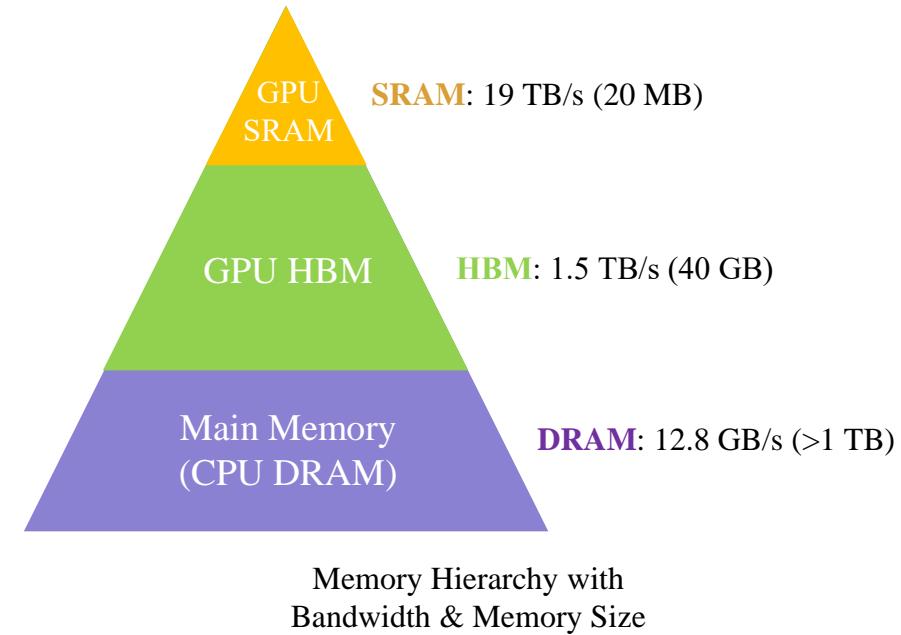
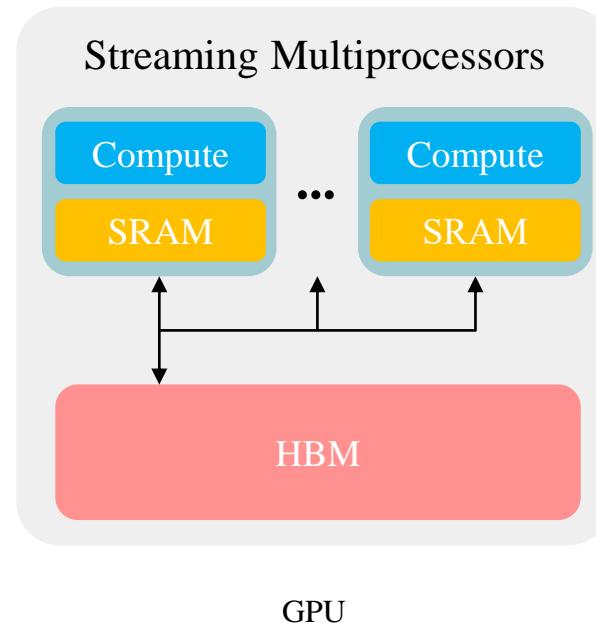
2
3
1

3x1

1	3	5
---	---	---

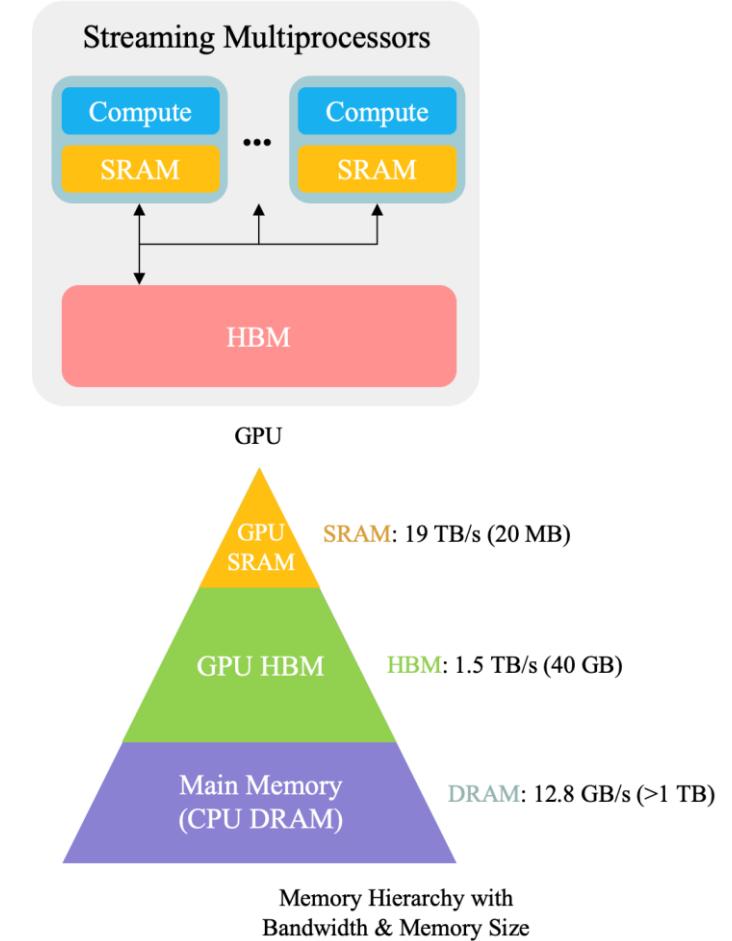
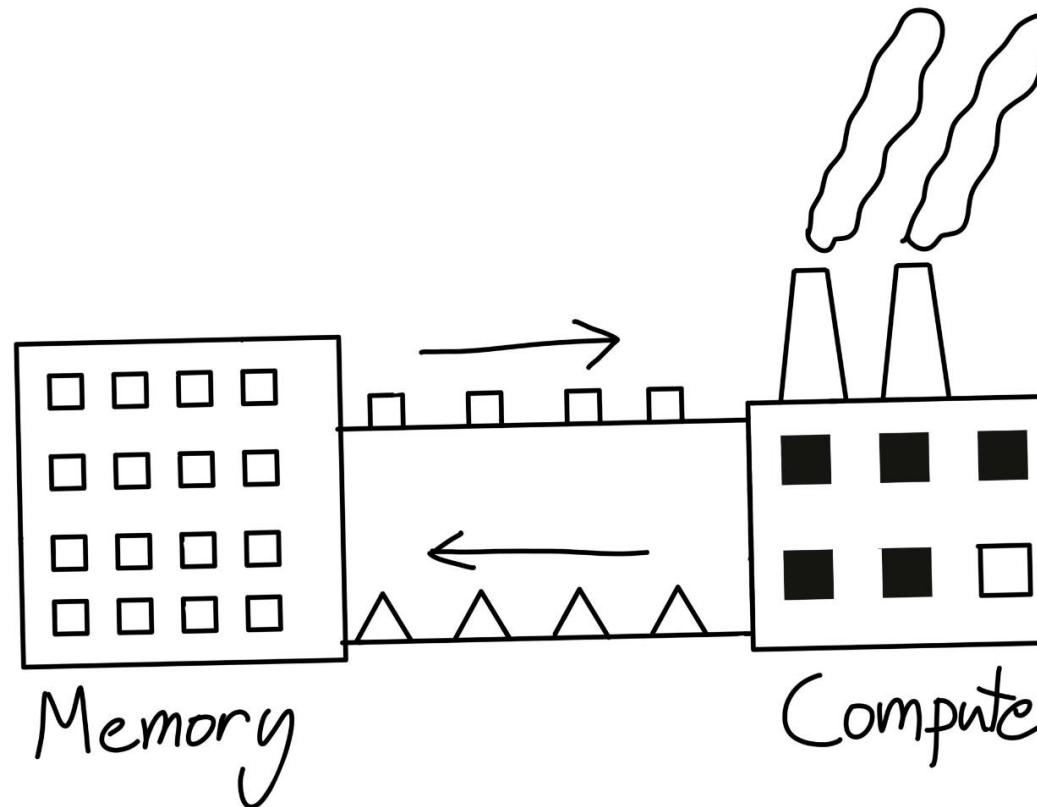
1x3

