

# Mambda

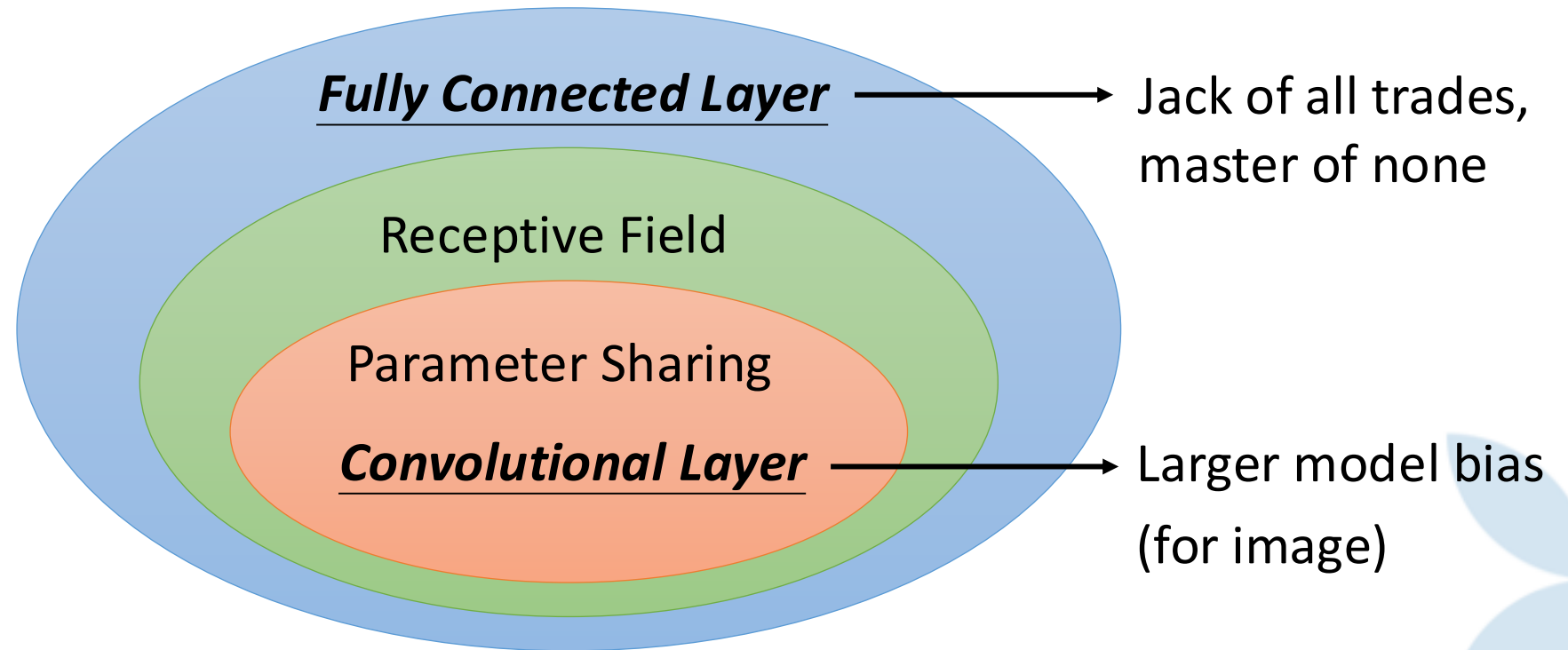
**Qiang Sun**

University of Toronto



# Every architecture has its advantages

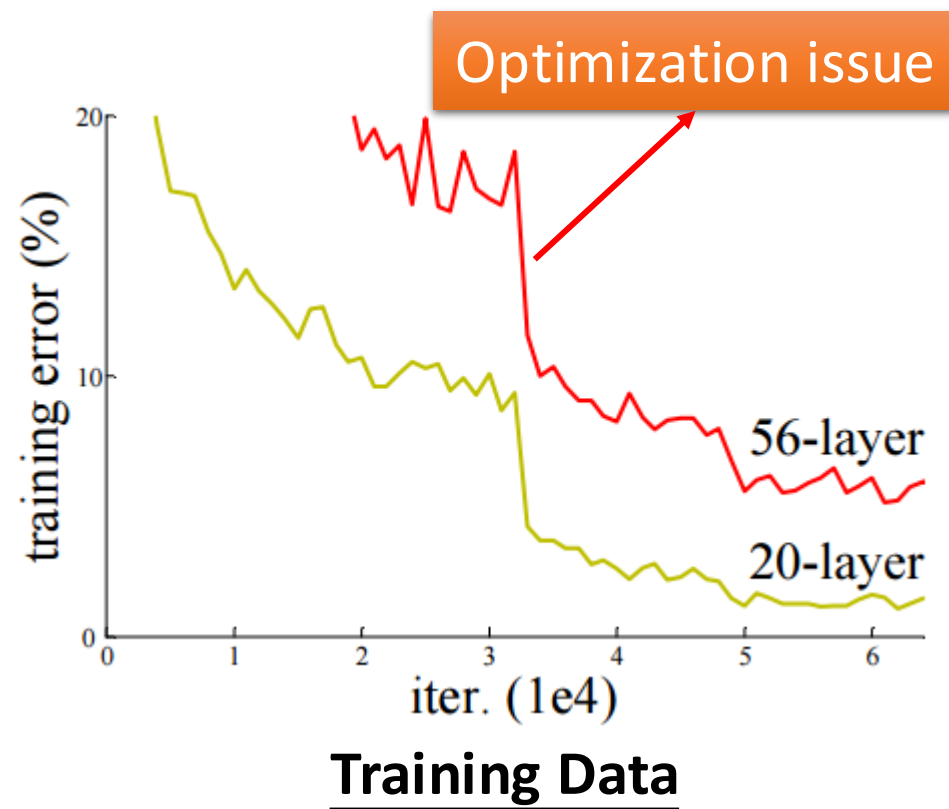
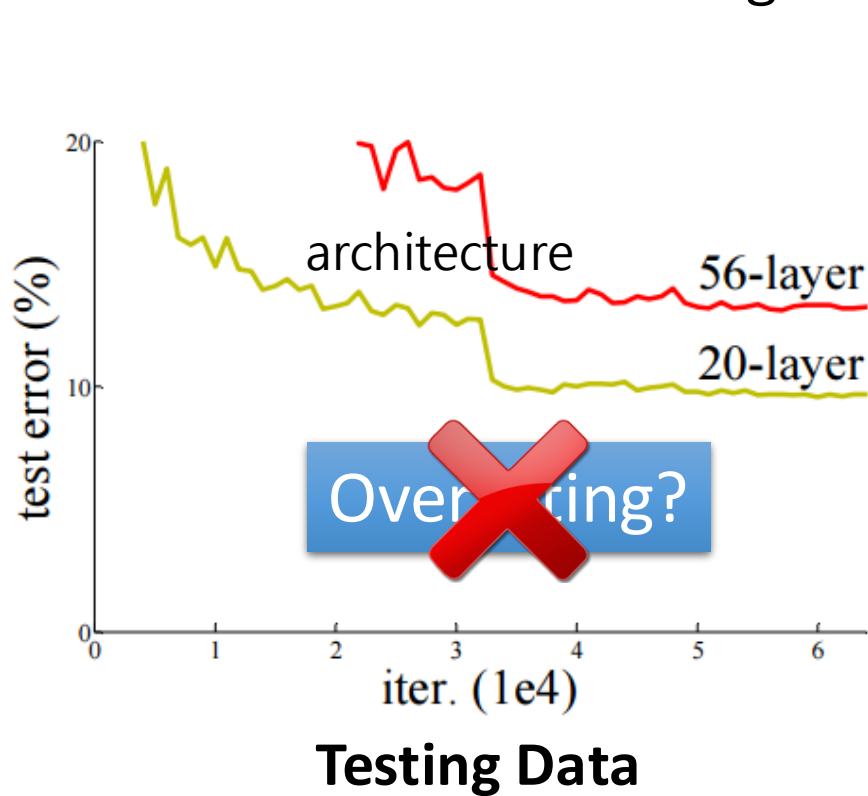
- CNN's advantages ?



Using image priors, reduce unnecessary parameters, **prevent overfitting**

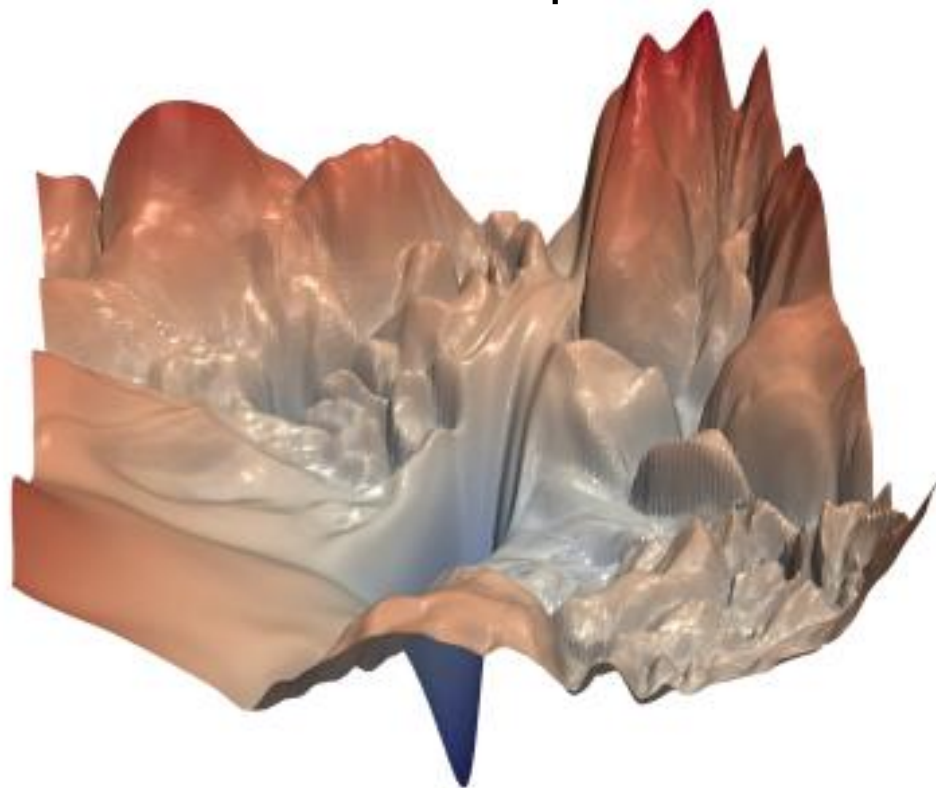
# Every architecture has its advantages

- Residual Connection's advantage?

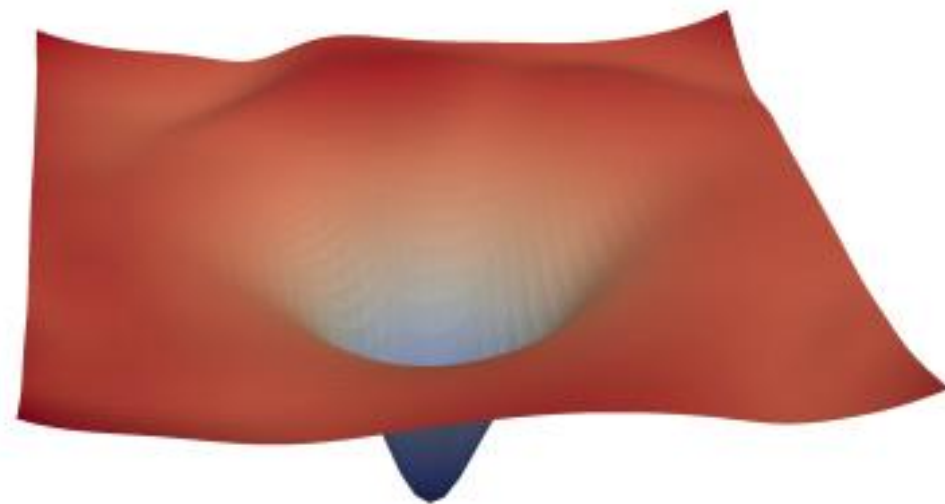


# Every architecture has its advantages

- Residual Connection? Optimize easier and better



(a) without skip connections

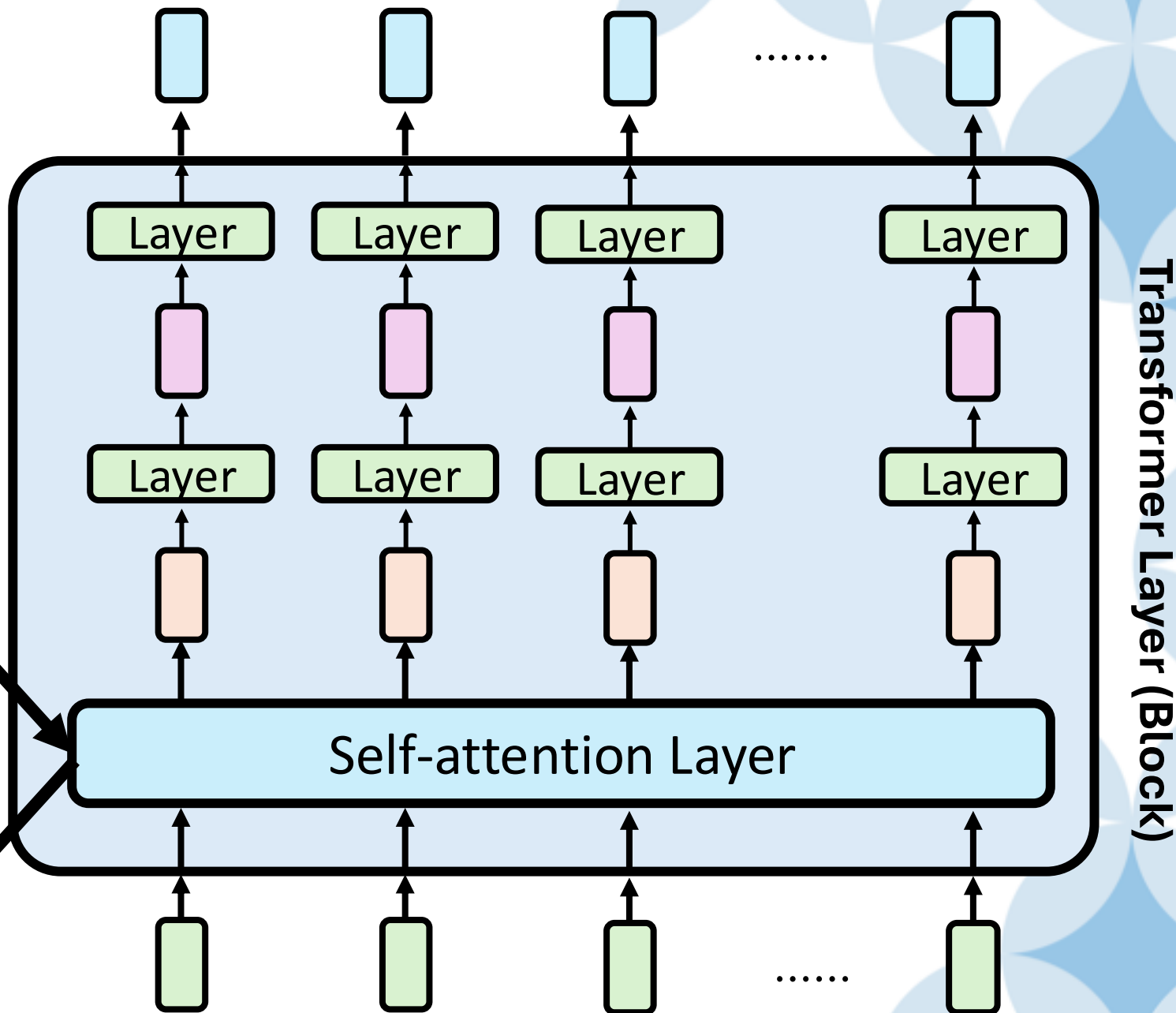


(b) with skip connections

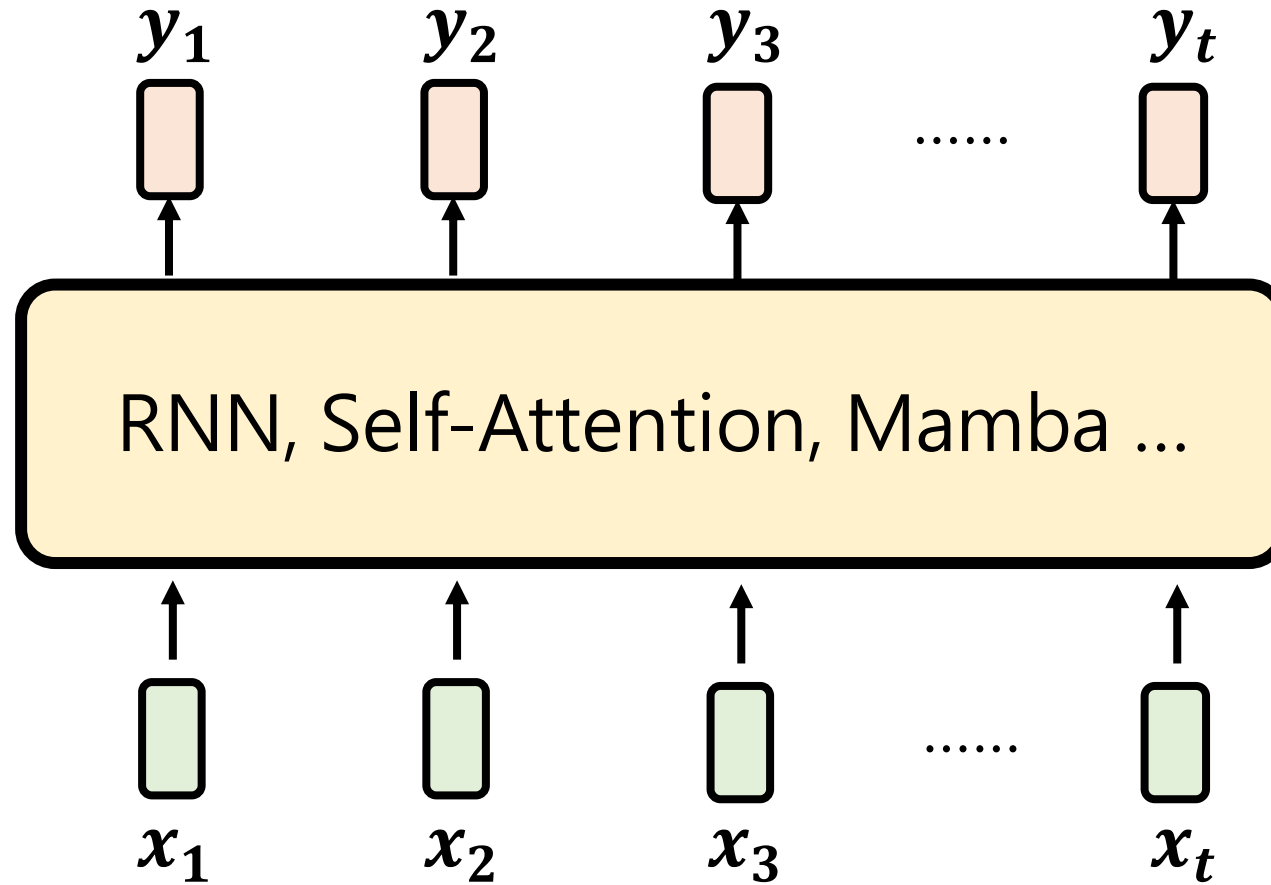
# Transformer?