How does GPU work?



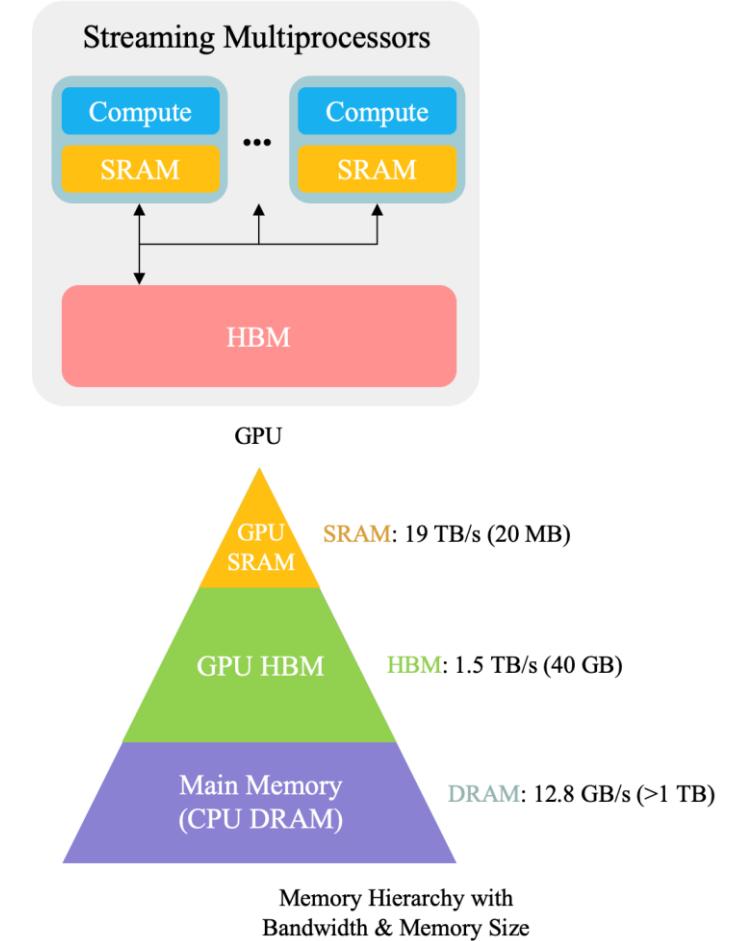
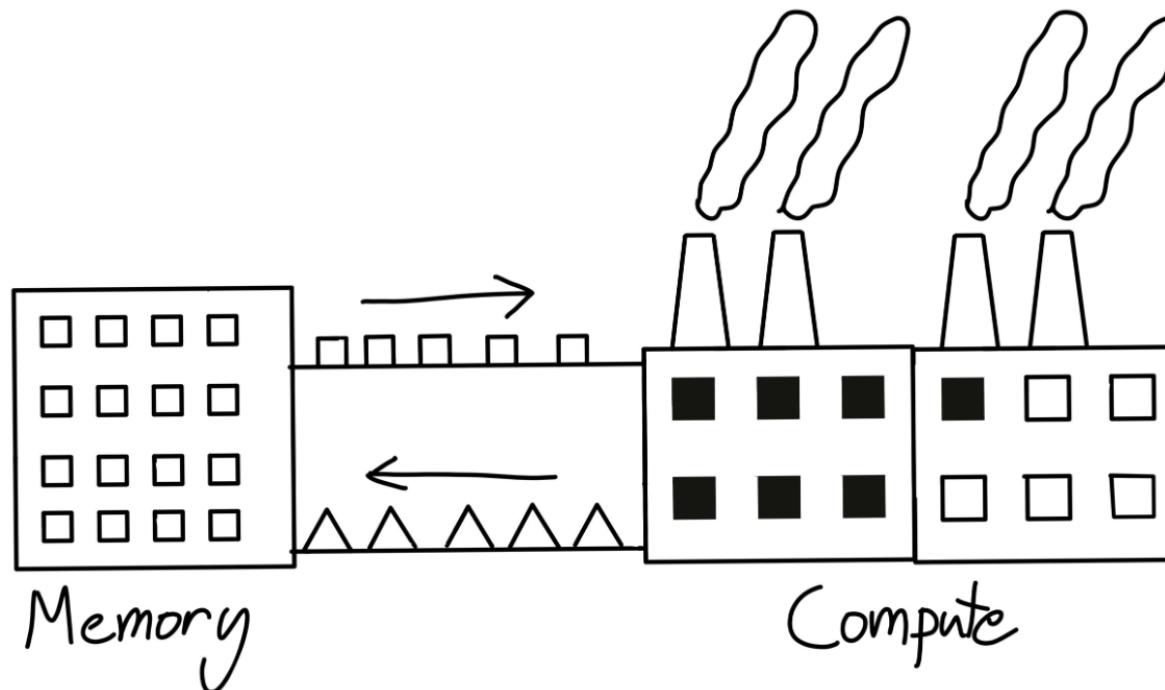
How does GPU work?