RNN (LSTM)

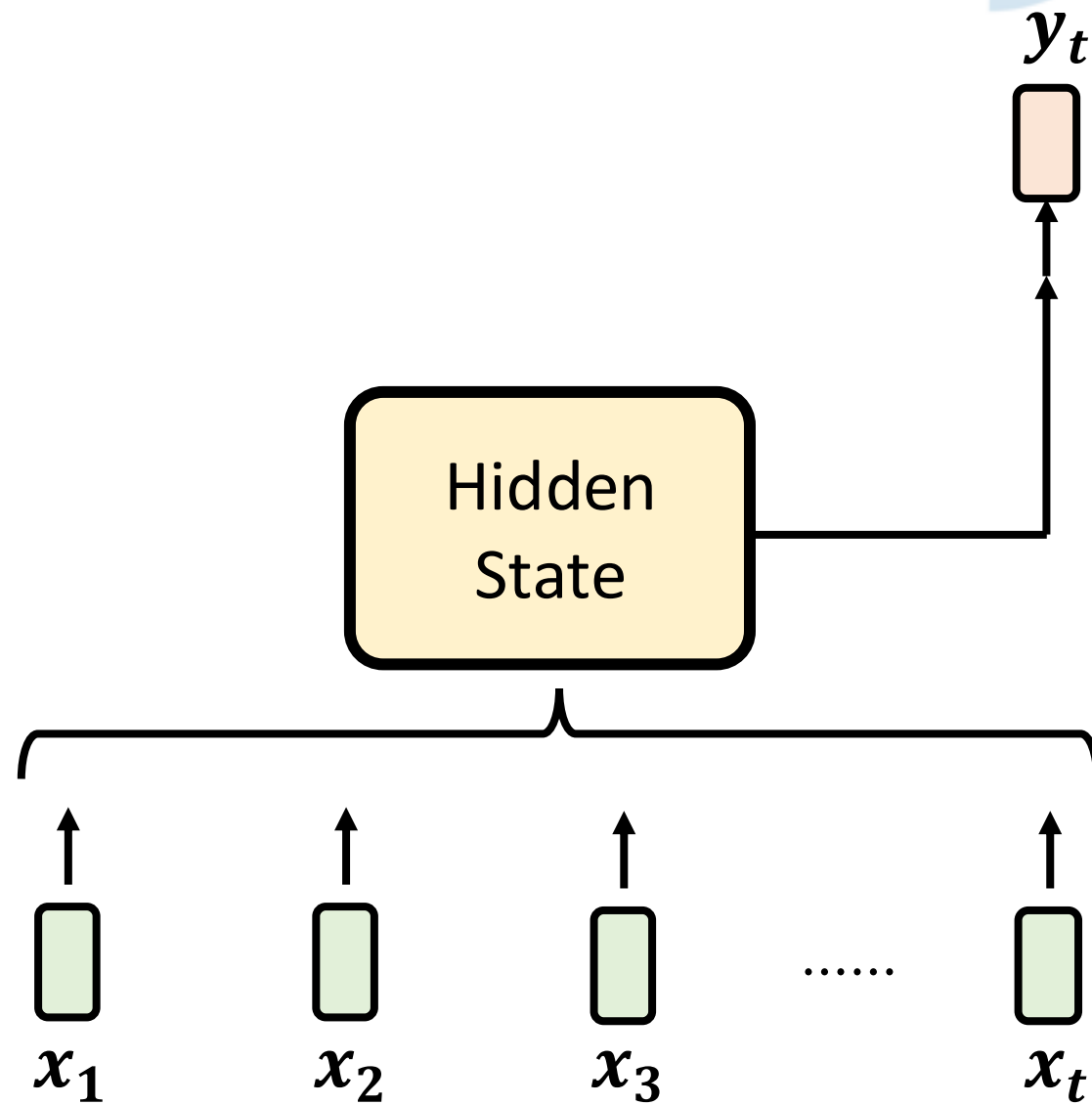
Mamba (and its friends)



# The key question to address



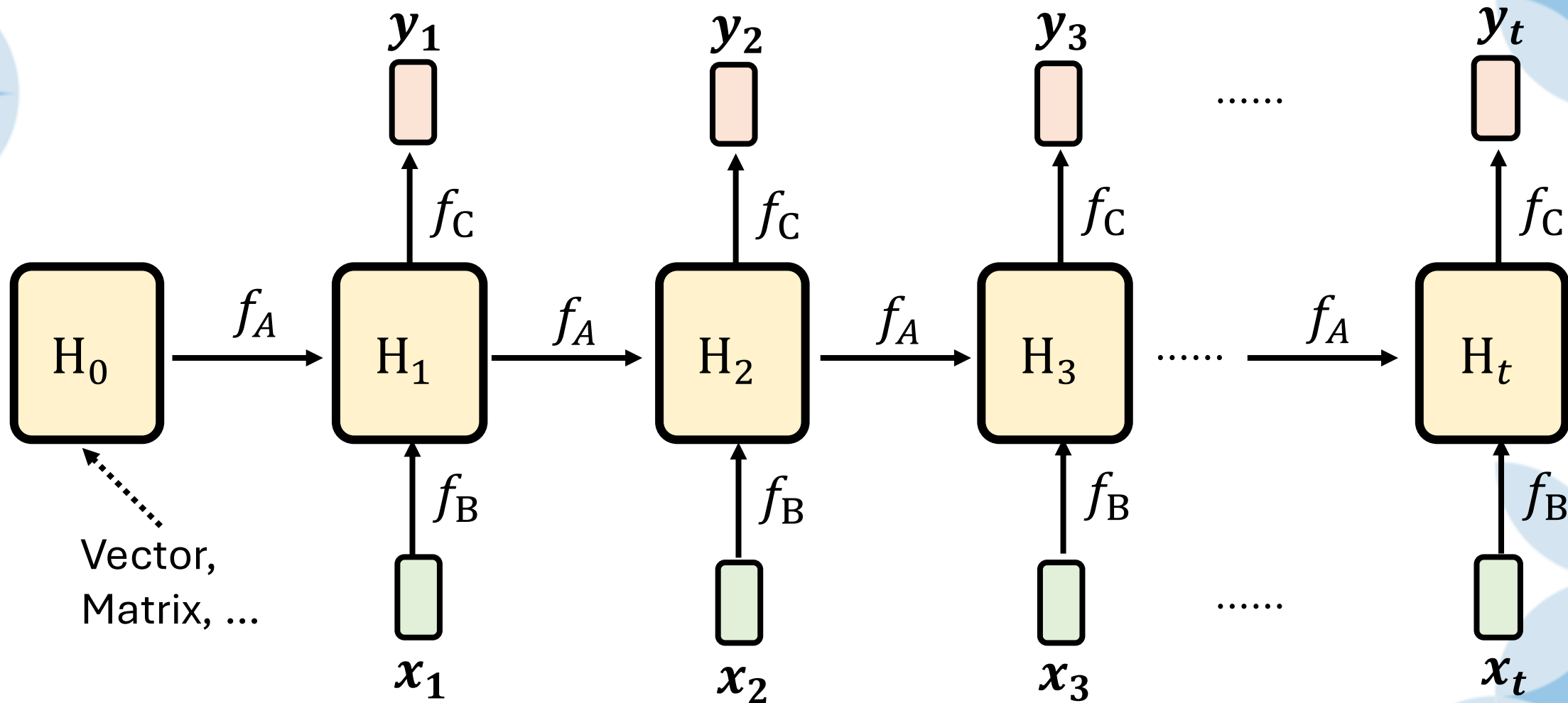
# RNN-Style



# RNN-Style

$$H_t = f_A(H_{t-1}) + f_B(x_t)$$

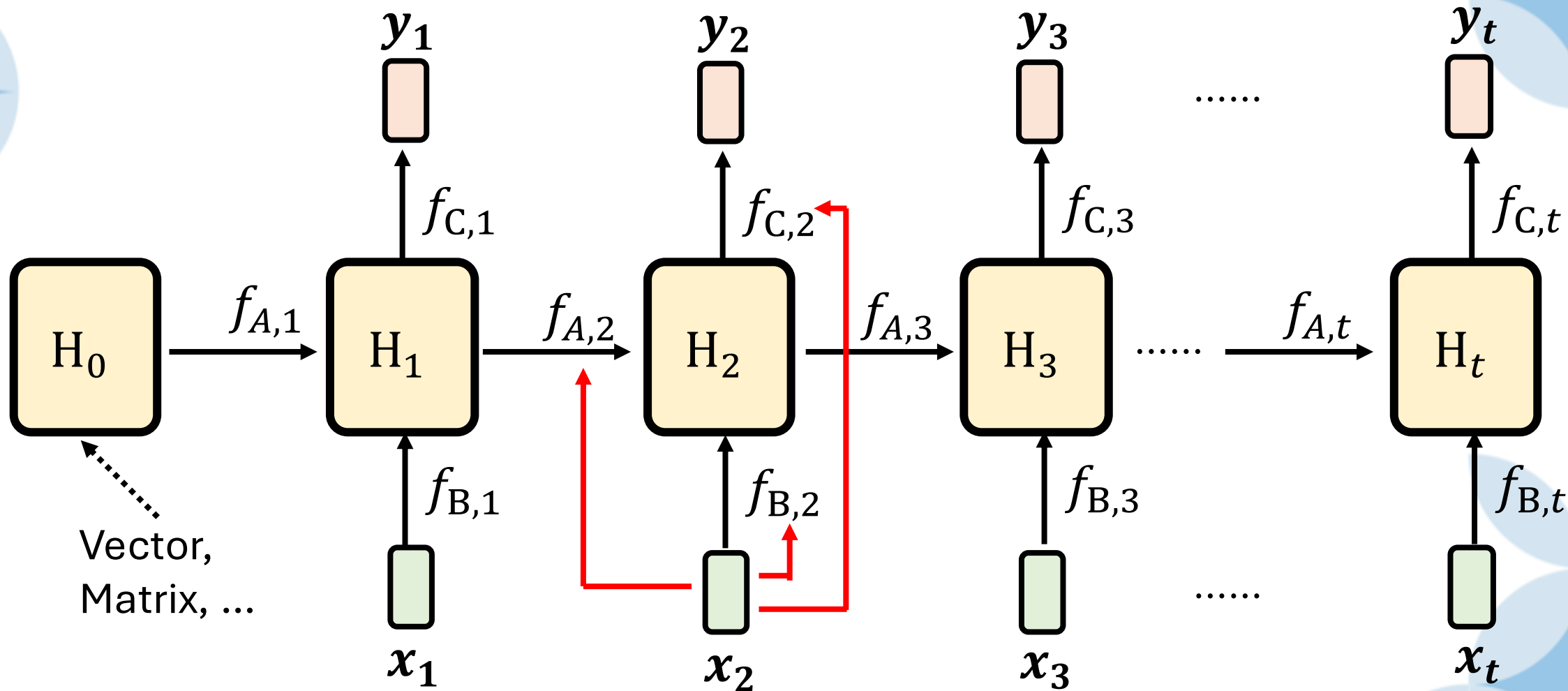
$$y_t = f_C(H_t)$$



# RNN-Style

$$H_t = f_{A,t}(H_{t-1}) + f_{B,t}(x_t)$$

$$y_t = f_{C,t}(H_t)$$



Vector,  
Matrix, ...

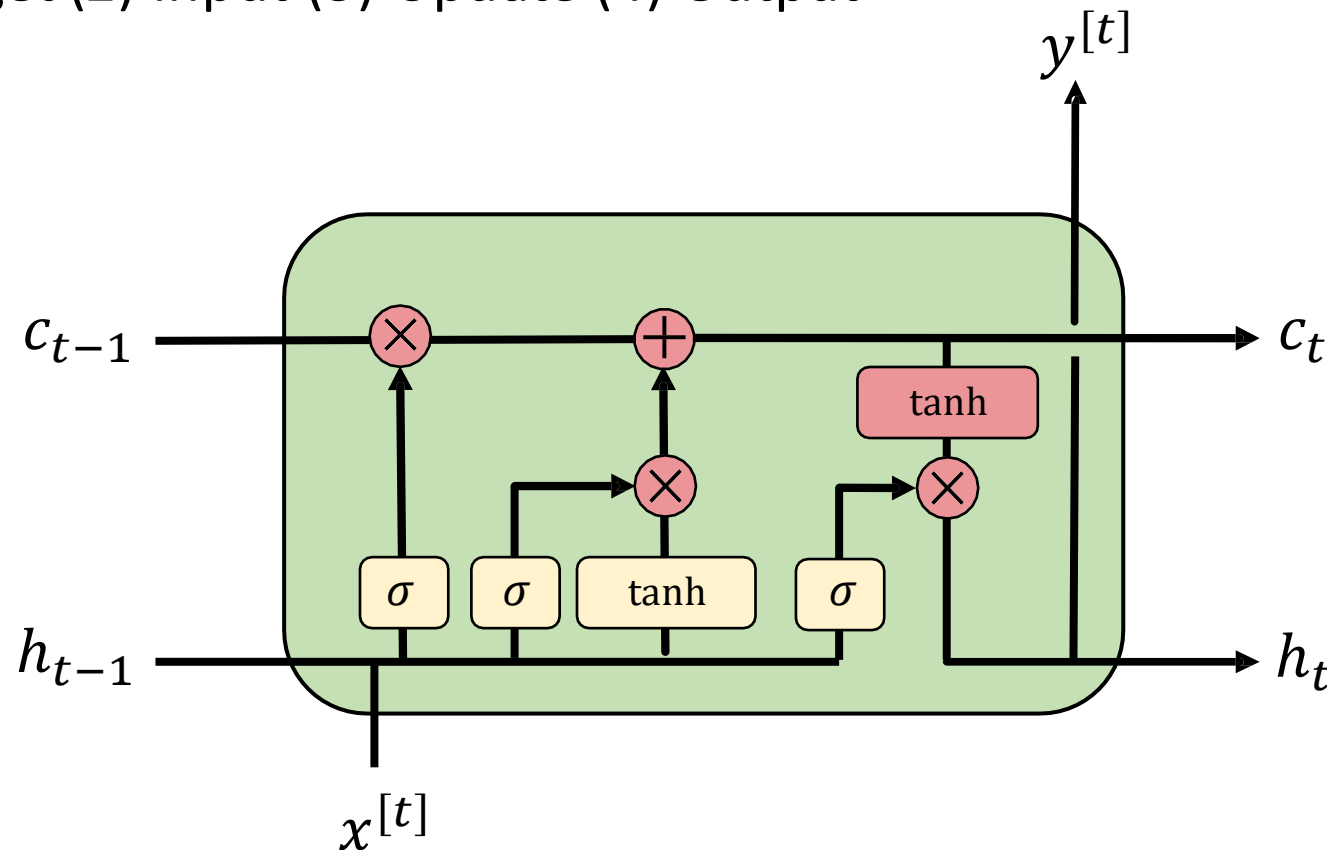
# RNN-Style

## LSTM

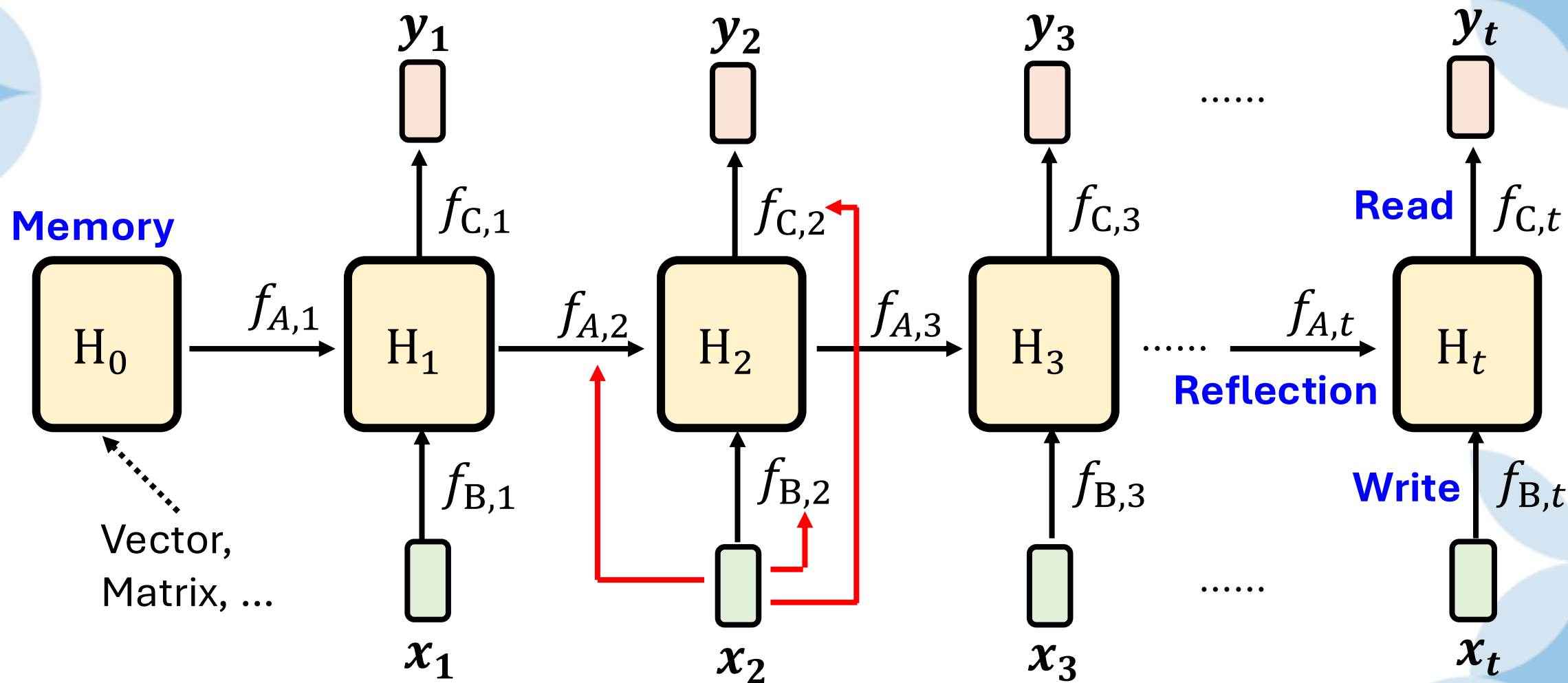
$$H_t = f_{A,t}(H_{t-1}) + f_{B,t}(x_t)$$

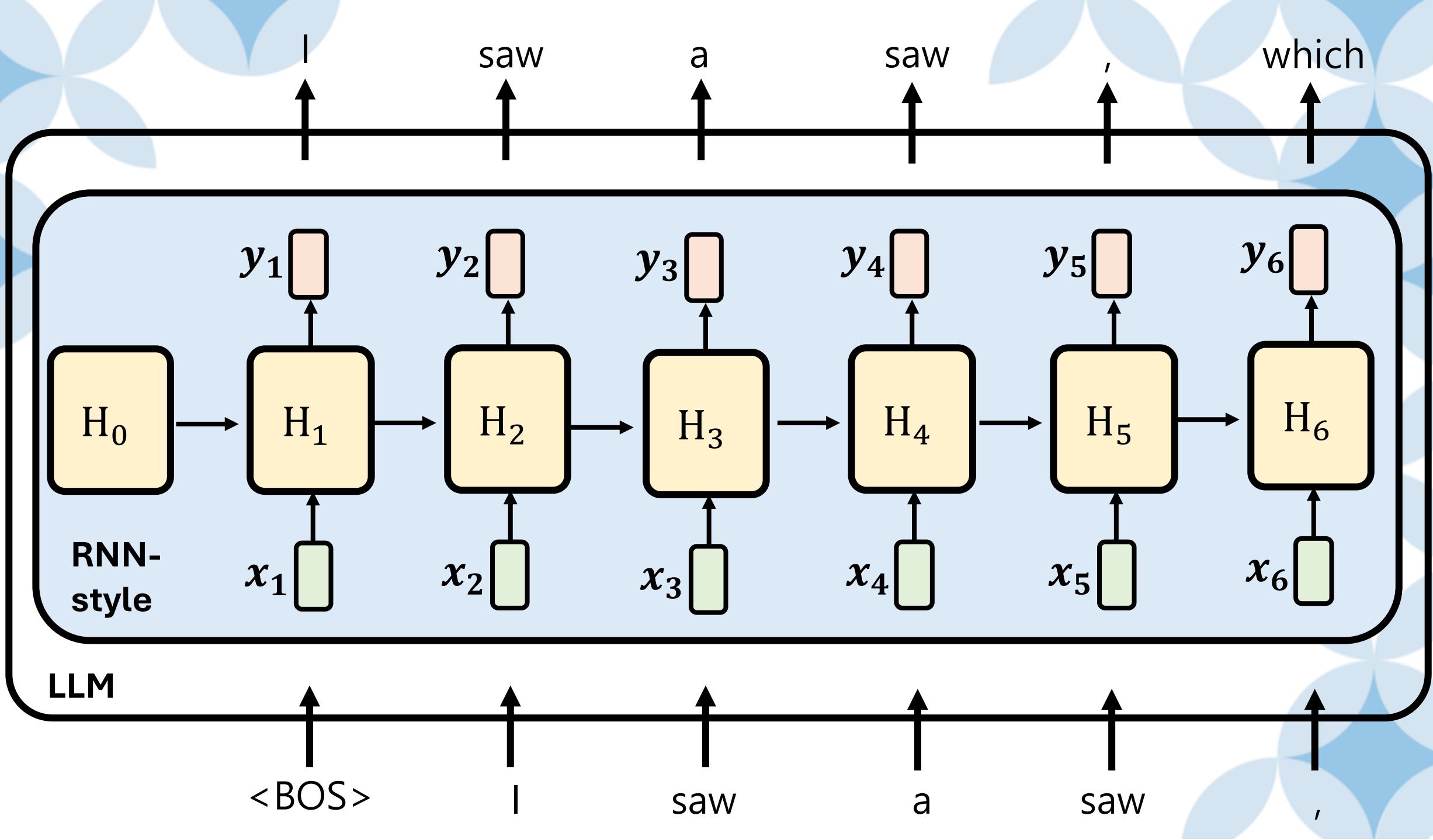
$$y_t = f_{C,t}(H_t)$$

- How do LSTM work: 4 Stages
  - (1) Forget (2) Input (3) Update (4) Output

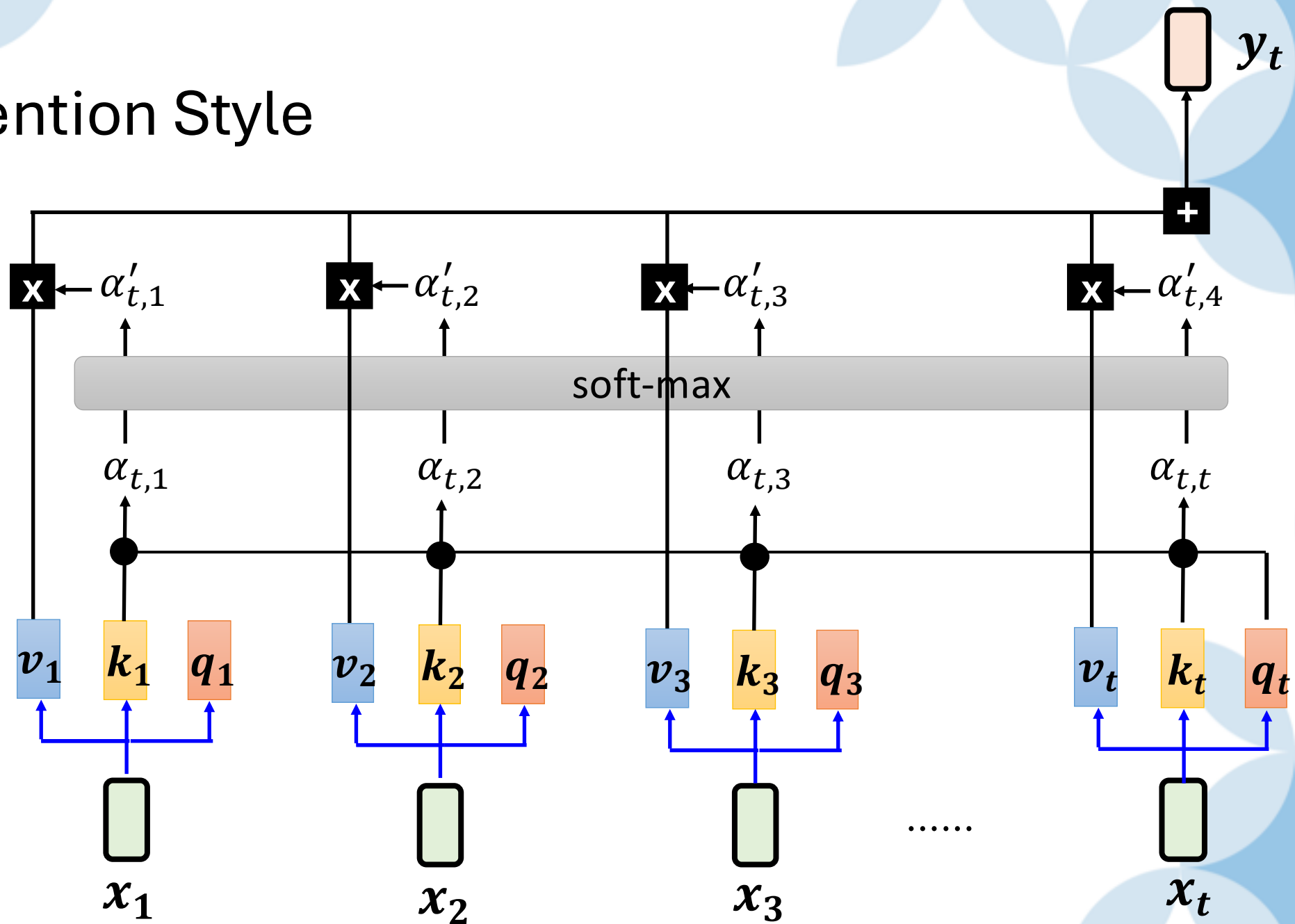


# RNN-Style: An explanation

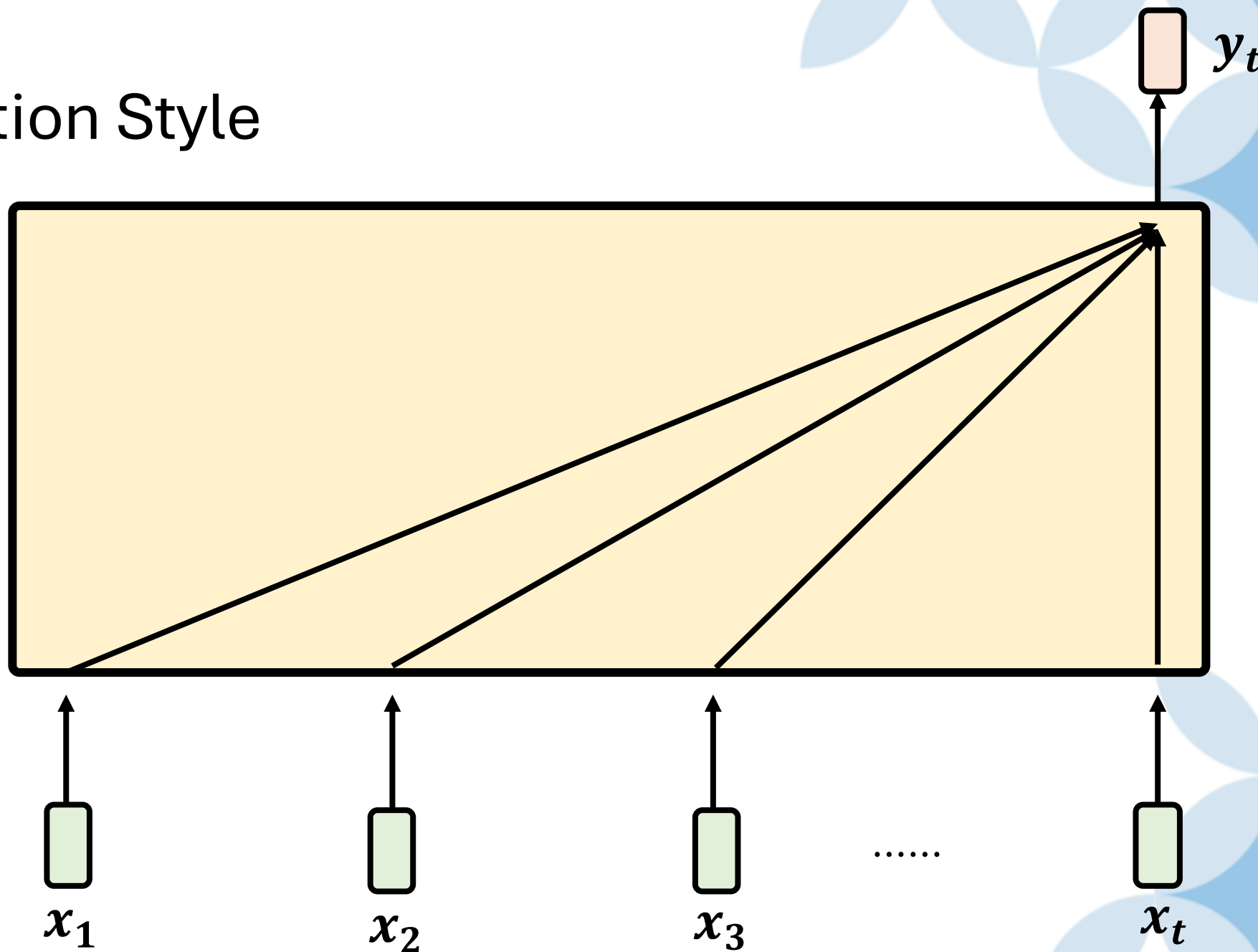




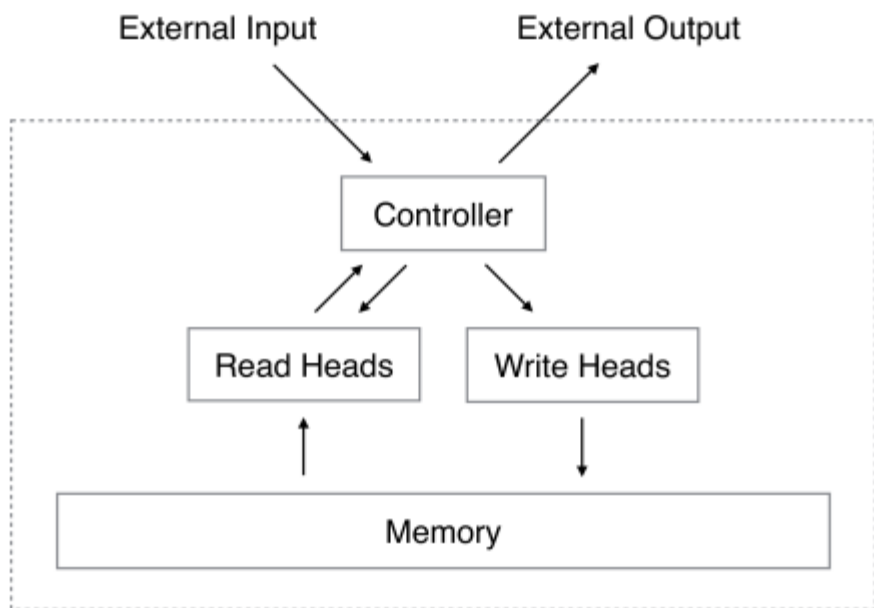
# Self-Attention Style



# Self-Attention Style



# The concept of attention has been around for a long time



## Neural Turing Machine

<https://arxiv.org/abs/1410.5401>

Joe went to the garden then Fred picked up the milk; Joe moved to the bathroom and Fred dropped the milk, and then Dan moved to the living\_room.  
Where is Dan? A: living room I believe  
Where is Joe? A: the bathroom  
Fred moved to the bedroom and Joe went to the kitchen then Joe took the milk there and Dan journeyed to the bedroom; Joe discarded the milk.  
Where is the milk now ? A: the milk is in the kitchen  
Where is Dan now? A: I think he is in the bedroom  
Joe took the milk there, after that Mike travelled to the office, then Joe went to the living\_room, next Dan went back to the kitchen and Joe travelled to the office.  
Where is Joe now? A: I think Joe is in the office

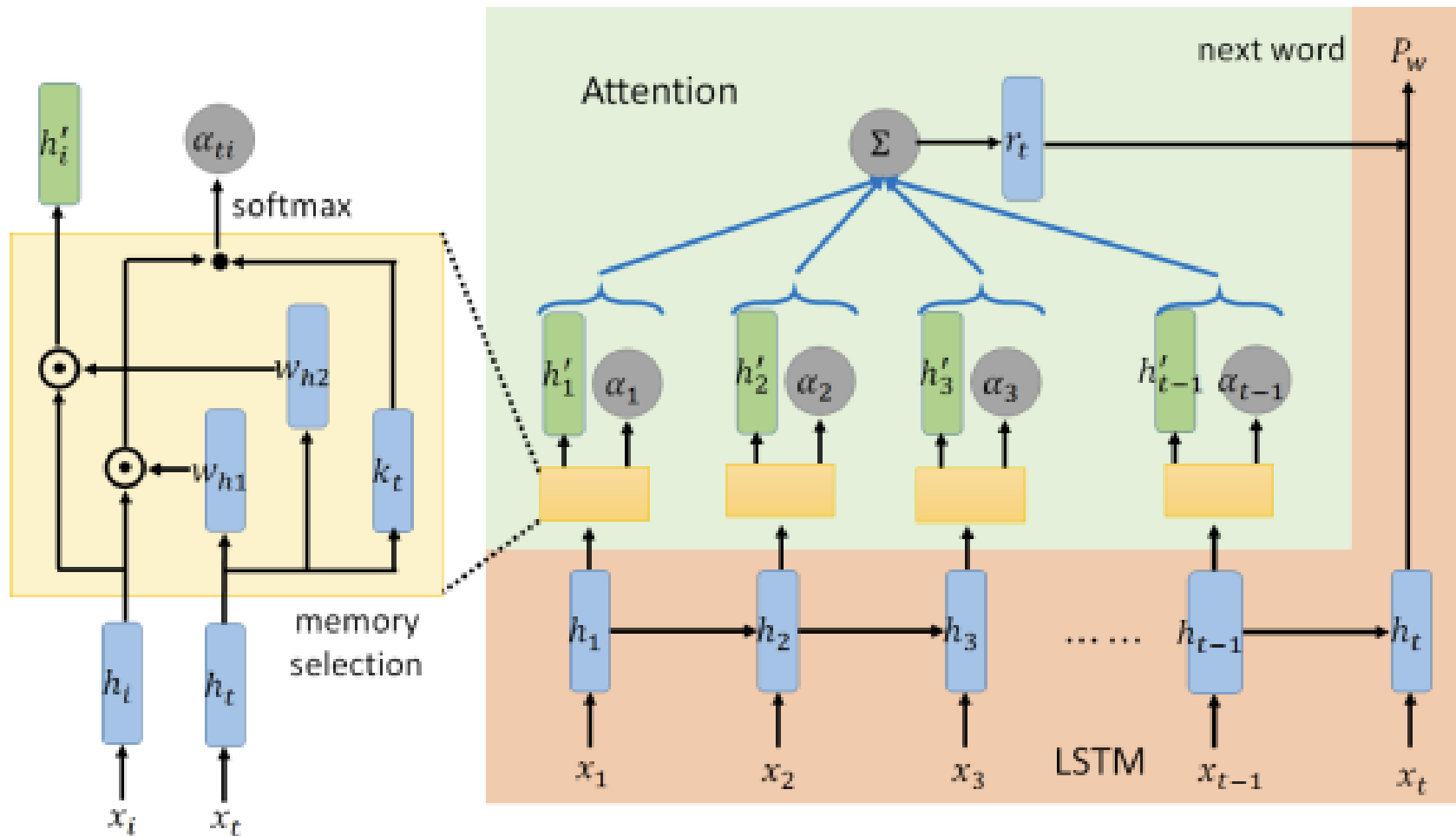
## Memory Networks

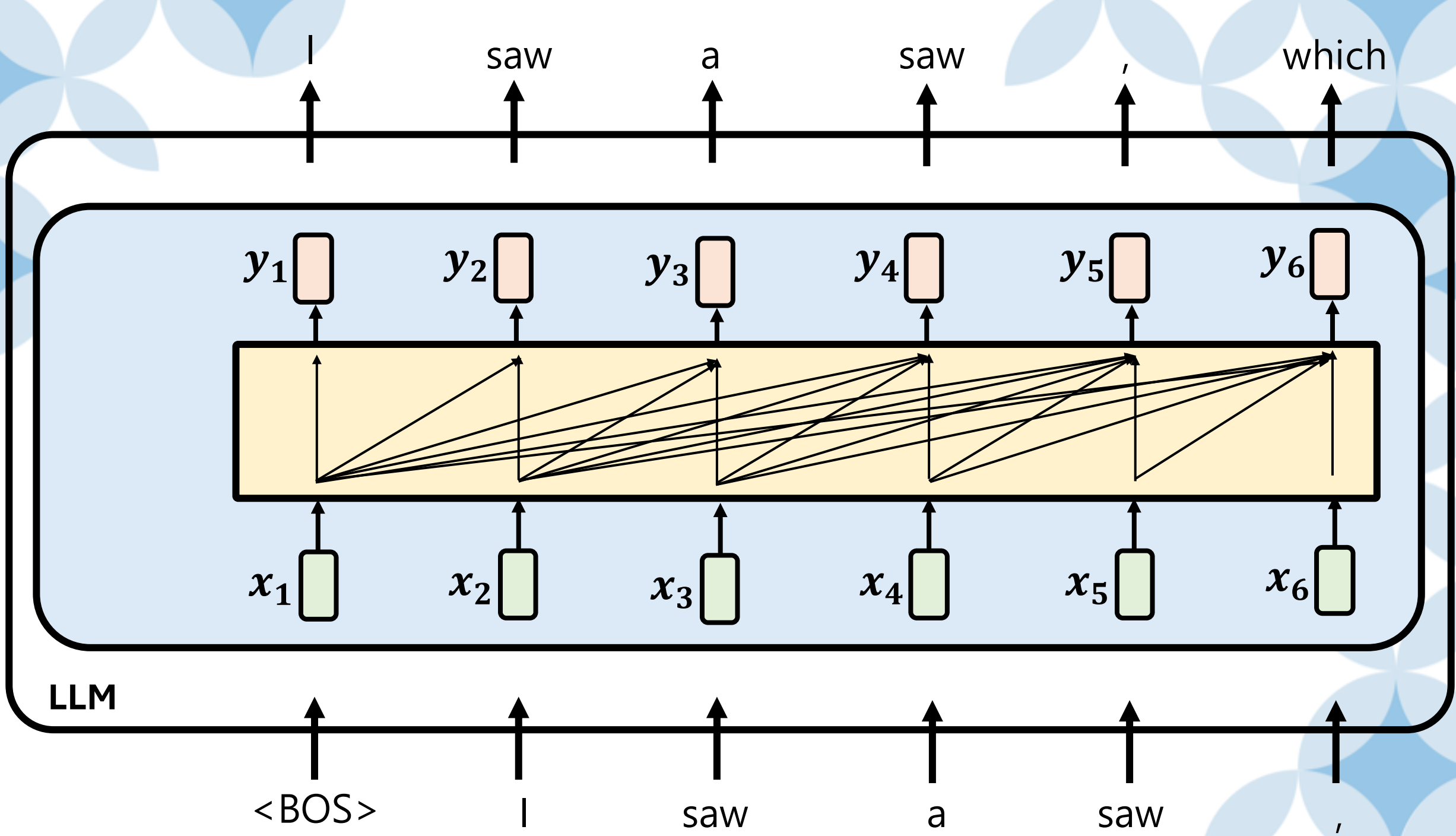
<https://arxiv.org/pdf/1410.3916>

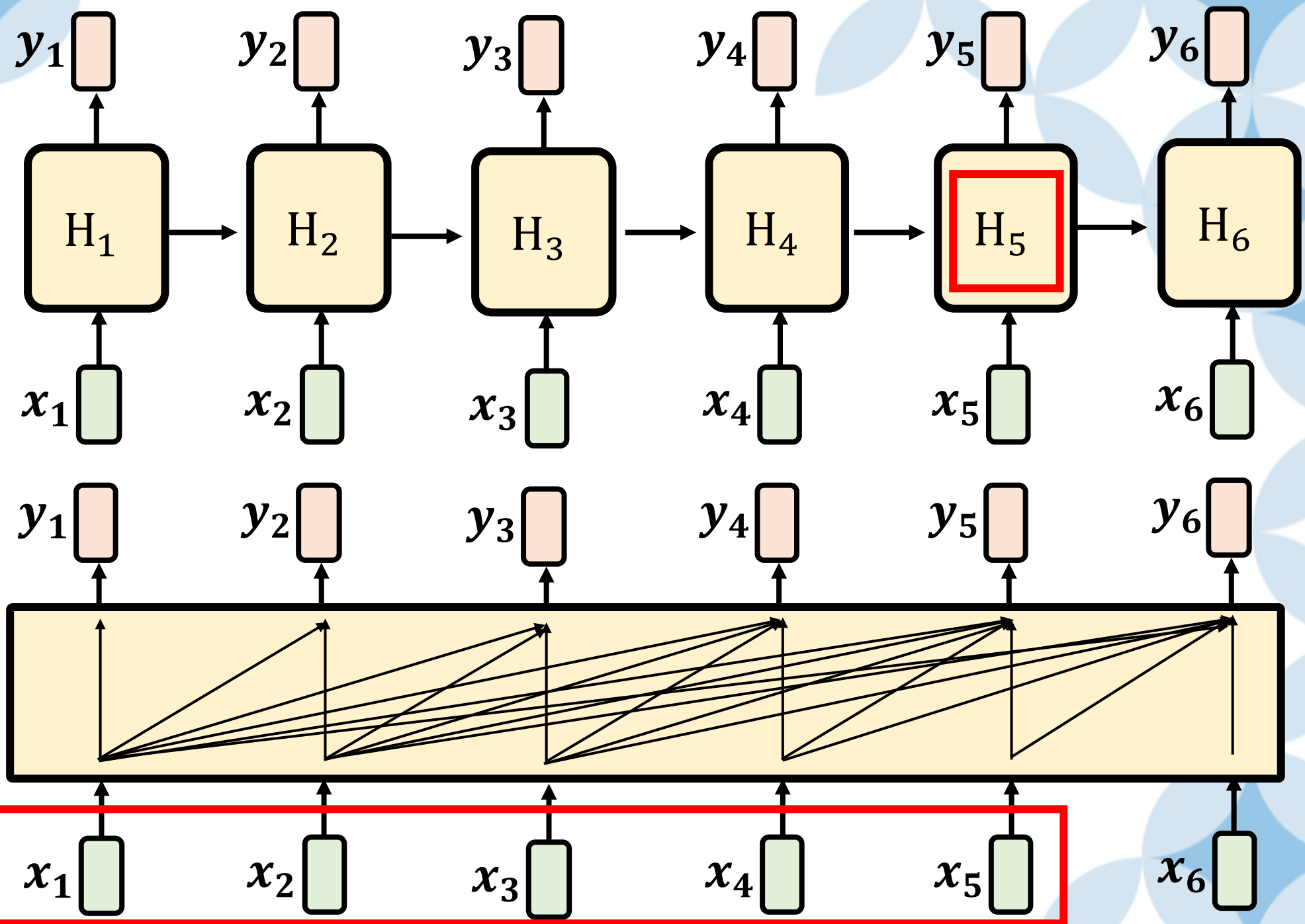
# The concept of attention has been around for a long time

Attention-based Memory Selection Recurrent Network for Language Modeling

<https://arxiv.org/abs/1611.08656>







Same compute in each step

RNN can not rem. long history?