Memory & Compute



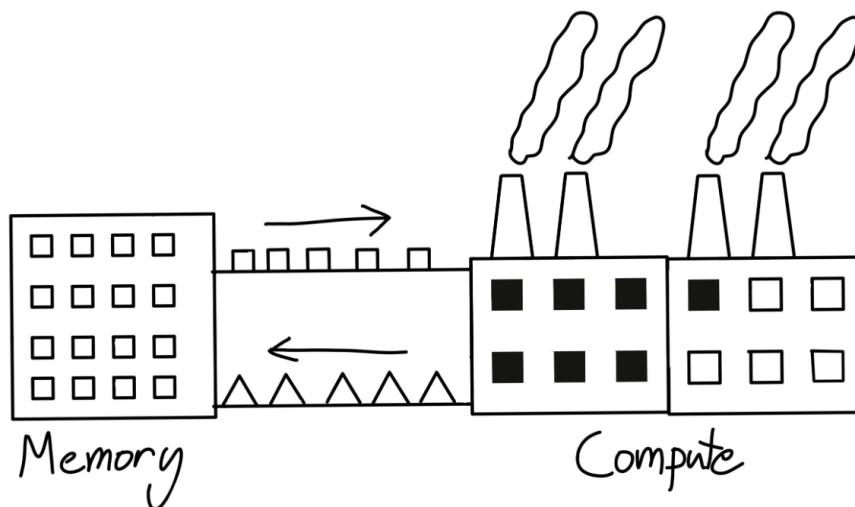
How does GPU work?

Memory & Compute



How does GPU work?

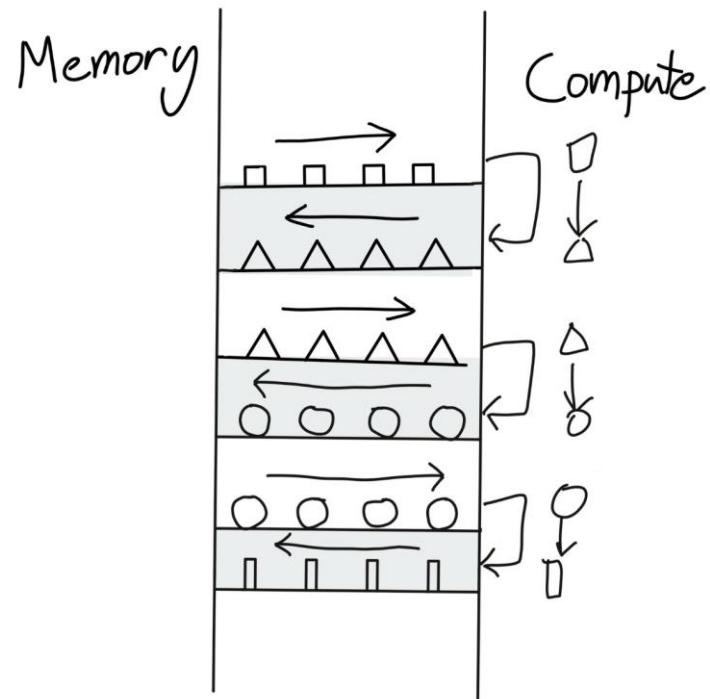
Memory & Compute



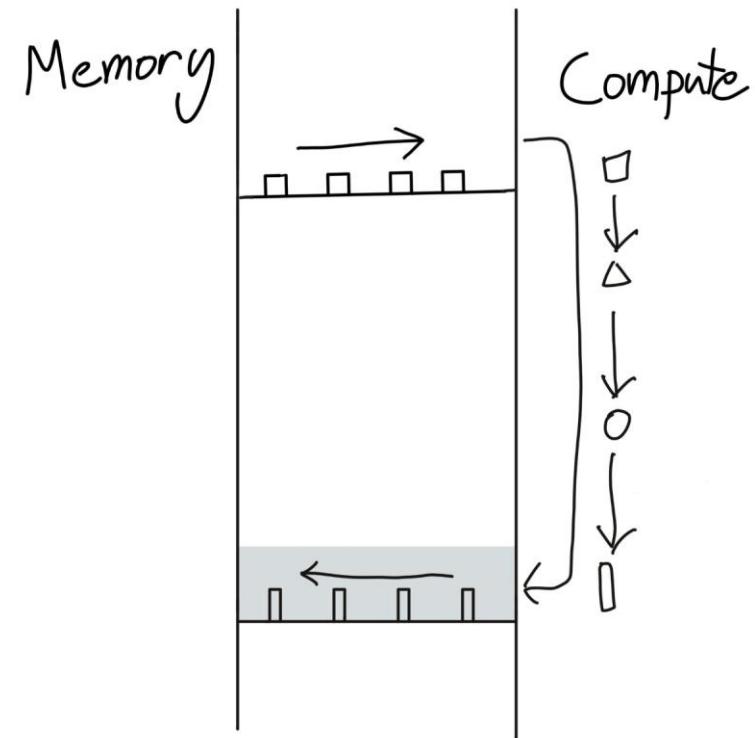
	A100 80GB PCIe	A100 80GB SXM
FP64		9.7 TFLOPS
FP64 Tensor Core		19.5 TFLOPS
FP32		19.5 TFLOPS
Tensor Float 32 (TF32)		156 TFLOPS 312 TFLOPS*
BFLOAT16 Tensor Core		312 TFLOPS 624 TFLOPS*
FP16 Tensor Core		312 TFLOPS 624 TFLOPS*
INT8 Tensor Core		624 TOPS 1248 TOPS*
GPU Memory	80GB HBM2e	80GB HBM2e
GPU Memory Bandwidth	1,935 GB/s	2,039 GB/s
Max Thermal Design Power (TDP)	300W	400W ***

Mamba Contribution #2: Hardware-aware Algorithm

Kernel Fusion



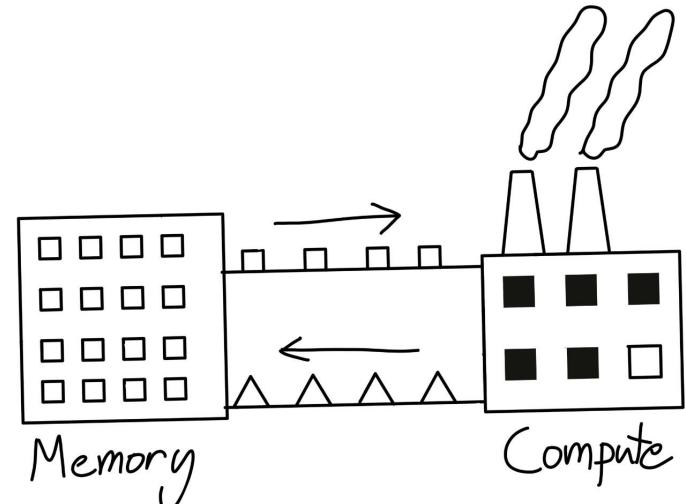
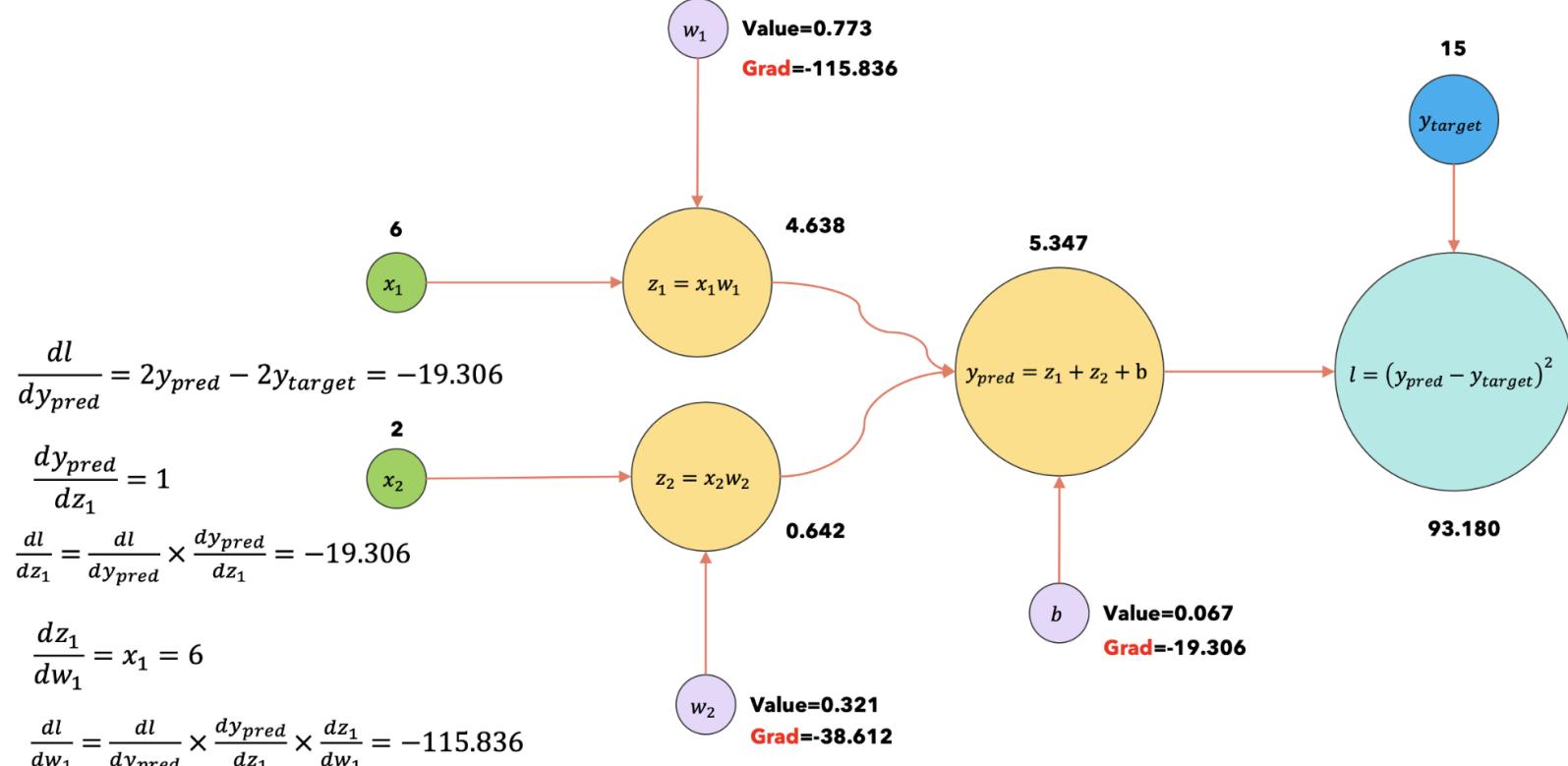
```
1 x1 = x.cos() # Read from x in global memory, write to x1
2 x2 = x1.cos() # Read from x1 in global memory, write to x2
3 x3 = x2.cos() # Read from x2 in global memory, write to x3
```



```
1 x3 = x.cos().cos().cos() # Read from x in global memory, write to x3
```