Longer input, more compute.

# Attention Is All You Need

*Not invent Attention, but remove modules other than Attention*

**Ashish Vaswani\***  
Google Brain  
avaswani@google.com

**Noam Shazeer\***  
Google Brain  
noam@google.com

**Niki Parmar\***  
Google Research  
nikip@google.com

**Jakob Uszkoreit\***  
Google Research  
usz@google.com

**Llion Jones\***  
Google Research  
llion@google.com

**Aidan N. Gomez\* †**  
University of Toronto  
aidan@cs.toronto.edu

**Łukasz Kaiser\***  
Google Brain  
lukaszkaizer@google.com

**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com

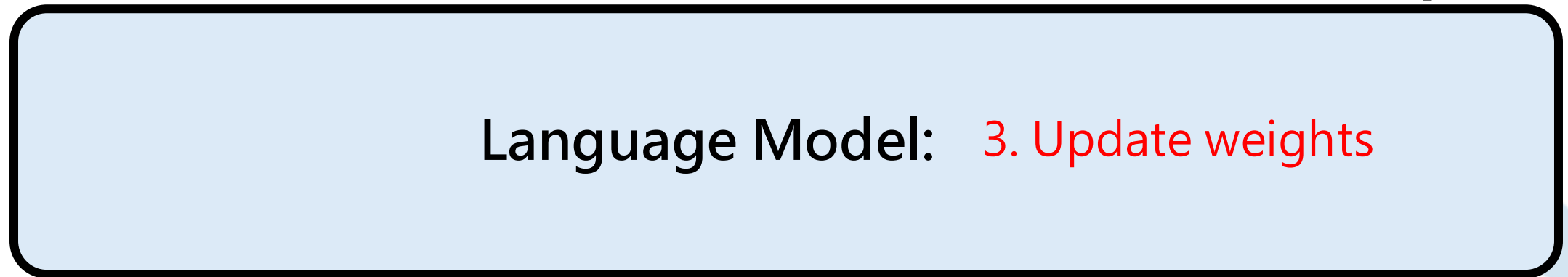
In this work we propose the Transformer, a model architecture eschewing recurrence and instead relying entirely on an attention mechanism to draw global dependencies between input and output. The Transformer allows for significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs.

# Training language models (finding parameters/weights)

- Obtain before updating weights

$$\{z_1, z_2, \dots, z_{t-1}\} \rightarrow z_t$$

1. Obtain current pred.  $???$
2. calculate diff  $\updownarrow$   $z_t$



$z_1$

$z_2$

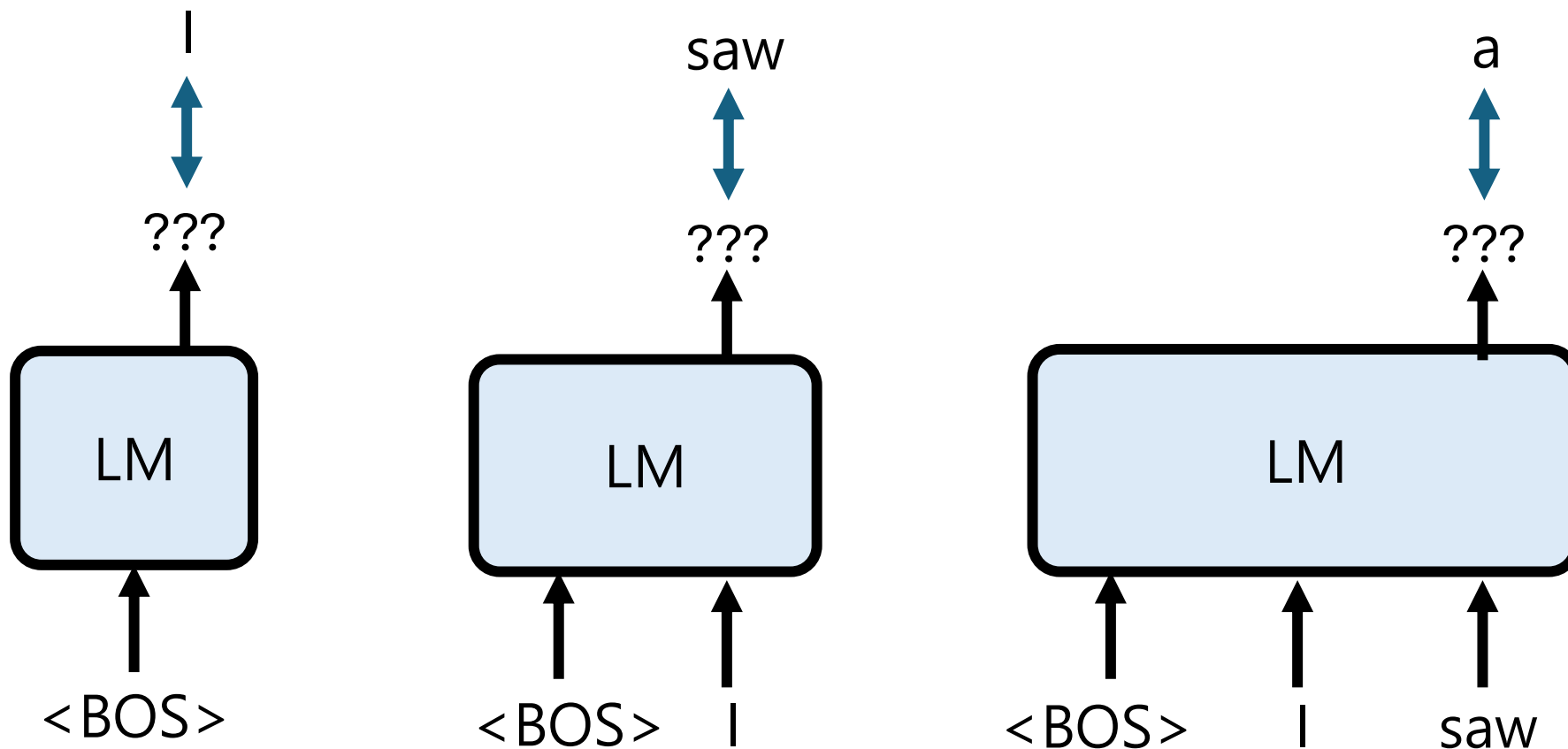
$z_3$

.....

$z_{t-1}$

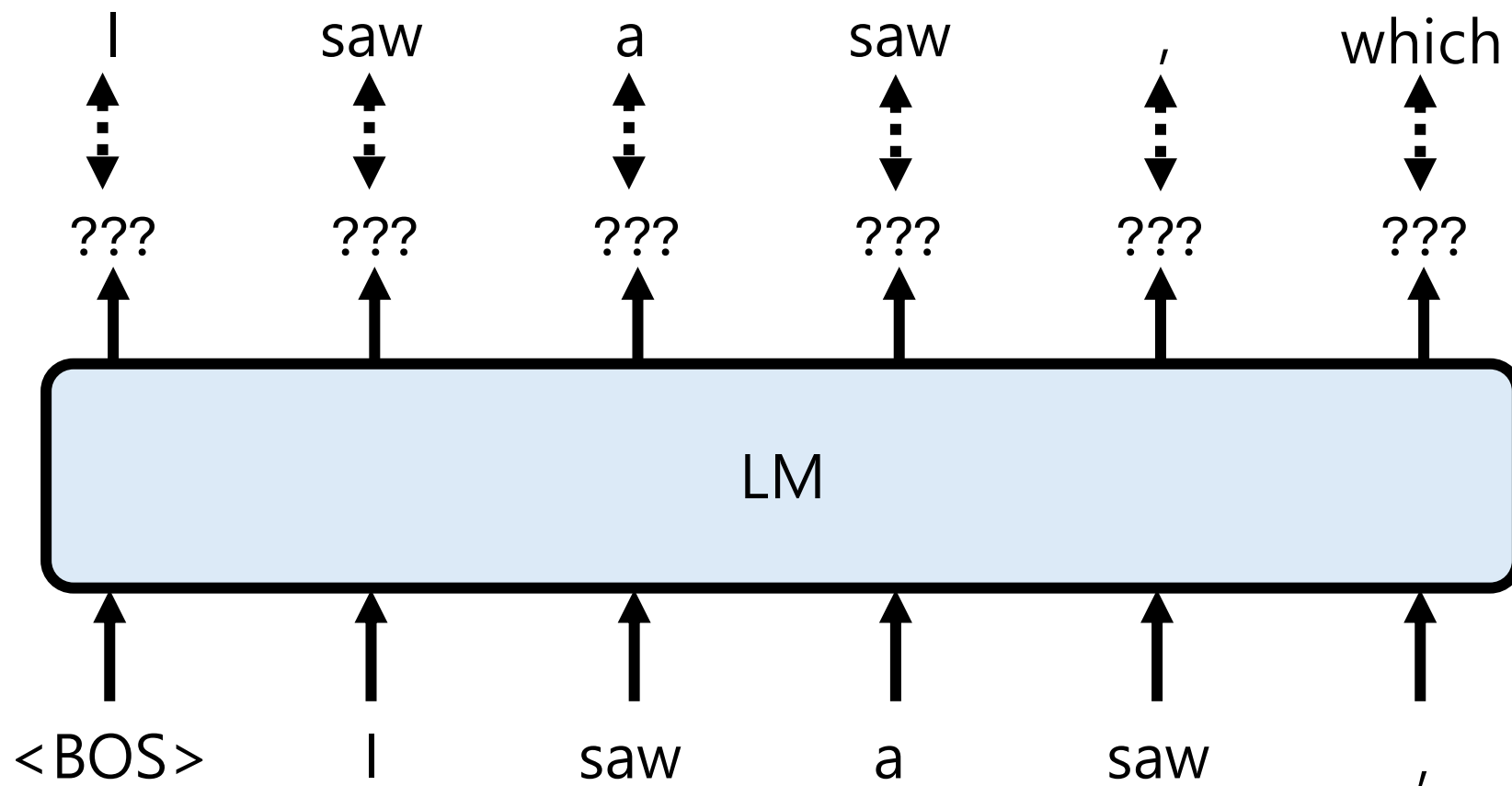
# Training language models

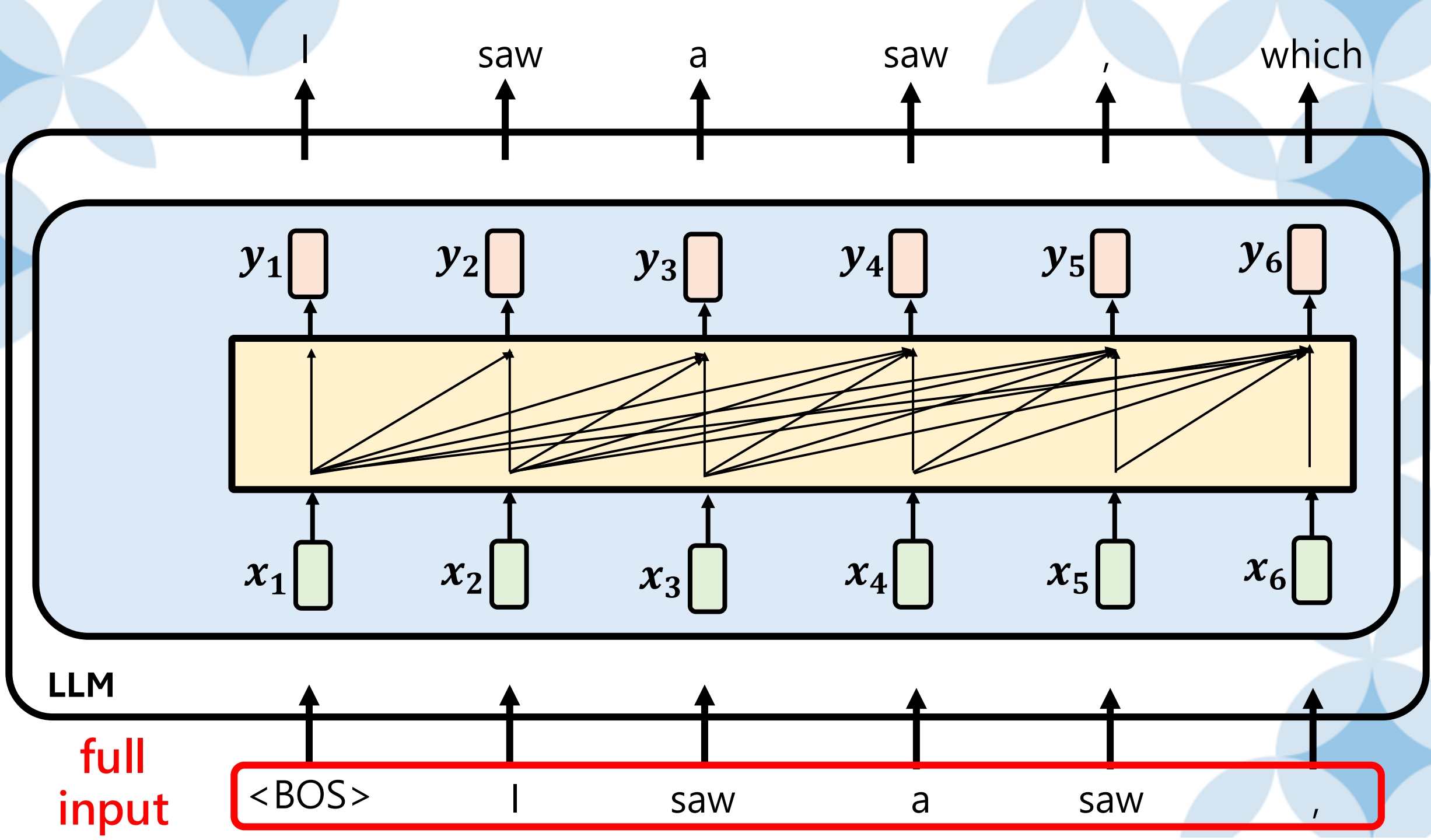
Assume we want model to say 「I saw a saw, which.....」

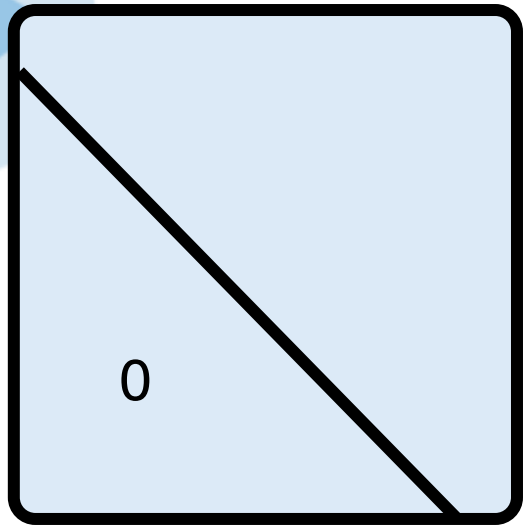


# Training language models

Assume we want model to say 「I saw a saw, which.....」

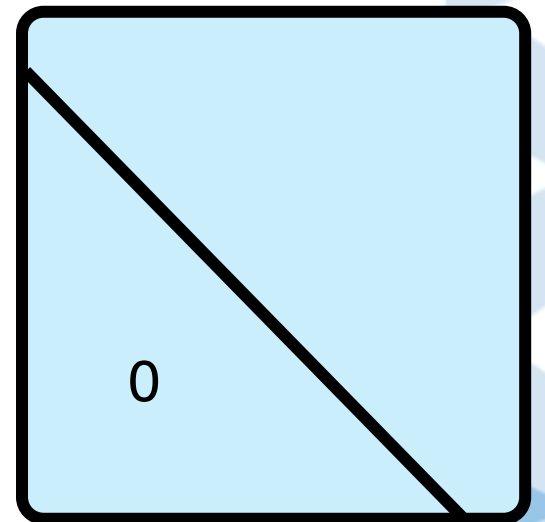
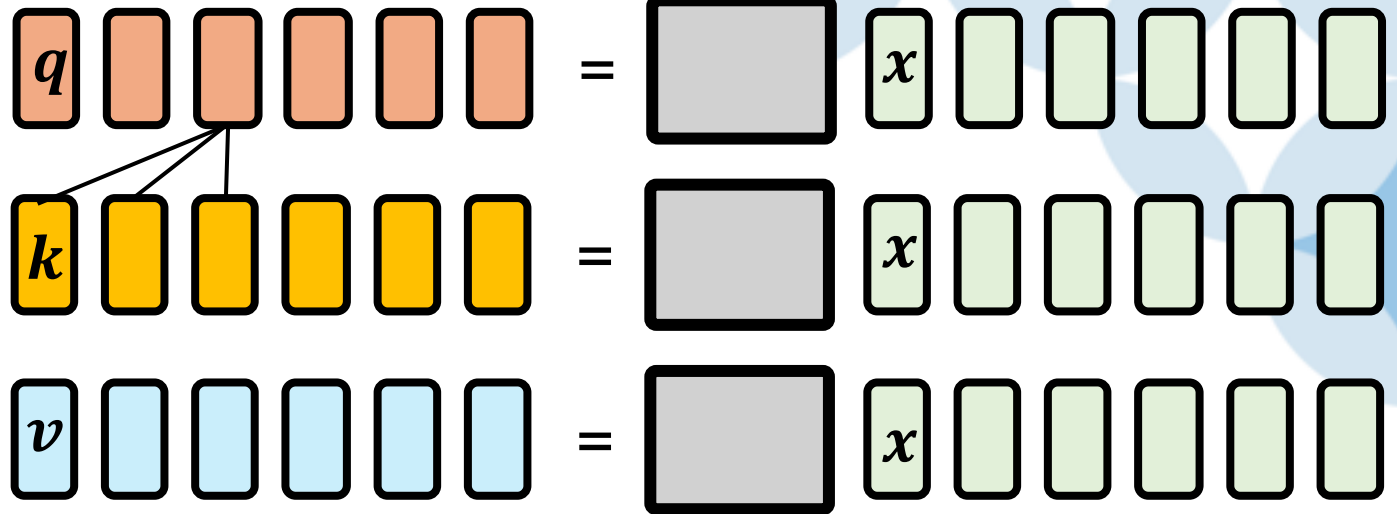
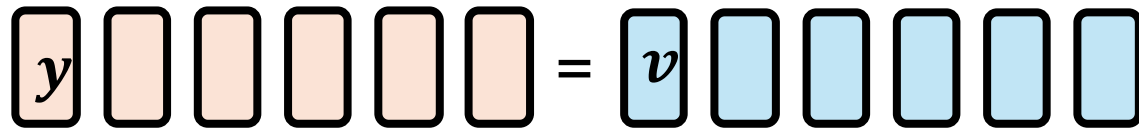
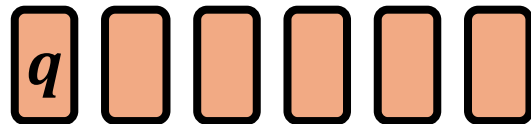
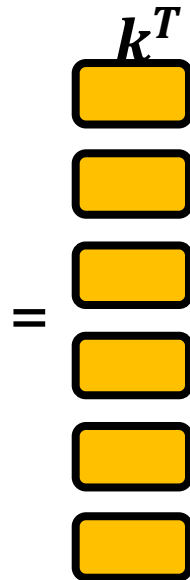