Mamba Contribution #2: Hardware-aware Algorithm

Recomputation



Mamba Contribution #2: Hardware-aware Algorithm

Parallel Scan

Prefix-sum

Initial array

9	6	7	10	8	7
---	---	---	----	---	---

Pre-fix sum

9	→	15	22	32	40	47
---	---	----	----	----	----	----

Parallel Scan

Model input

x_0	x_1	x_2	x_3	x_4	x_5
-------	-------	-------	-------	-------	-------

Scan output

h_0	→	h_1	h_2	h_3	h_4	h_5
-------	---	-------	-------	-------	-------	-------

Recurrent

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

$$y_t = \mathbf{C}h_t + \mathbf{D}x_t$$

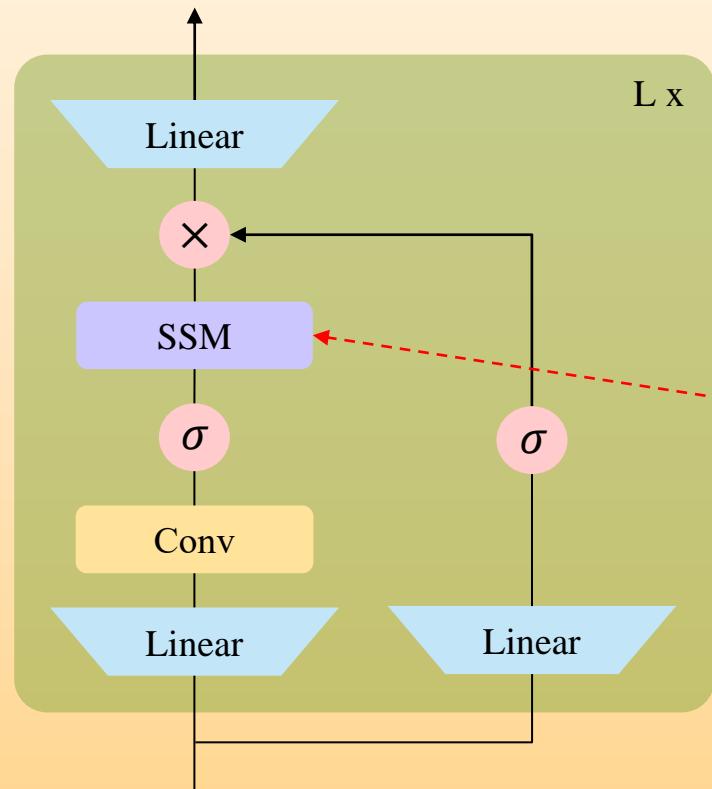
Mamba Contribution #2: Hardware-aware Algorithm

Parallel Scan

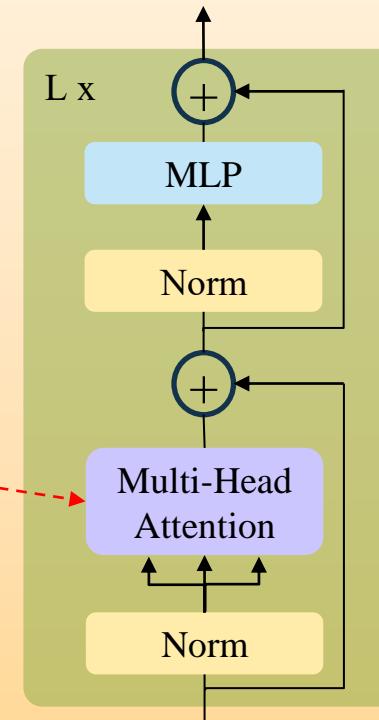
Position	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	
Sequential Prefix-Sum (what we expect as output from the parallel computation)	1	5	11	14	24	26	33	34	41	50	58	66	76	85	91	101	
Array of values	1	4	6	3	10	2	7	1	7	9	8	8	10	9	6	10	
		5		9		12		8		16		16		19		16	
			14					20				32				35	
								34								67	
											66					101	
						26				50				85			
					11		24		33		41		58		76		91
Parallel Prefix-Sum (output of parallel computation)	1	5	11	14	24	26	33	34	41	50	58	66	76	85	91	101	
Do the two computations match?	YES																