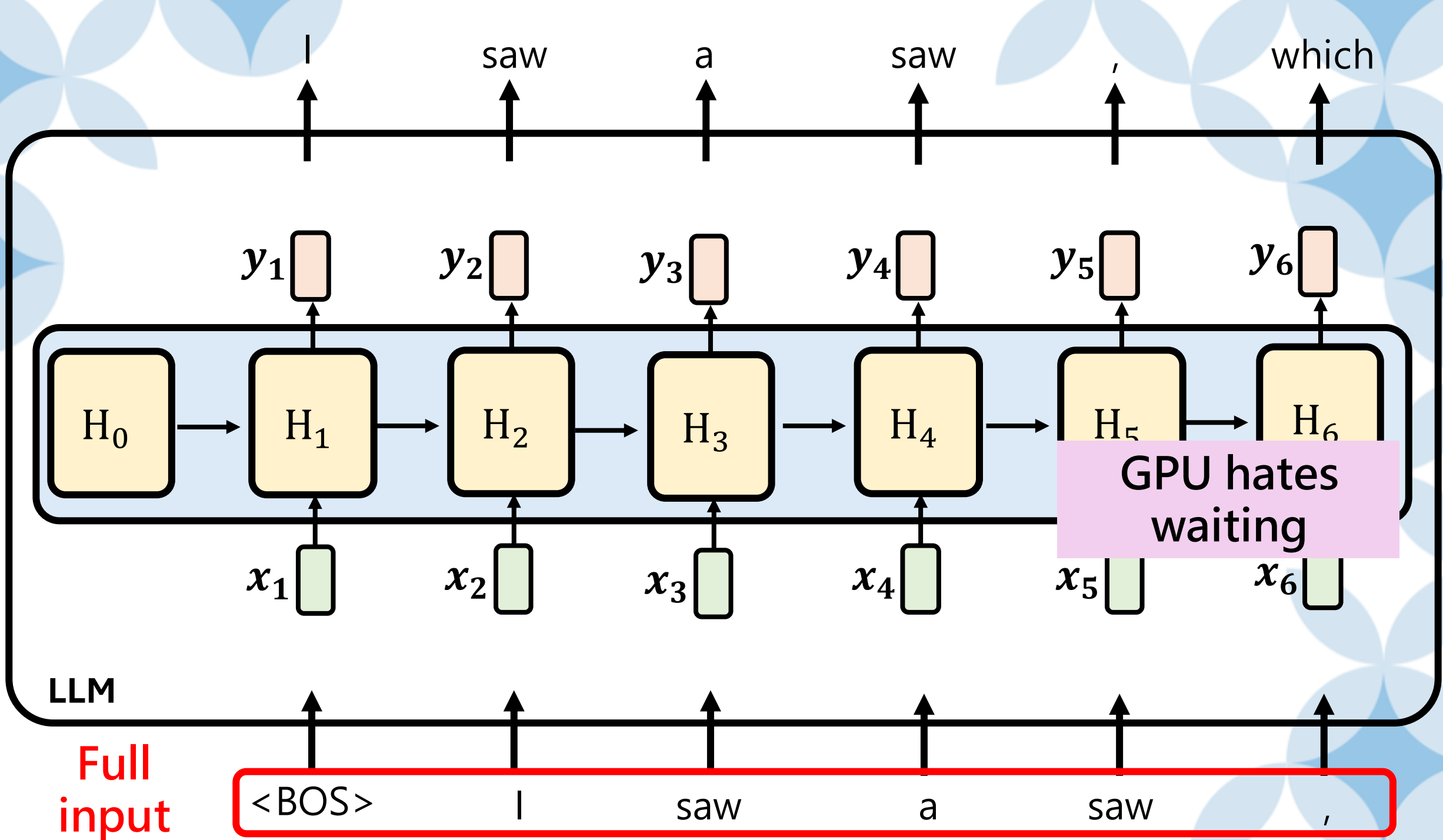




softmax

GPU-friendly computation



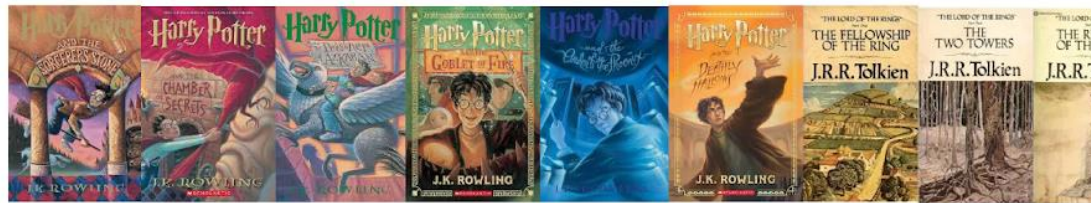


# Self-attention vs. RNN-style

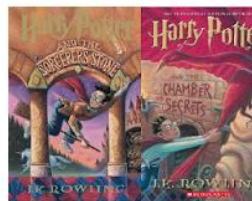
	Self-attention	RNN
Inference	compute, memory increases as the input increases	Fixed compute and memory requirements
Training	Easy to parallelize	Hard to parallelize (?)

# Google's Gemini 1.5 can (almost) fit the entire Harry Potter + Lord of the Ring series in its 2 million context window

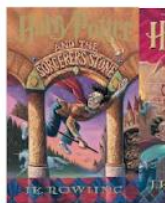
Gemini 1.5 2M  
(June 2024)



Claude 2.1  
(July 2023)



GPT-4 Turbo  
(March 2023)



GPT-3.5 Turbo  
(March 2022)



RAG and AI agents both require language models to handle very long sequences.

Videos and audios are longer sequences than texts

Source of image: <https://www.artfish.ai/p/long-context-llms>

# Is it possible to parallelize RNN training?

$$f_{A,1}(H_0) = 0$$

$$H_t = f_{A,t}(H_{t-1}) + f_{B,t}(\mathbf{x}_t)$$

$$\mathbf{y}_t = f_{C,t}(H_t)$$

$$\left\{ \begin{array}{l} H_1 = f_{A,1}(H_0) + f_{B,1}(\mathbf{x}_1) = f_{B,1}(\mathbf{x}_1) \\ H_2 = f_{A,2}(H_1) + f_{B,2}(\mathbf{x}_2) = f_{A,2}(f_{B,1}(\mathbf{x}_1)) + f_{B,2}(\mathbf{x}_2) \\ H_3 = f_{A,3}(H_2) + f_{B,3}(\mathbf{x}_3) = f_{A,3}(f_{A,2}(f_{B,1}(\mathbf{x}_1)) + f_{B,2}(\mathbf{x}_2)) + f_{B,3}(\mathbf{x}_3) \\ \vdots \\ H_t = f_{A,t}(H_{t-1}) + f_{B,t}(\mathbf{x}_t) = \underbrace{f_{A,t}(f_{A,t-1} \dots f_{A,3}(f_{A,2}(f_{B,1}(\mathbf{x}_1) \dots))}_{\dots} + f_{B,t}(\mathbf{x}_t) \end{array} \right.$$

# Is it possible to parallelize RNN training?

$$f_{A,1}(H_0) = 0$$

$$\left\{ \begin{array}{l} H_1 = H_0 + f_{B,1}(\mathbf{x}_1) \\ H_2 = H_1 + f_{B,2}(\mathbf{x}_2) \\ H_3 = H_2 + f_{B,3}(\mathbf{x}_3) \\ \vdots \\ H_t = H_{t-1} + f_{B,t}(\mathbf{x}_t) \end{array} \right. \quad \begin{array}{l} = f_{B,1}(\mathbf{x}_1) \\ = f_{B,1}(\mathbf{x}_1) + f_{B,2}(\mathbf{x}_2) \\ = f_{B,1}(\mathbf{x}_1) + f_{B,2}(\mathbf{x}_2) + f_{B,3}(\mathbf{x}_3) \\ \dots \\ = f_{B,1}(\mathbf{x}_1) + f_{B,2}(\mathbf{x}_2) + f_{B,3}(\mathbf{x}_3) \quad \dots \dots + f_{B,t}(\mathbf{x}_t) \end{array}$$

$$H_t = H_{t-1} + f_{B,t}(\mathbf{x}_t)$$

$$\mathbf{y}_t = f_{C,t}(H_t)$$

$H_t$  is a  $d \times d$  matrix

$$f_{B,t}(\mathbf{x}_t) = D_t$$

# Is it possible to parallelize RNN training?

$$f_{A,1}(H_0) = 0$$

$$\left\{ \begin{array}{l} H_1 = D_1 \\ H_2 = D_1 + D_2 \\ H_3 = D_1 + D_2 + D_3 \\ \vdots \\ H_t = D_1 + D_2 + \dots + D_t \end{array} \right. \quad \begin{array}{l} \mathbf{y}_1 = D_1 \mathbf{q}_1 \\ \mathbf{y}_2 = D_1 \mathbf{q}_2 + D_2 \mathbf{q}_2 \\ \mathbf{y}_3 = D_1 \mathbf{q}_3 + D_2 \mathbf{q}_3 + D_3 \mathbf{q}_3 \\ \dots \\ \mathbf{y}_t = D_1 \mathbf{q}_t + D_2 \mathbf{q}_t + \dots + D_t \mathbf{q}_t \end{array}$$

$$H_t = H_{t-1} + f_{B,t}(\mathbf{x}_t)$$

$$\mathbf{y}_t = f_{C,t}(H_t)$$

$H_t$  is a  $d \times d$  matrix

$$f_{B,t}(\mathbf{x}_t) = D_t$$

$$f_{C,t}(H_t) = H_t \mathbf{q}_t$$

$$\mathbf{q}_t = W_Q \mathbf{x}_t$$

# Is it possible to parallelize RNN training?

$$f_{A,1}(H_0) = 0$$

$$\left\{ \begin{array}{l} \mathbf{y}_1 = D_1 \mathbf{q}_1 \\ \mathbf{y}_2 = D_1 \mathbf{q}_2 + D_2 \mathbf{q}_2 \\ \mathbf{y}_3 = D_1 \mathbf{q}_3 + D_2 \mathbf{q}_3 + D_3 \mathbf{q}_3 \\ \vdots \\ \mathbf{y}_t = D_1 \mathbf{q}_t + D_2 \mathbf{q}_t + \dots + D_t \mathbf{q}_t \end{array} \right.$$

$$H_t = H_{t-1} + f_{B,t}(\mathbf{x}_t)$$

$$\mathbf{y}_t = f_{C,t}(H_t)$$

$H_t$  is a  $d \times d$  matrix

$$f_{B,t}(\mathbf{x}_t) = D_t$$

$$D_t = \mathbf{v}_t \mathbf{k}_t^T \quad \begin{array}{l} \mathbf{v}_t = W_v \mathbf{x}_t \\ \mathbf{k}_t = W_k \mathbf{x}_t \end{array}$$

$$f_{C,t}(H_t) = H_t \mathbf{q}_t$$

$$\mathbf{q}_t = W_Q \mathbf{x}_t$$

# Is it possible to parallelize RNN training?

$$f_{A,1}(H_0) = 0$$

$$\left\{ \begin{array}{l} y_1 = v_1 k_1^T q_1 \\ y_2 = v_1 k_1^T q_2 + v_2 k_2^T q_2 \\ y_3 = v_1 k_1^T q_3 + v_2 k_2^T q_3 + v_3 k_3^T q_3 \\ \vdots \\ y_t = v_1 k_1^T q_t + v_2 k_2^T q_t + \dots + v_t k_t^T q_t \end{array} \right.$$

$$H_t = H_{t-1} + f_{B,t}(x_t)$$

$$y_t = f_{C,t}(H_t)$$

$H_t$  is a  $d \times d$  matrix

$$f_{B,t}(x_t) = D_t$$

$$D_t = v_t k_t^T \quad \begin{array}{l} v_t = W_v x_t \\ k_t = W_k x_t \end{array}$$

$$f_{C,t}(H_t) = H_t q_t$$

$$q_t = W_Q x_t$$

# Is it possible to parallelize RNN training?

$$f_{A,1}(H_0) = 0$$

$$\mathbf{y}_t = \mathbf{v}_1 \mathbf{k}_1^T \mathbf{q}_t + \mathbf{v}_2 \mathbf{k}_2^T \mathbf{q}_t + \dots + \mathbf{v}_t \mathbf{k}_t^T \mathbf{q}_t$$

$$= \mathbf{v}_1 a_{t,1} + \mathbf{v}_2 a_{t,2} + \dots + \mathbf{v}_t a_{t,t}$$

$$= a_{t,1} \mathbf{v}_1 + a_{t,2} \mathbf{v}_2 + \dots + a_{t,t} \mathbf{v}_t$$

This is self-attention! (w/o softmax)

⇒ Linear Attention

$$H_t = H_{t-1} + f_{B,t}(\mathbf{x}_t)$$

$$\mathbf{y}_t = f_{C,t}(H_t)$$

$H_t$  is a  $d \times d$  matrix

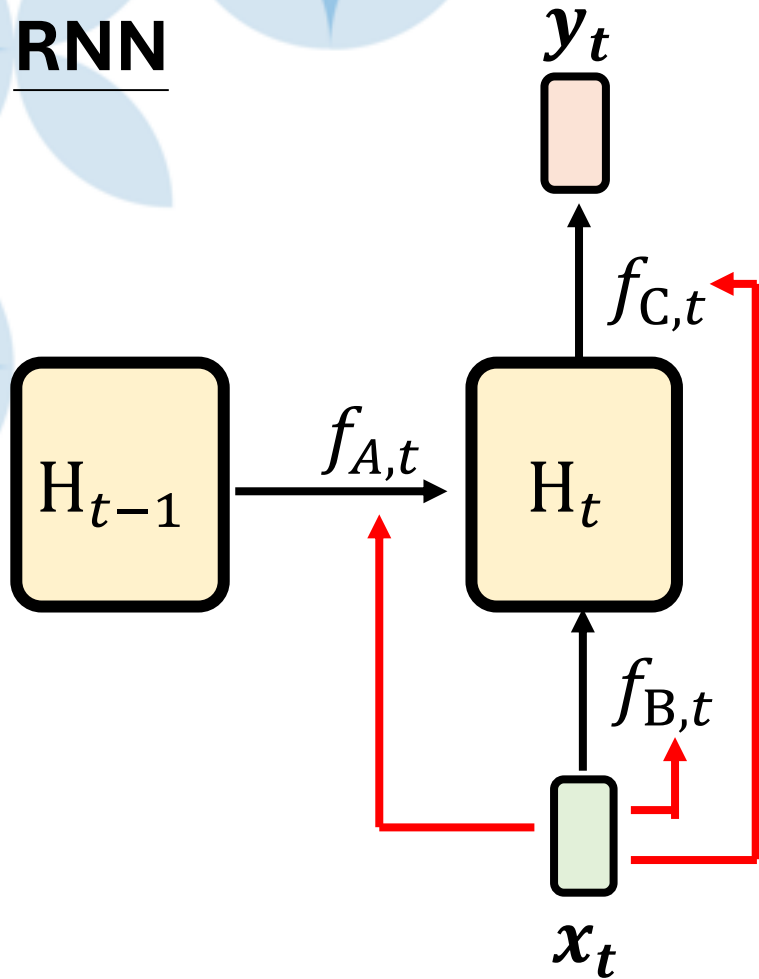
$$f_{B,t}(\mathbf{x}_t) = D_t$$

$$D_t = \mathbf{v}_t \mathbf{k}_t^T \quad \begin{array}{l} \mathbf{v}_t = W_v \mathbf{x}_t \\ \mathbf{k}_t = W_k \mathbf{x}_t \end{array}$$

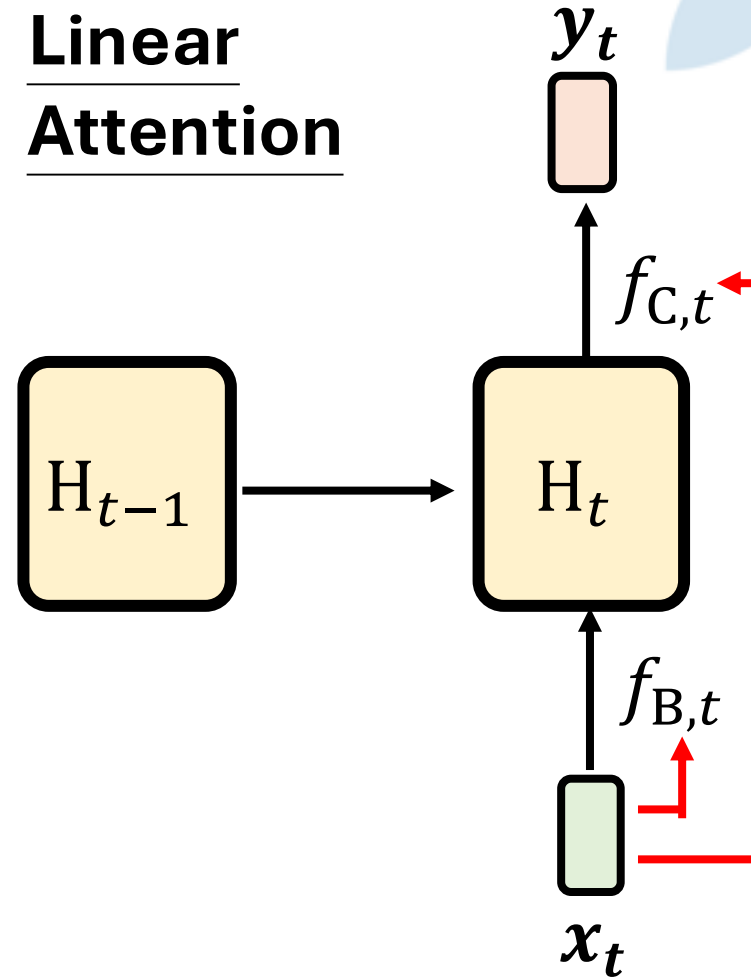
$$f_{C,t}(H_t) = H_t \mathbf{q}_t$$

$$\mathbf{q}_t = W_Q \mathbf{x}_t$$

## RNN



## Linear Attention



$$f_{C,t}(H_t) = H_t \mathbf{q}_t$$
$$\mathbf{q}_t = W_Q \mathbf{x}_t$$

$$f_{B,t}(\mathbf{x}_t) = \mathbf{v}_t \mathbf{k}_t^T$$

$$\mathbf{v}_t = W_v \mathbf{x}_t$$

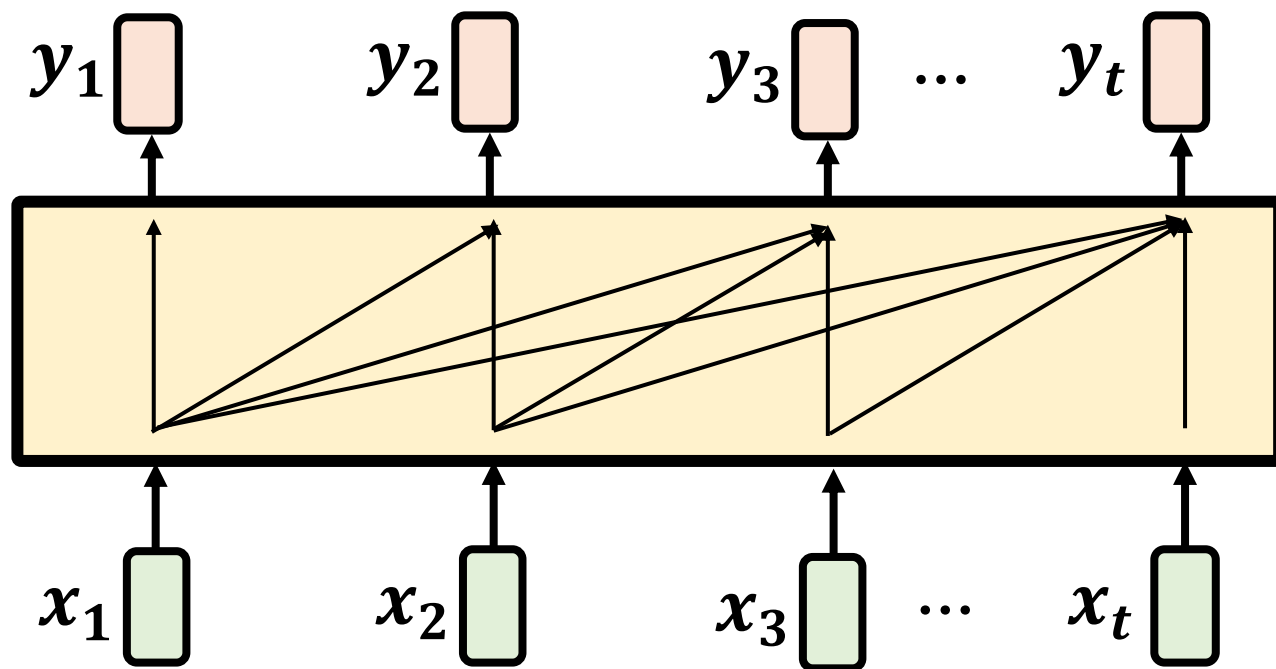
$$\mathbf{k}_t = W_k \mathbf{x}_t$$

- Linear Attention is generalized RNN w/o “Reflection”  $f_{A,t}$
- Linear Attention is Self-attention w/o Softmax