Mamba Block

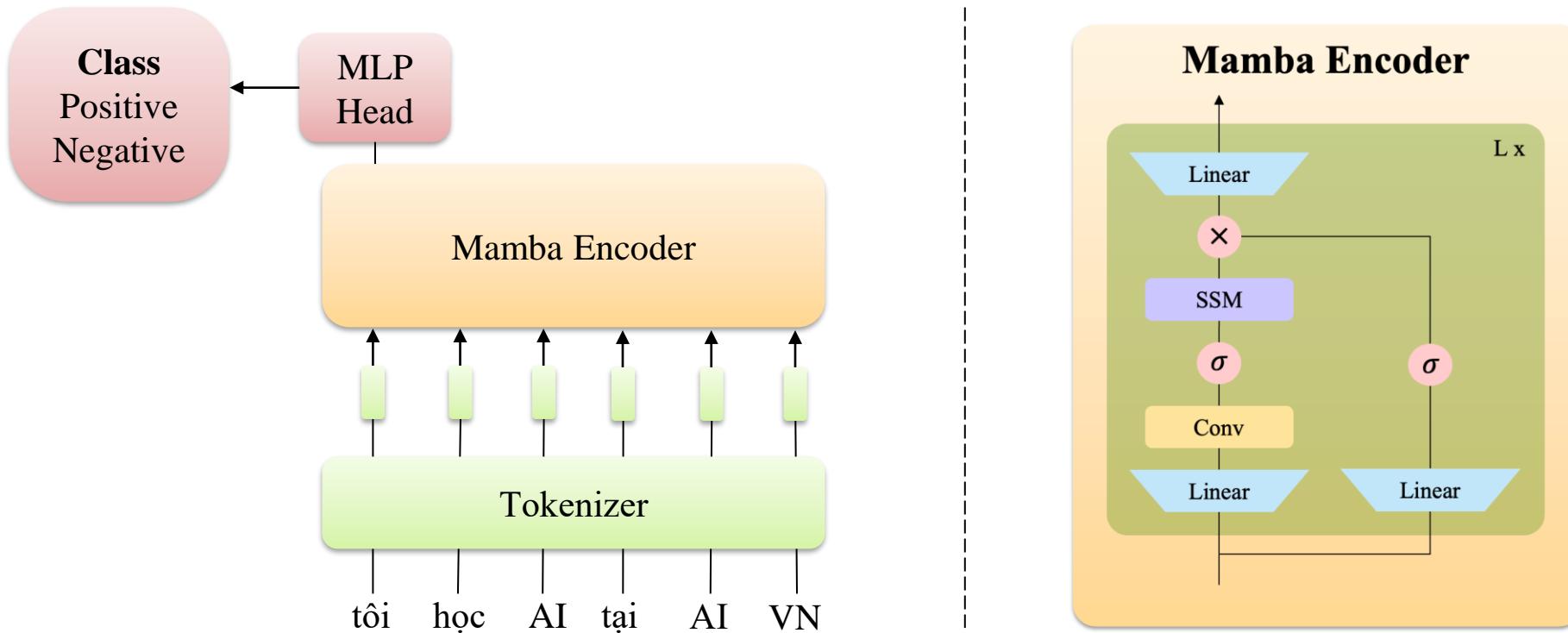
Mamba Encoder



Transformer Encoder

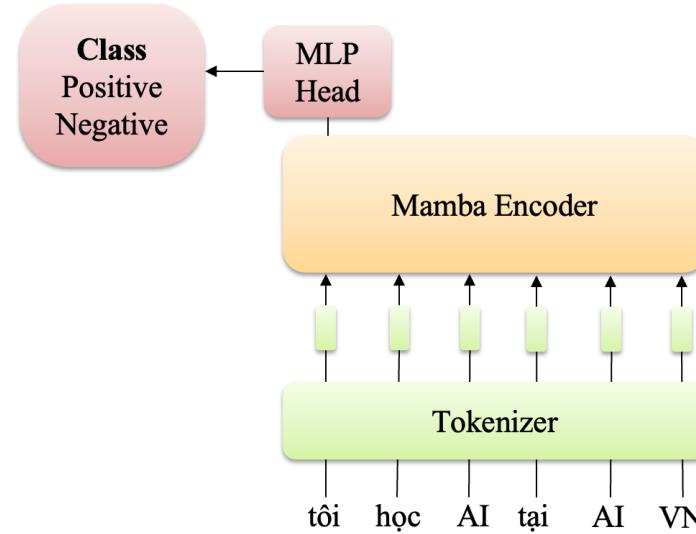
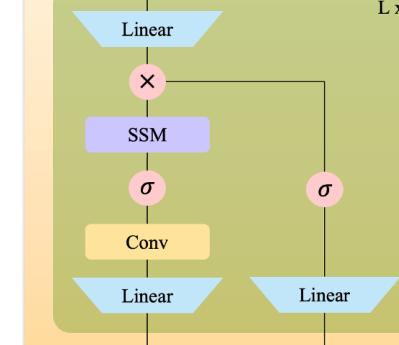


Mamba Architecture

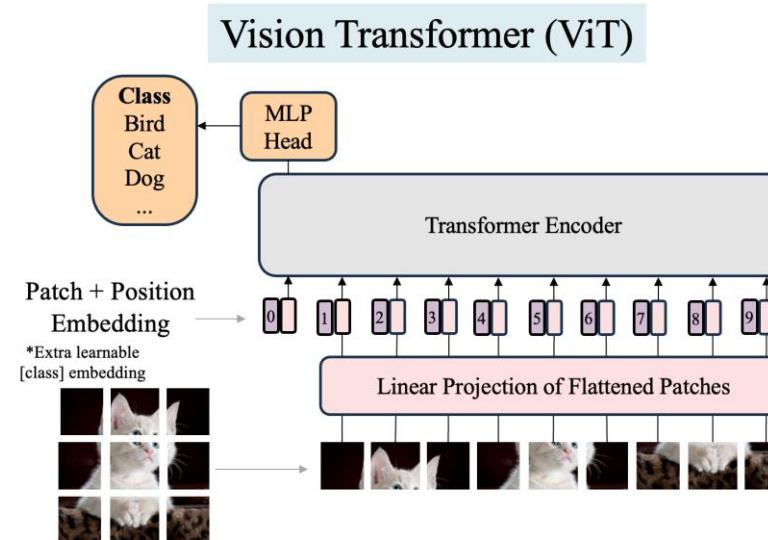


Mamba Arch. & Transformer Arch.

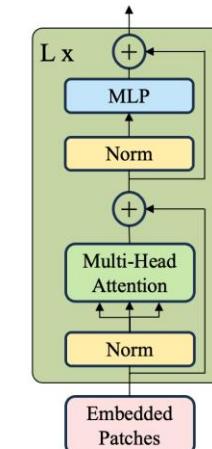
Mamba

**Mamba Encoder**

ViT



Transformer Encoder



Outline

1. Motivation

2. State Space Models

3. Mamba

4. Coding

Custom Trainer
(MambaTrainer)**Mamba Model**
(MambaTextClassification)**Mamba Backbone****Mamba Head**
(MambaClassificationHead)

Custom Trainer (MambaTrainer)

Mamba Model (MambaTextClassification)

Mamba Backbone

Mamba Head (MambaClassificationHead)

```
1 # Định nghĩa lớp head để phân loại
2 class MambaClassificationHead(nn.Module):
3     def __init__(self, d_model, num_classes, **kwargs):
4         super(MambaClassificationHead, self).__init__()
5         # Sử dụng một lớp tuyến tính để thực hiện phân loại dựa trên đầu vào
6         # có kích thước d_model và num_classes cần phân loại.
7         self.classification_head = nn.Linear(d_model, num_classes, **kwargs)
8
9     def forward(self, hidden_states):
10        return self.classification_head(hidden_states)
```

Custom Trainer (MambaTrainer)

Mamba Model (MambaTextClassification)

Mamba Backbone

Mamba Head (MambaClassificationHead)

```
1 class MambaTextClassification(MambaLMHeadModel):
2     def __init__(self,
3                  config: MambaConfig,
4                  initializer_cfg=None,
5                  device=None,
6                  dtype=None,
7                  ) -> None:
8         super().__init__(config, initializer_cfg, device, dtype)
9
10        # Tạo một đầu phân loại sử dụng MambaClassificationHead với kích thước đầu vào là d_model và số lớp là 2.
11        self.classification_head = MambaClassificationHead(d_model=config.d_model, num_classes=2)
12
13        del self.lm_head
14
15    def forward(self, input_ids, attention_mask=None, labels=None):
16        # Truyền input_ids qua mô hình gốc để nhận hidden_states.
17        hidden_states = self.backbone(input_ids)
18
19        # Lấy trung bình của hidden_states theo chiều thứ 2 để tạo ra [CLS] feature đại diện
20        mean_hidden_states = hidden_states.mean(dim=1)
21
22        # Đưa mean_hidden_states qua đầu phân loại để nhận logits.
23        logits = self.classification_head(mean_hidden_states)
24
25        if labels is None:
26            ClassificationOutput = namedtuple("ClassificationOutput", ["logits"])
27            return ClassificationOutput(logits=logits)
28        else:
29            ClassificationOutput = namedtuple("ClassificationOutput", ["loss", "logits"])
30
31
32        # Sử dụng hàm mất mát CrossEntropyLoss để tính loss.
33        loss_fct = nn.CrossEntropyLoss()
34        loss = loss_fct(logits, labels)
35
36        return ClassificationOutput(loss=loss, logits=logits)
```

Custom Trainer (MambaTrainer)

Mamba Model (MambaTextClassification)

Mamba Backbone

Mamba Head (MambaClassificationHead)

```
1 class MambaTextClassification(MambaLMHeadModel):
2     ...
3     def predict(self, text, tokenizer, id2label=None):
4         input_ids = torch.tensor(tokenizer(text)['input_ids'], device='cuda')[None]
5         with torch.no_grad():
6             logits = self.forward(input_ids).logits[0]
7             label = np.argmax(logits.cpu().numpy())
8
9         if id2label is not None:
10             return id2label[label]
11         else:
12             return label
13
14     @classmethod
15     def from_pretrained(cls, pretrained_model_name, device=None, dtype=None, **kwargs):
16         # Tải pretrained từ mô hình đã được huấn luyện trước đó.
17         config_data = load_config_hf(pretrained_model_name)
18         config = MambaConfig(**config_data)
19
20         # Khởi tạo mô hình từ cấu hình và chuyển nó đến thiết bị và kiểu dữ liệu mong muốn.
21         model = cls(config, device=device, dtype=dtype, **kwargs)
22
23         # Tải trạng thái mô hình đã được huấn luyện trước đó.
24         model_state_dict = load_state_dict_hf(pretrained_model_name, device=device, dtype=dtype)
25         model.load_state_dict(model_state_dict, strict=False)
26
27         # In ra các tham số embedding mới được khởi tạo.
28         print("Newly initialized embedding:", set(model.state_dict().keys()) - set(model_state_dict.keys()))
29         return model
```

Custom Trainer (MambaTrainer)

Mamba Model (MambaTextClassification)