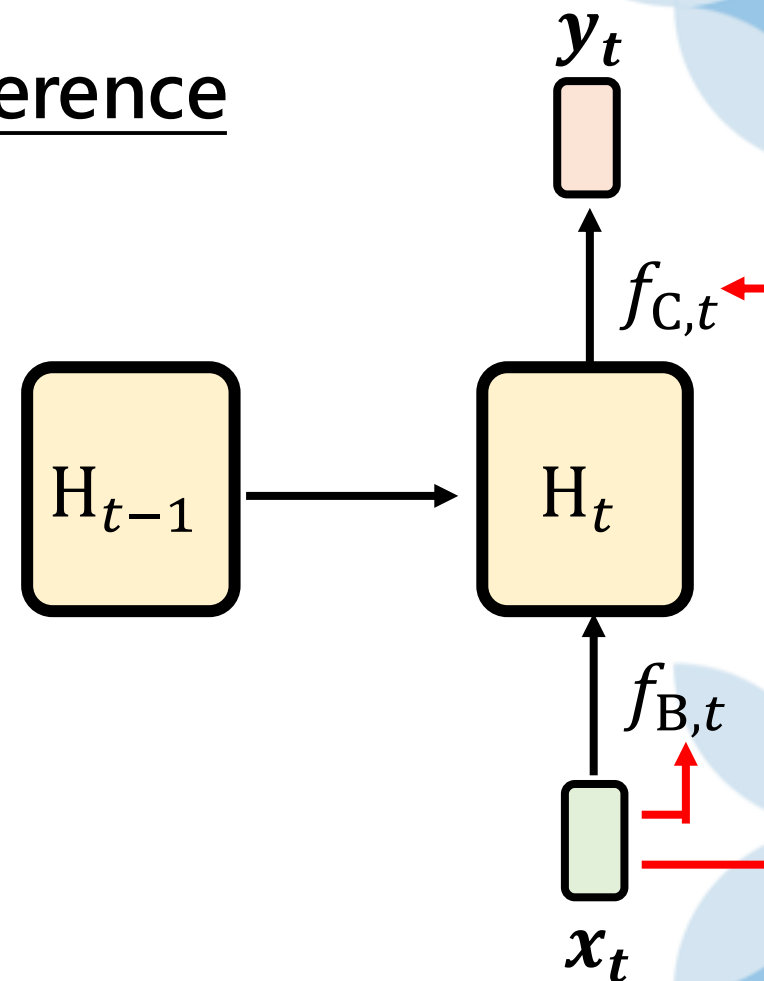
# Linear Attention

at training, it works like self-attention  
at inference, it works like RNN

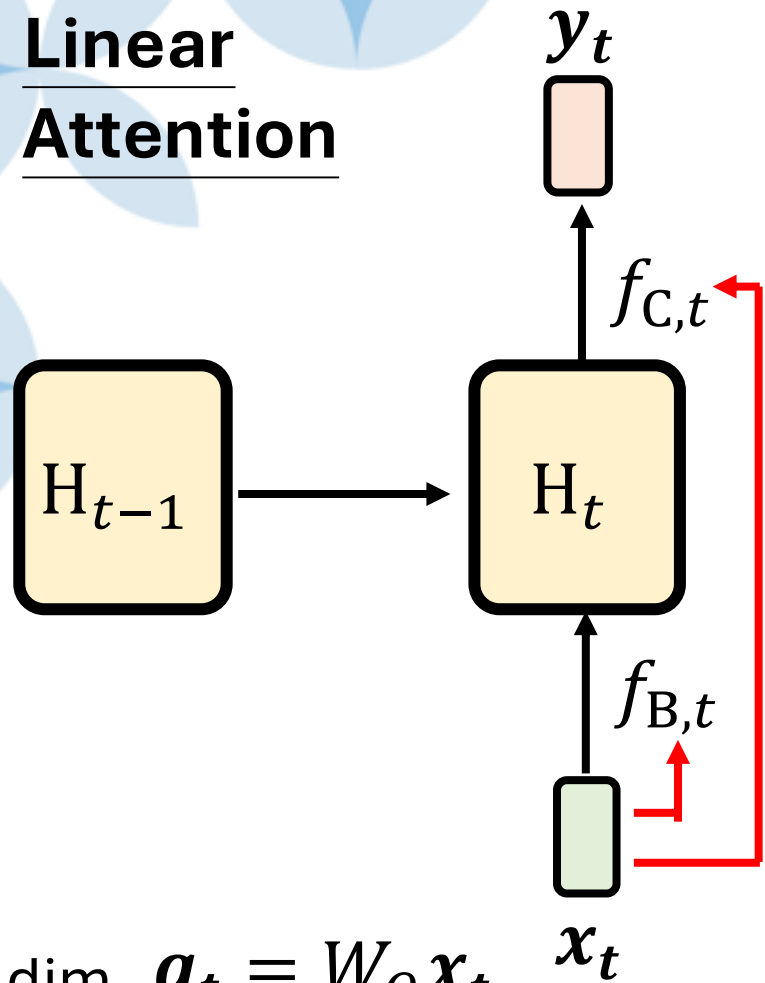
## Training



## Inference



# Linear Attention



$$d \text{ dim } \mathbf{q}_t = W_Q \mathbf{x}_t$$

$$d \text{ dim } \mathbf{k}_t = W_k \mathbf{x}_t$$

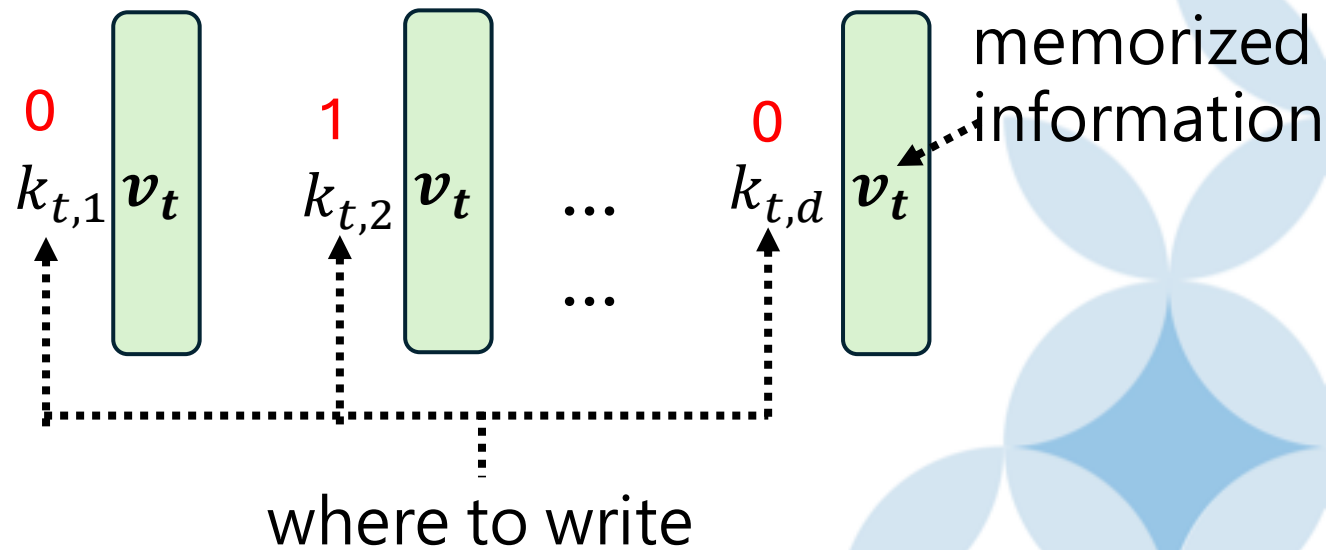
$$d' \text{ dim } \mathbf{v}_t = W_v \mathbf{x}_t$$

$$H_t = H_{t-1} + f_{B,t}(\mathbf{x}_t) \quad f_{B,t}(\mathbf{x}_t) = \mathbf{v}_t \mathbf{k}_t^T$$

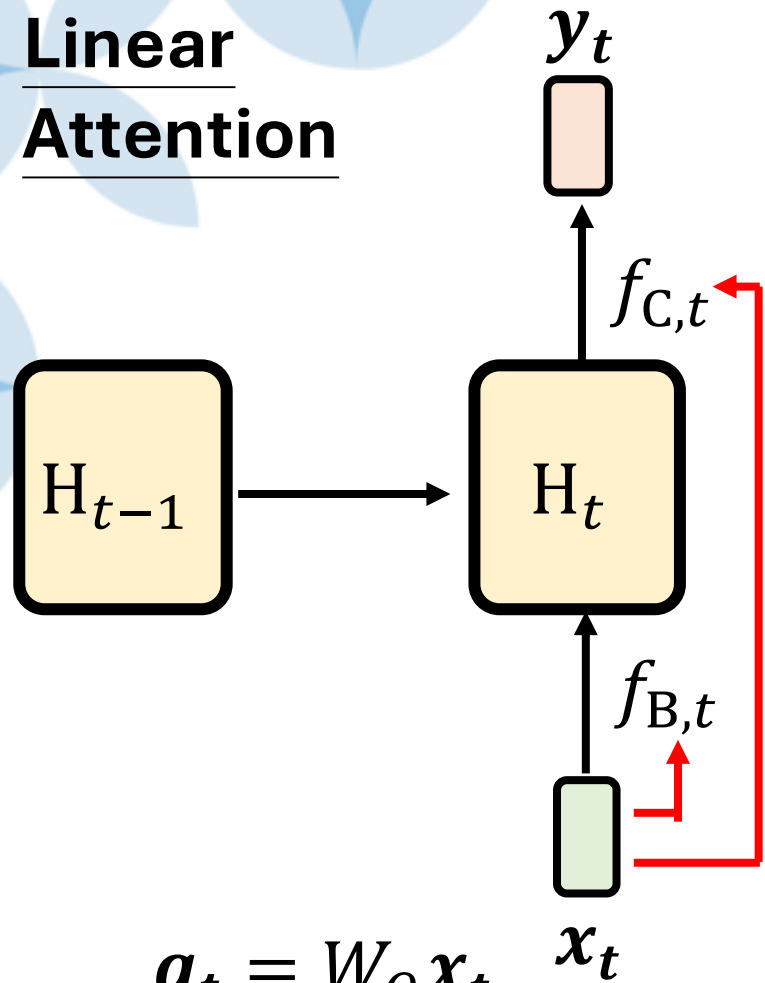
$$\mathbf{y}_t = f_{C,t}(H_t) \quad f_{C,t}(H_t) = H_t \mathbf{q}_t$$

$$H_t = H_{t-1} + d' \begin{matrix} d \\ \mathbf{v}_t \mathbf{k}_t^T \end{matrix}$$

Write  $\mathbf{v}_t$  into 2nd column of H



# Linear Attention

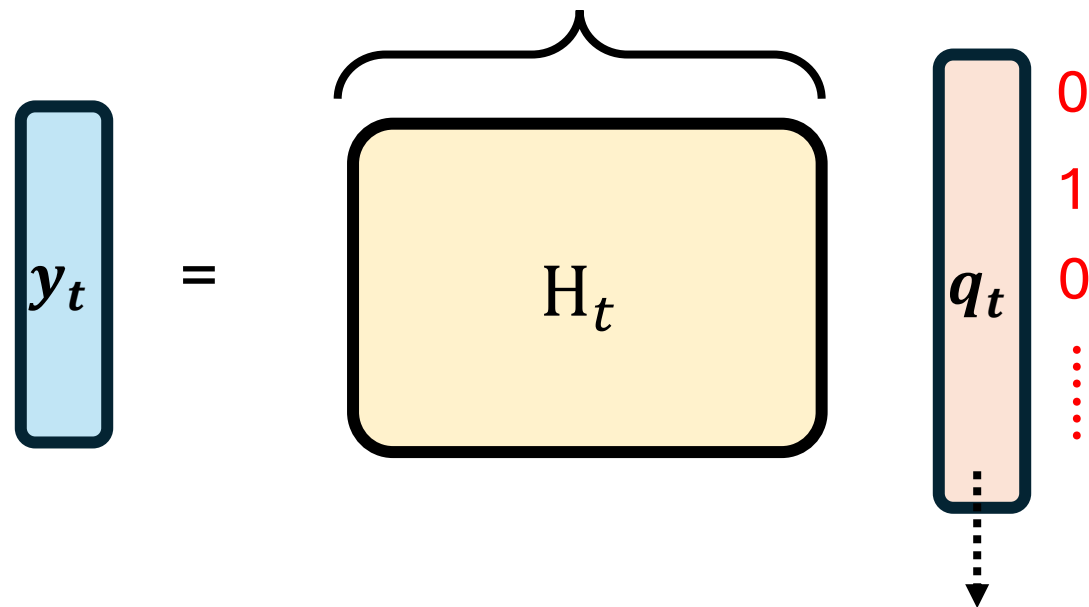


$$\begin{aligned} d & \text{ dim} & \mathbf{q}_t &= W_Q \mathbf{x}_t \\ d & \text{ dim} & \mathbf{k}_t &= W_k \mathbf{x}_t \\ d' & \text{ dim} & \mathbf{v}_t &= W_v \mathbf{x}_t \end{aligned}$$

$$H_t = H_{t-1} + f_{B,t}(\mathbf{x}_t) \quad f_{B,t}(\mathbf{x}_t) = \mathbf{v}_t \mathbf{k}_t^T$$

$$\mathbf{y}_t = f_{C,t}(H_t) \quad f_{C,t}(H_t) = H_t \mathbf{q}_t$$

Different information into different Columns



How much information to take from each column

# This is not a new idea

---

## Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention

---

Angelos Katharopoulos<sup>1,2</sup> Apoorv Vyas<sup>1,2</sup> Nikolaos Pappas<sup>3</sup> François Fleuret<sup>2,4\*</sup>

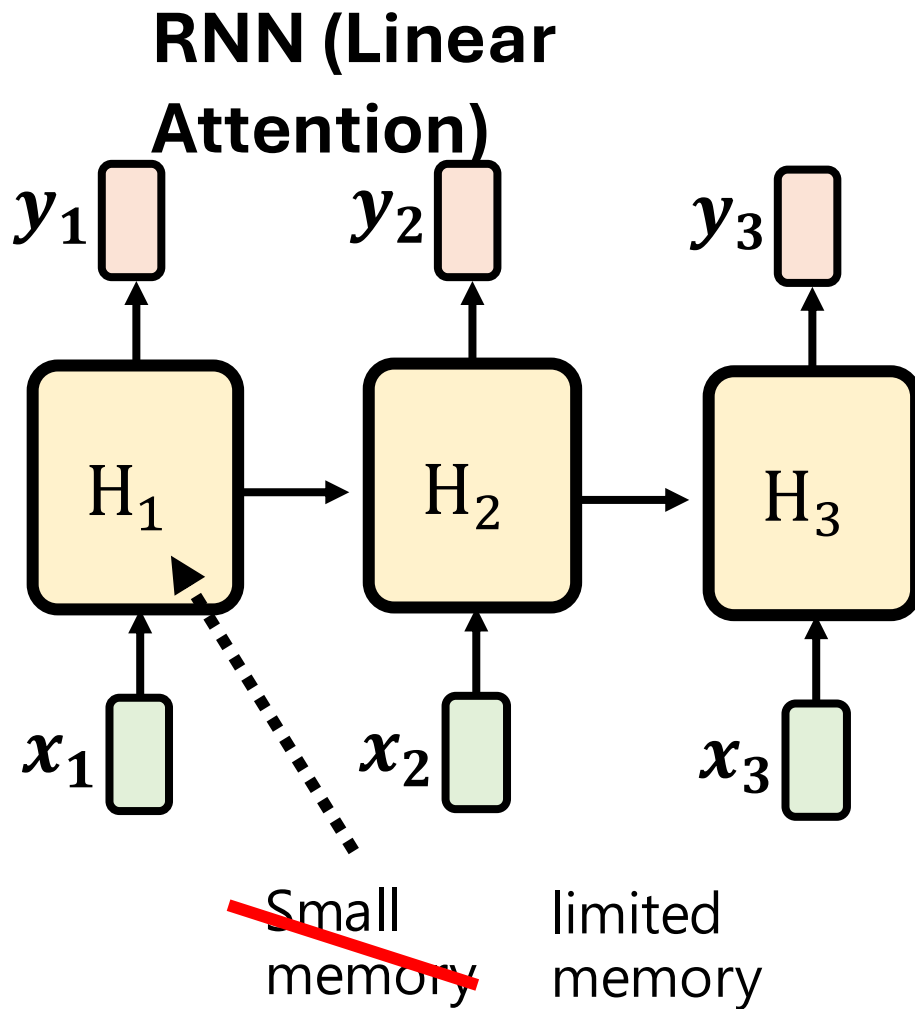
### Abstract

Transformers achieve remarkable performance in several tasks but due to their quadratic complexity, with respect to the input's length, they are prohibitively slow for very long sequences. To address this limitation, we express the self-attention as a linear dot-product of kernel feature maps and make use of the associativity property of matrix

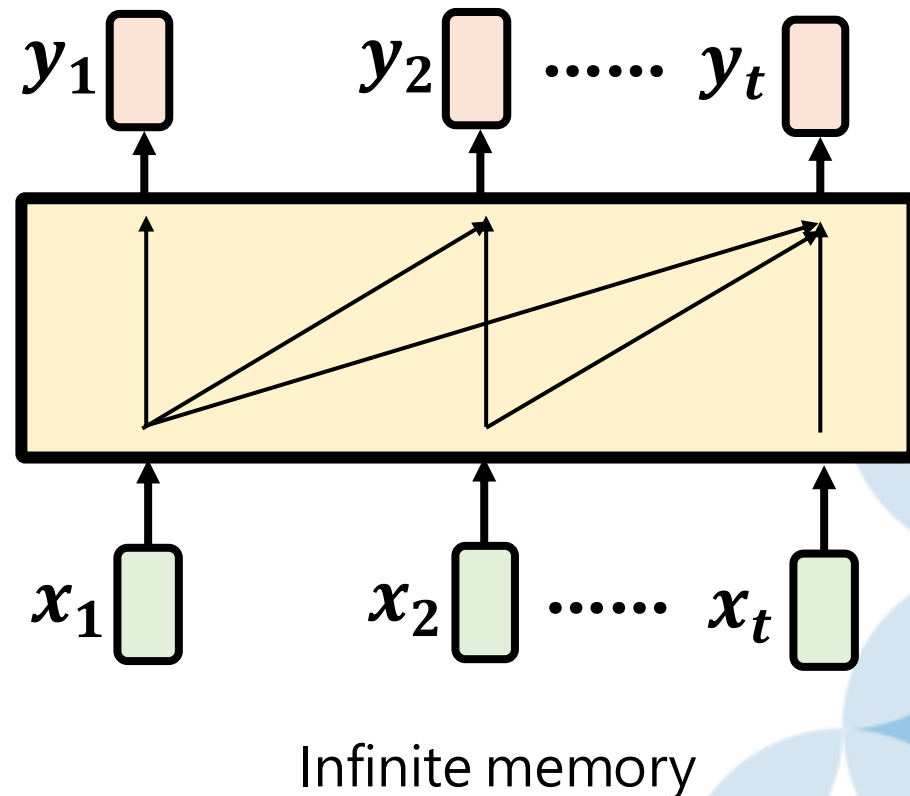
by the global receptive field of self-attention, which processes contexts of  $N$  inputs with a quadratic memory and time complexity  $\mathcal{O}(N^2)$ . As a result, in practice transformers are slow to train and their context is *limited*. This disrupts temporal coherence and hinders the capturing of long-term dependencies. Dai et al. (2019) addressed the latter by attending to memories from previous contexts albeit at the expense of computational efficiency.