Mamba Backbone

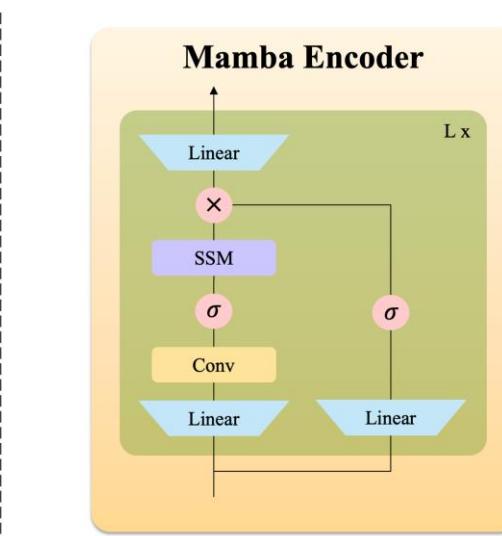
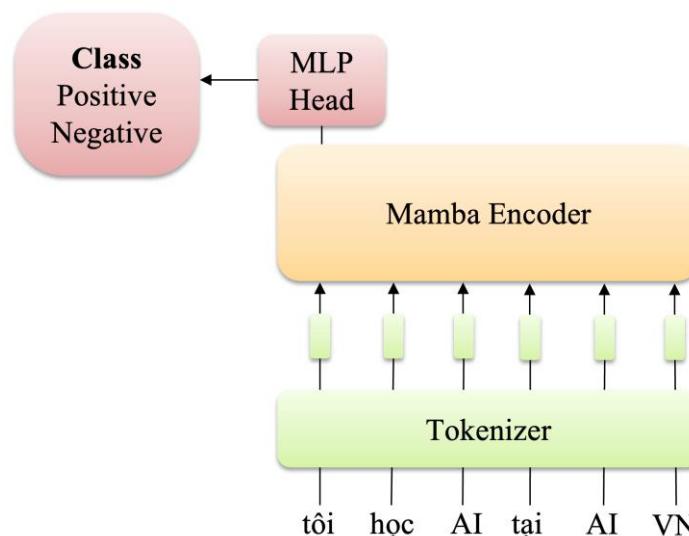
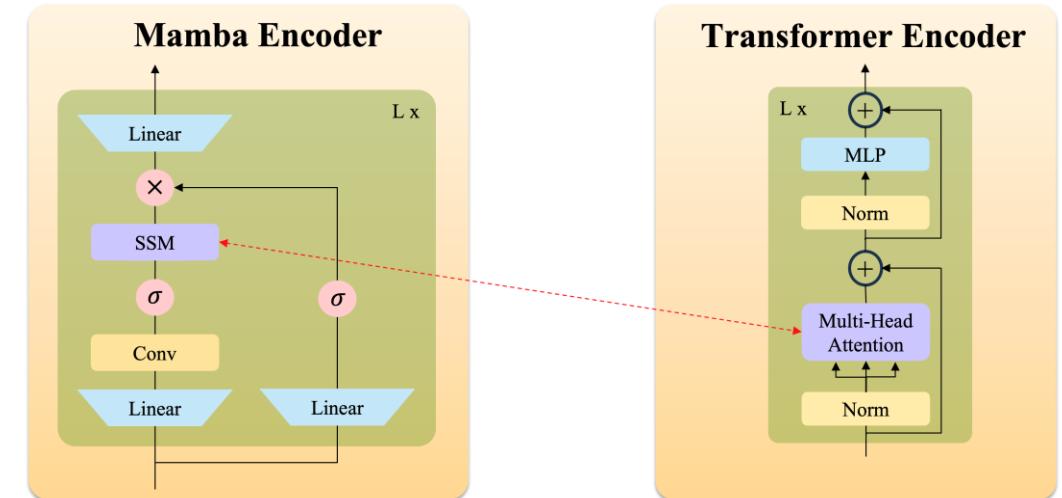
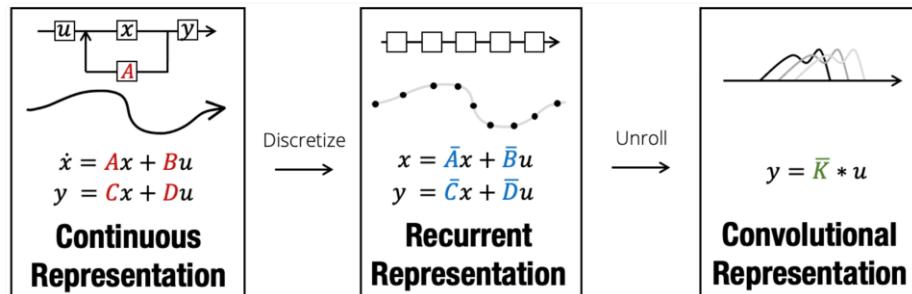
Mamba Head (MambaClassificationHead)



```
1 # Định nghĩa một lớp con MambaTrainer kế thừa từ lớp Trainer.
2 class MambaTrainer(Trainer):
3
4     # Định nghĩa hàm compute_loss để tính toán hàm mất mát trong quá trình huấn luyện.
5     def compute_loss(self, model, inputs, return_outputs=False):
6         # Lấy giá trị input_ids và labels từ inputs.
7         input_ids = inputs.pop("input_ids")
8         labels = inputs.pop('labels')
9
10    # Gọi hàm forward của mô hình với input_ids và labels để nhận các kết quả.
11    outputs = model(input_ids=input_ids, labels=labels)
12
13    # Lấy giá trị loss từ kết quả của mô hình.
14    loss = outputs.loss
15
16    # Trả về cả loss và outputs nếu return_outputs là True, ngược lại chỉ trả về loss.
17    return (loss, outputs) if return_outputs else loss
18
19    # Định nghĩa hàm save_model để lưu model trong quá trình huấn luyện.
20    def save_model(self, output_dir = None, _internal_call = False):
21        # Kiểm tra nếu thư mục lưu trữ không được chỉ định, sử dụng thư mục mặc định từ đối số 'args'.
22        if output_dir is None:
23            output_dir = self.args.output_dir
24
25        # Nếu thư mục đầu ra không tồn tại, tạo mới nó.
26        if not os.path.exists(output_dir):
27            os.makedirs(output_dir)
28
29        # Lưu trạng thái của mô hình PyTorch vào file 'pytorch_model.bin' trong thư mục đầu ra.
30        torch.save(self.model.state_dict(), f"{output_dir}/pytorch_model.bin")
31
32        # Lưu trạng thái của tokenizer vào thư mục đầu ra.
33        self.tokenizer.save_pretrained(output_dir)
34
35        # Lưu cấu hình của mô hình vào file 'config.json' trong thư mục đầu ra.
36        with open(f'{output_dir}/config.json', 'w') as f:
37            json.dump(self.model.config.to_dict(), f)
```

Summary

	RNN	Transformer	Ideal Model
Parallelizable	no	yes	yes
Training Stage	$O(n)$	$O(n^2)$	$O(n)$
Inference Stage	$O(n)$	$O(n^2)$	$O(n)$
Inference per Token	$O(1)$	$O(n)$	$O(1)$



Thanks!

Any questions?