Aug 2020

RNN (Linear Attention) < Transformer (Self-attention with Softmax) ?

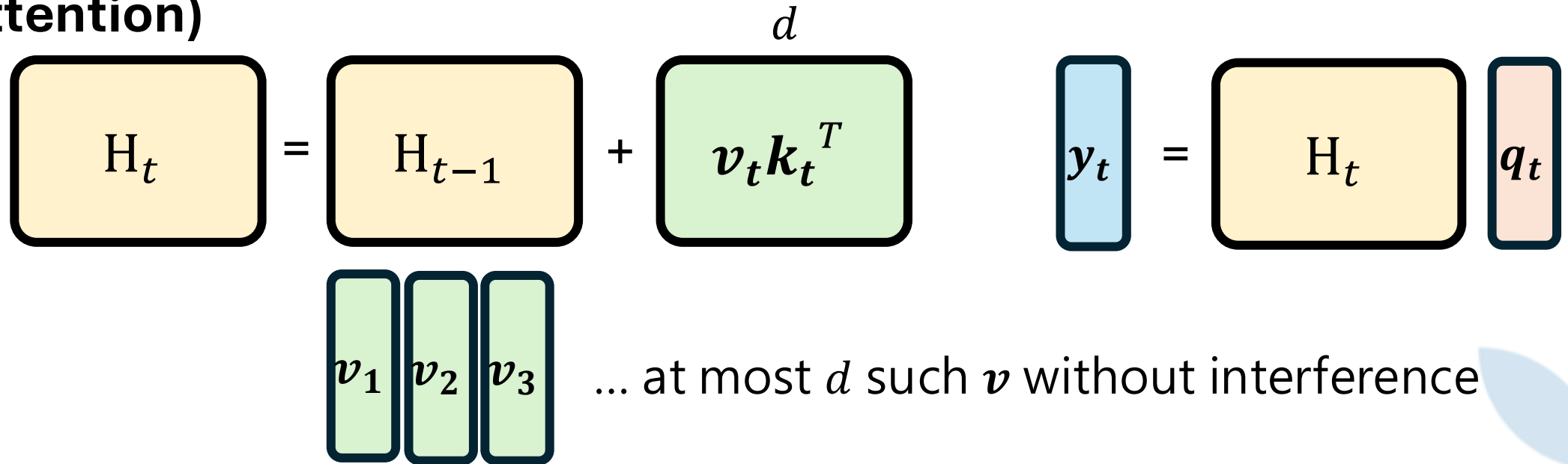


**Transformer (Self-attention with softmax)**



# RNN (Linear Attention) < Transformer (Self-attention with Softmax) ?

## RNN (Linear Attention)

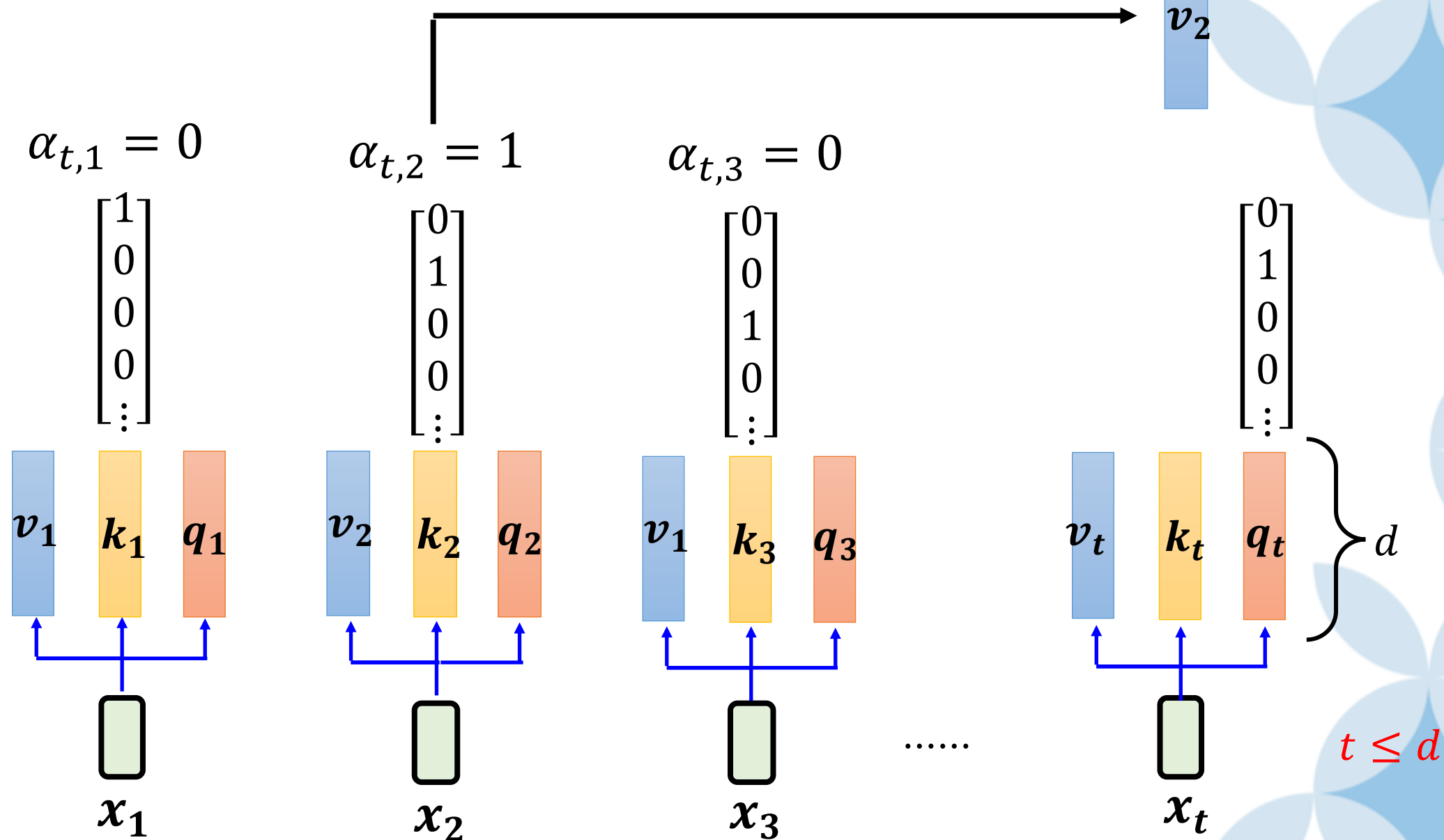


$$k_1^T = [1 \quad 0 \quad \dots]$$

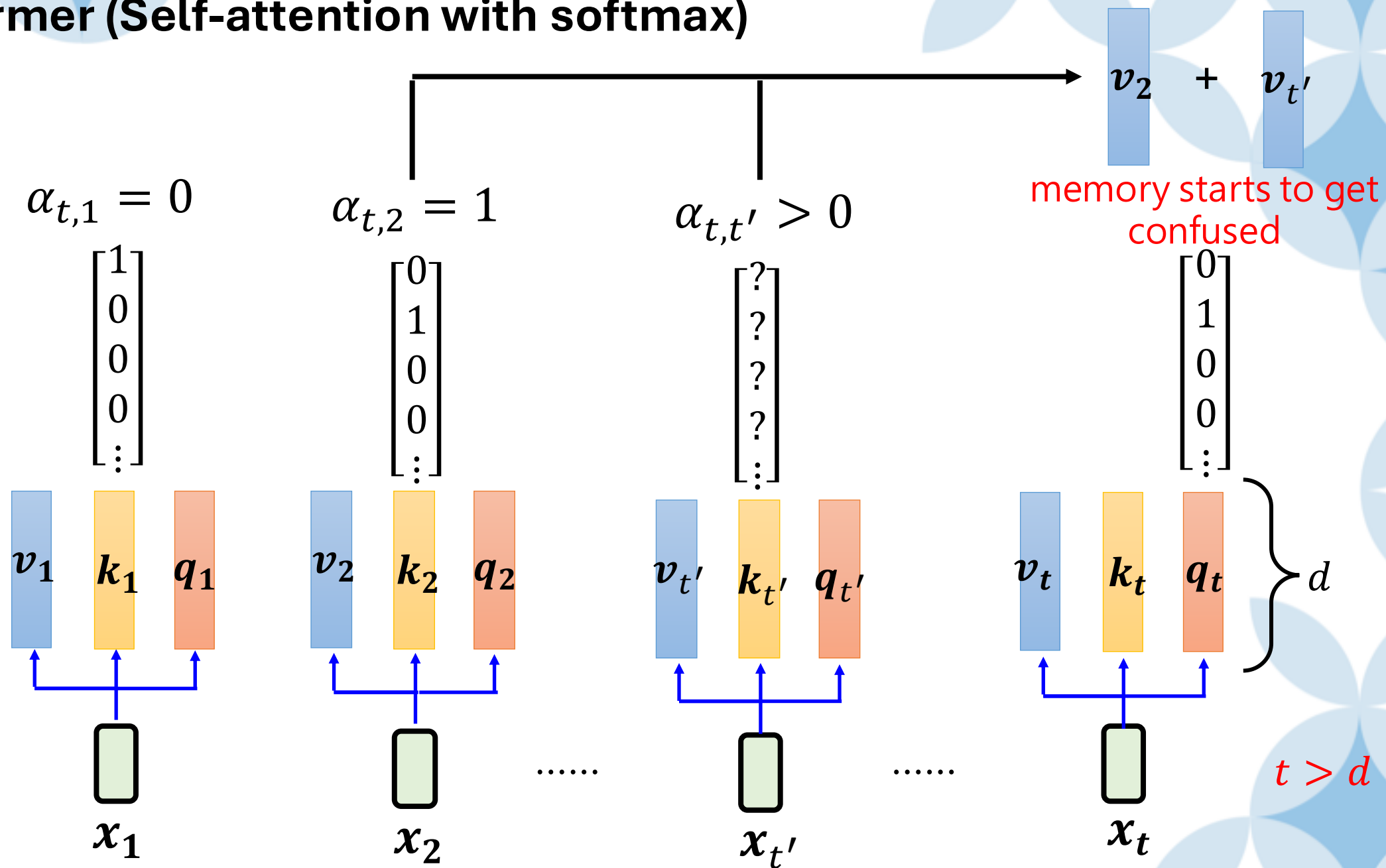
$$k_2^T = [0 \quad 1 \quad \dots]$$

$$k_3^T = [0 \quad 0 \quad 1 \quad \dots]$$

# Transformer (Self-attention with softmax)



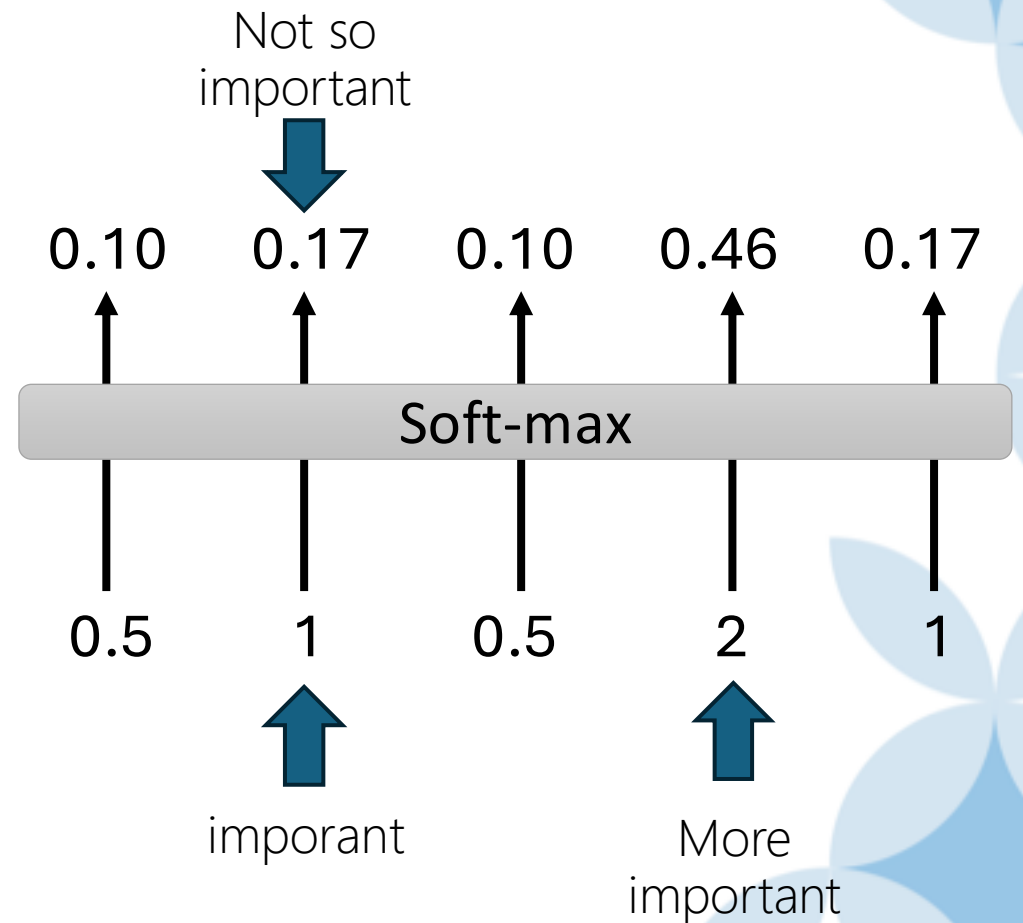
# Transformer (Self-attention with softmax)



# RNN (Linear Attention) < Transformer (Self-attention with Softmax) ?

$$H_t = H_{t-1} + f_{B,t}(x_t)$$

Linear Attention memory unchanged



# + Reflection: gradually forget

## Linear Attention

$$\mathbf{H}_t = \mathbf{H}_{t-1} + \mathbf{v}_t \mathbf{k}_t^T$$

$$\mathbf{y}_t = \mathbf{H}_t \mathbf{q}_t$$

$$\mathbf{v}_t = W_v \mathbf{x}_t$$

$$\mathbf{k}_t = W_k \mathbf{x}_t$$

$$\mathbf{q}_t = W_Q \mathbf{x}_t$$

## Retention Network (RetNet)

$$\mathbf{H}_t = \gamma \mathbf{H}_{t-1} + \mathbf{v}_t \mathbf{k}_t^T$$

$$\mathbf{y}_t = \mathbf{H}_t \mathbf{q}_t$$

$$\mathbf{v}_t = W_v \mathbf{x}_t$$

$$\mathbf{k}_t = W_k \mathbf{x}_t$$

$$\mathbf{q}_t = W_Q \mathbf{x}_t$$



# + Reflection: forget as needed

<https://arxiv.org/abs/2405.05254>

Retention Network (RetNet)

$$\mathbf{H}_t = \gamma \mathbf{H}_{t-1} + \mathbf{v}_t \mathbf{k}_t^T$$

$$\mathbf{y}_t = \mathbf{H}_t \mathbf{q}_t$$

$$\mathbf{v}_t = W_v \mathbf{x}_t$$

$$\mathbf{k}_t = W_k \mathbf{x}_t$$

$$\mathbf{q}_t = W_Q \mathbf{x}_t$$

Gated Retention

$$\mathbf{H}_t = \gamma_t \mathbf{H}_{t-1} + \mathbf{v}_t \mathbf{k}_t^T$$

$$\mathbf{y}_t = \mathbf{H}_t \mathbf{q}_t$$

$$\mathbf{v}_t = W_v \mathbf{x}_t$$

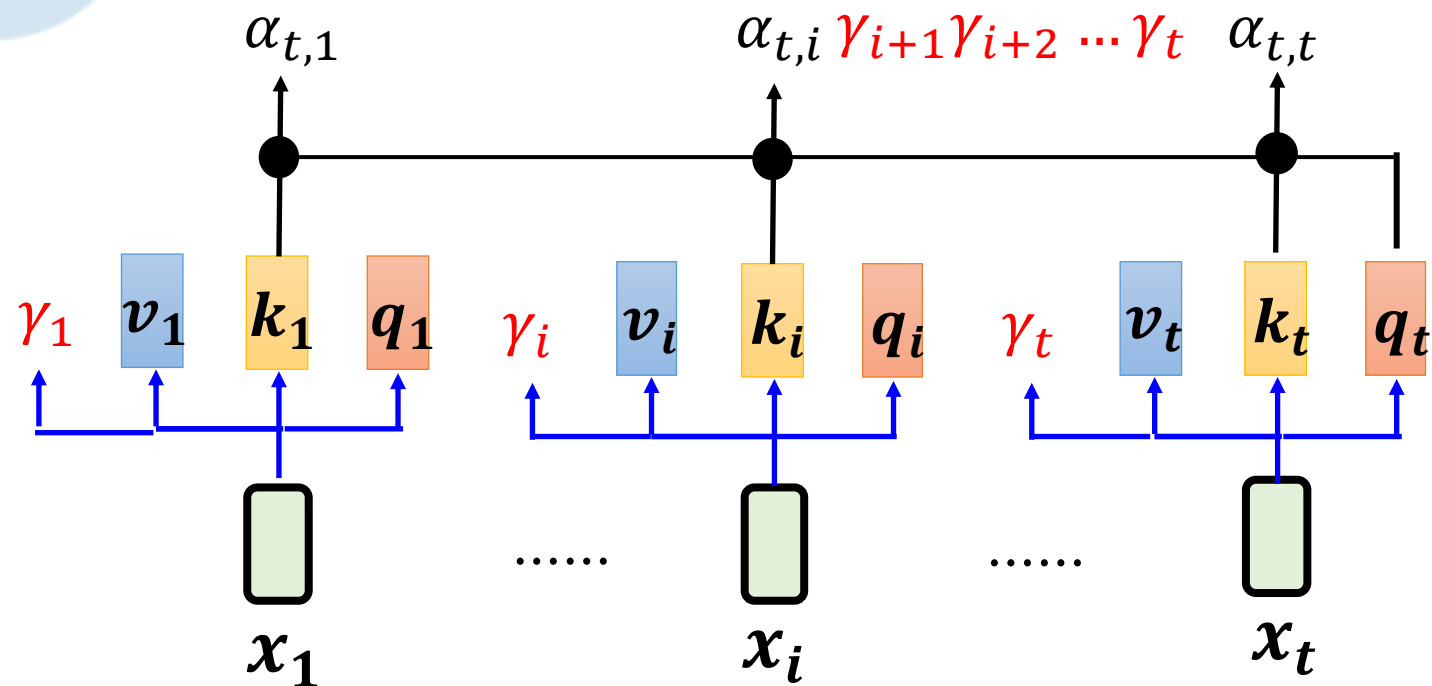
$$\mathbf{k}_t = W_k \mathbf{x}_t$$

$$\mathbf{q}_t = W_Q \mathbf{x}_t$$

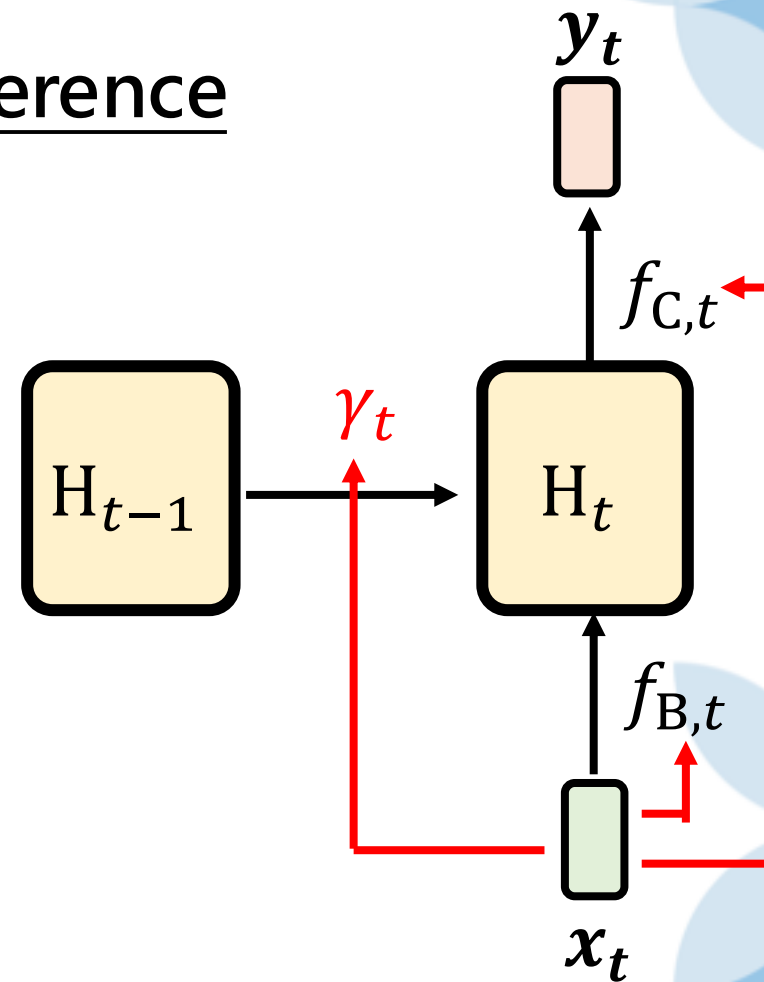
$$\gamma_t = \text{sigmoid}(W_\gamma \mathbf{x}_t)$$

# + Reflection: gradually forget

## Training



## Inference



# More complex Reflection

$$H_t = G_t \odot H_{t-1} + v_t k_t^T$$

$$G_t = e_t s_t^T$$

$$G_t = \mathbf{1} s_t^T$$

$$\mathbf{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$s_t^T = [0 \quad 1 \quad 0.1 \quad \dots \dots]$$

$$\mathbf{1} \odot H_{t-1}$$

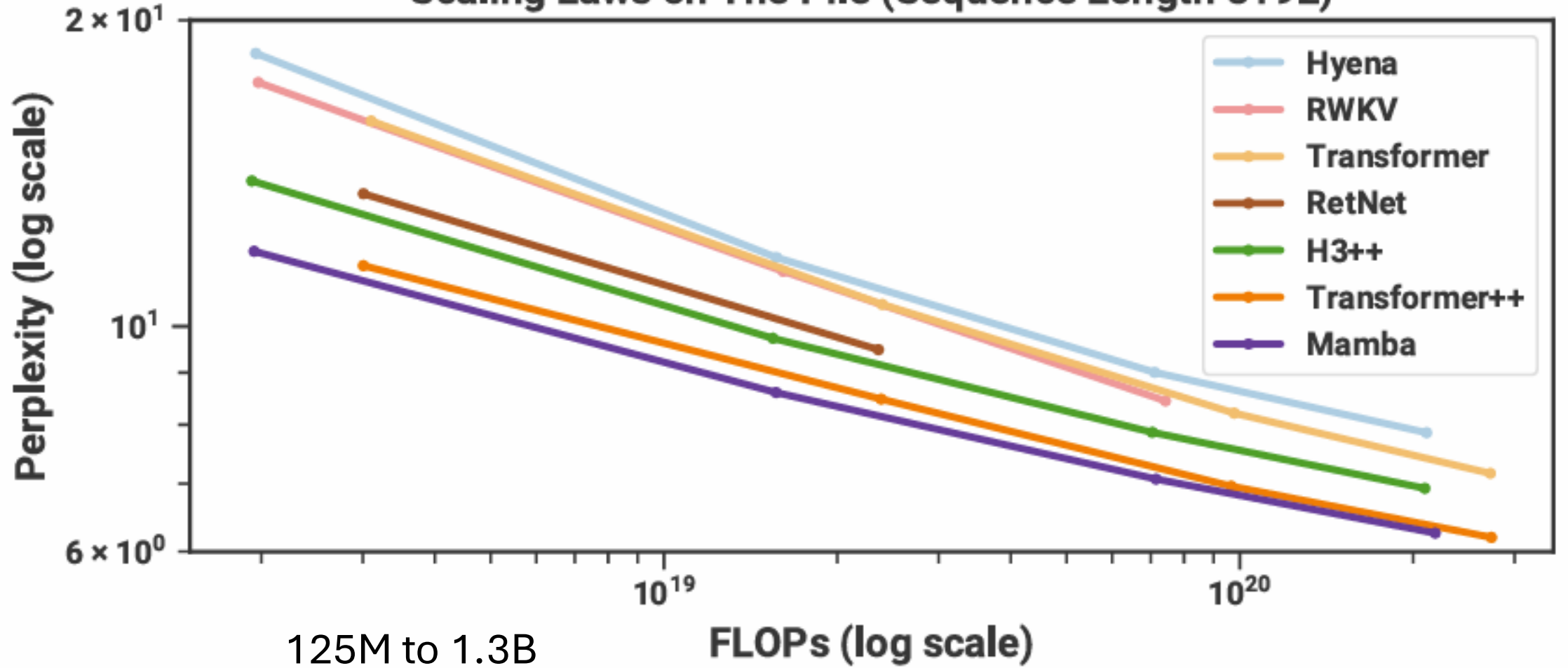
$s_{t,1}$	$s_{t,2}$	$s_{t,3}$	$s_{t,d}$
$s_{t,1}$	$s_{t,2}$	$s_{t,3}$	$s_{t,d}$
$\vdots$	$\vdots$	$\vdots$	$\dots \vdots$
$s_{t,1}$	$s_{t,2}$	$s_{t,3}$	$s_{t,d}$

0
1
0.1

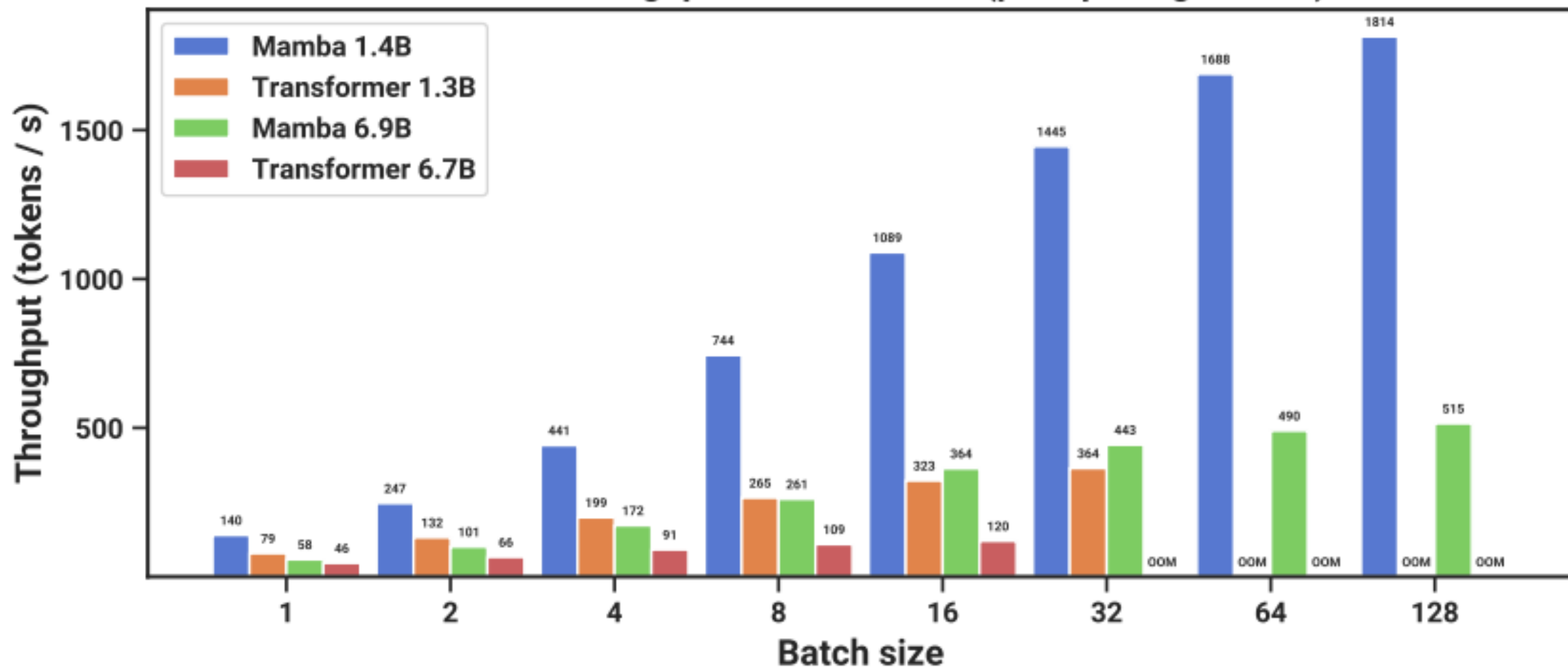
remove
keep
reduce

Model	Recurrence	Memory read-out
Linear Attention [48, 47]	$\mathbf{S}_t = \mathbf{S}_{t-1} + \mathbf{v}_t \mathbf{k}_t^\top$	$\mathbf{o}_t = \mathbf{S}_t \mathbf{q}_t$
+ Kernel	$\mathbf{S}_t = \mathbf{S}_{t-1} + \mathbf{v}_t \phi(\mathbf{k}_t)^\top$	$\mathbf{o}_t = \mathbf{S}_t \phi(\mathbf{q}_t)$
+ Normalization	$\mathbf{S}_t = \mathbf{S}_{t-1} + \mathbf{v}_t \phi(\mathbf{k}_t)^\top, \mathbf{z}_t = \mathbf{z}_{t-1} + \phi(\mathbf{k}_t)$	$\mathbf{o}_t = \mathbf{S}_t \phi(\mathbf{q}_t) / (\mathbf{z}_t^\top \phi(\mathbf{q}_t))$
DeltaNet [101]	$\mathbf{S}_t = \mathbf{S}_{t-1} (\mathbf{I} - \beta_t \mathbf{k}_t \mathbf{k}_t^\top) + \beta_t \mathbf{v}_t \mathbf{k}_t^\top$	$\mathbf{o}_t = \mathbf{S}_t \mathbf{q}_t$
Gated RFA [81]	$\mathbf{S}_t = g_t \mathbf{S}_{t-1} + (1 - g_t) \mathbf{v}_t \mathbf{k}_t^\top, \mathbf{z}_t = g_t \mathbf{z}_{t-1} + (1 - g_t) \mathbf{k}_t$	$\mathbf{o}_t = \mathbf{S}_t \mathbf{q}_t / (\mathbf{z}_t^\top \mathbf{q}_t)$
S4 [32, 106]	$\mathbf{S}_t = \mathbf{S}_{t-1} \odot \exp(-(\boldsymbol{\alpha} \mathbf{1}^\top) \odot \exp(\mathbf{A})) + \mathbf{B} \odot (\mathbf{v}_t \mathbf{1}^\top)$	$\mathbf{o}_t = (\mathbf{S}_t \odot \mathbf{C}) \mathbf{1} + \mathbf{d} \odot \mathbf{v}_t$
ABC [82]	$\mathbf{S}_t^k = \mathbf{S}_{t-1}^k + \mathbf{k}_t \phi_t^\top, \mathbf{S}_t^v = \mathbf{S}_{t-1}^v + \mathbf{v}_t \phi_t^\top$	$\mathbf{o}_t = \mathbf{S}_t^v \text{softmax}(\mathbf{S}_t^k \mathbf{q}_t)$
DFW [63]	$\mathbf{S}_t = \mathbf{S}_{t-1} \odot (\beta_t \boldsymbol{\alpha}_t^\top) + \mathbf{v}_t \mathbf{k}_t^\top$	$\mathbf{o}_t = \mathbf{S}_t \mathbf{q}_t$
RetNet [108]	$\mathbf{S}_t = \gamma \mathbf{S}_{t-1} + \mathbf{v}_t \mathbf{k}_t^\top$	$\mathbf{o}_t = \mathbf{S}_t \mathbf{q}_t$
Mamba [31]	$\mathbf{S}_t = \mathbf{S}_{t-1} \odot \exp(-(\boldsymbol{\alpha}_t \mathbf{1}^\top) \odot \exp(\mathbf{A})) + (\boldsymbol{\alpha}_t \odot \mathbf{v}_t) \mathbf{k}_t^\top$	$\mathbf{o}_t = \mathbf{S}_t \mathbf{q}_t + \mathbf{d} \odot \mathbf{v}_t$
GLA [124]	$\mathbf{S}_t = \mathbf{S}_{t-1} \odot (\mathbf{1} \boldsymbol{\alpha}_t^\top) + \mathbf{v}_t \mathbf{k}_t^\top = \mathbf{S}_{t-1} \text{Diag}(\boldsymbol{\alpha}_t) + \mathbf{v}_t \mathbf{k}_t^\top$	$\mathbf{o}_t = \mathbf{S}_t \mathbf{q}_t$
RWKV-6 [79]	$\mathbf{S}_t = \mathbf{S}_{t-1} \text{Diag}(\boldsymbol{\alpha}_t) + \mathbf{v}_t \mathbf{k}_t^\top$	$\mathbf{o}_t = (\mathbf{S}_{t-1} + (\mathbf{d} \odot \mathbf{v}_t) \mathbf{k}_t^\top) \mathbf{q}_t$
HGRN-2 [92]	$\mathbf{S}_t = \mathbf{S}_{t-1} \text{Diag}(\boldsymbol{\alpha}_t) + \mathbf{v}_t (\mathbf{1} - \boldsymbol{\alpha}_t)^\top$	$\mathbf{o}_t = \mathbf{S}_t \mathbf{q}_t$
mLSTM [9]	$\mathbf{S}_t = f_t \mathbf{S}_{t-1} + i_t \mathbf{v}_t \mathbf{k}_t^\top, \mathbf{z}_t = f_t \mathbf{z}_{t-1} + i_t \mathbf{k}_t$	$\mathbf{o}_t = \mathbf{S}_t \mathbf{q}_t / \max\{1,  \mathbf{z}_t^\top \mathbf{q}_t \}$
Mamba-2 [19]	$\mathbf{S}_t = \gamma_t \mathbf{S}_{t-1} + \mathbf{v}_t \mathbf{k}_t^\top$	$\mathbf{o}_t = \mathbf{S}_t \mathbf{q}_t$
GSA [131]	$\mathbf{S}_t^k = \mathbf{S}_{t-1}^k \text{Diag}(\boldsymbol{\alpha}_t) + \mathbf{k}_t \phi_t^\top, \mathbf{S}_t^v = \mathbf{S}_{t-1}^v \text{Diag}(\boldsymbol{\alpha}_t) + \mathbf{v}_t \phi_t^\top$	$\mathbf{o}_t = \mathbf{S}_t^v \text{softmax}(\mathbf{S}_t^k \mathbf{q}_t)$
Gated DeltaNet [125]	$\mathbf{S}_t = \mathbf{S}_{t-1} \left( \boldsymbol{\alpha}_t (\mathbf{I} - \beta_t \mathbf{k}_t \mathbf{k}_t^\top) \right) + \beta_t \mathbf{v}_t \mathbf{k}_t^\top$	$\mathbf{o}_t = \mathbf{S}_t \mathbf{q}_t$

Scaling Laws on The Pile (Sequence Length 8192)



Inference throughput on A100 80GB (prompt length 2048)



# DeltaNet

<https://arxiv.org/abs/2406.06484>

$$H_t = H_{t-1}(I - \beta_t \mathbf{k}_t \mathbf{k}_t^T) + \beta_t \mathbf{v}_t \mathbf{k}_t^T$$

$$H_t = H_{t-1} + \mathbf{v}_t \mathbf{k}_t^T \longleftarrow \text{Gradient Descent, with different } L_t$$

$$H_t = H_{t-1} - \mathbf{v}_{t,old} \mathbf{k}_t^T + \mathbf{v}_t \mathbf{k}_t^T \quad \mathbf{v}_{t,old} = H_{t-1} \mathbf{k}_t$$

$$H_t = H_{t-1} - \beta_t \mathbf{v}_{t,old} \mathbf{k}_t^T + \beta_t \mathbf{v}_t \mathbf{k}_t^T$$

$$H_t = H_{t-1} - \beta_t H_{t-1} \mathbf{k}_t \mathbf{k}_t^T + \beta_t \mathbf{v}_t \mathbf{k}_t^T$$

**Gradient  
Descent**

$$H_t = H_{t-1} - \beta_t (\mathbf{H}_{t-1} \mathbf{k}_t - \mathbf{v}_t) \mathbf{k}_t^T$$

Parameter after update    before update    learning rate    gradient

$$L_t(H) = \frac{1}{2} \|H \mathbf{k}_t - \mathbf{v}_t\|^2$$

$$\nabla L_t(H_{t-1})$$

Update  $H$  s.t. info. retrieved using  $\mathbf{k}_t$  is close to  $\mathbf{v}_t$

**Titans: Learning to Memorize at Test Time**

<https://arxiv.org/abs/2501.00663>

Name	Modality	Affiliations	Sizes	Access Link
Mamba 1&2	Language	Carnegie Mellon University & Princeton University	130M-2.8B	1
Falcon Mamba 7B	Language	Technology Innovation Institute	7B	2
Mistral 7B	Language	Mistral AI & NVIDIA	7B	3
Jamba	Language	AI21 Lab	12B/52B	4
Vision Mamba	Vision	Huazhong University of Science and Technology	7M-98M	5
VideoMamba	Video	OpenGVLab, Shanghai AI Laboratory	28M-392M	6
Codestral Mamba	Code	Mistral AI	7B, 22B	7

1. <https://github.com/state-spaces/mamba>
2. <https://huggingface.co/tiiuae/falcon-mamba-7b>
3. <https://huggingface.co/mistralai/Mistral-7B-v0.1>
4. <https://huggingface.co/ai21labs/Jamba-v0.1>
5. <https://huggingface.co/hustvl/Vim-base-midclstok>
6. <https://huggingface.co/OpenGVLab/VideoMamba>
7. <https://mistral.ai/news/codestral-mamba/>

<https://arxiv.org/abs/2408.01129>

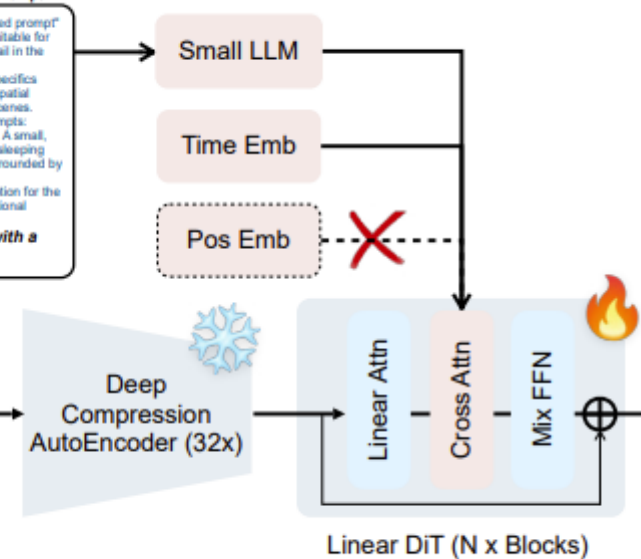
Minimax-01

<https://arxiv.org/abs/2501.08313>

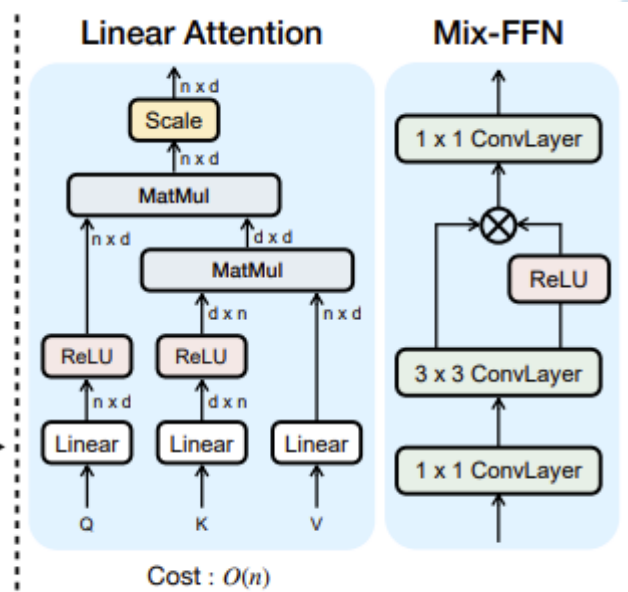


<Complex Human Instruction> <User Prompt>

Given a user prompt, generate an "Enhanced prompt" that provides detailed visual descriptions suitable for image generation. Evaluate the level of detail in the user prompt:  
 - If the prompt is simple, focus on adding specifics about colors, shapes, sizes, textures, and spatial relationships to create vivid and concrete scenes.  
 Examples of how to transform or refine prompts:  
 - User Prompt: A cat sleeping -> Enhanced: A small, fluffy white cat curled up in a round shape, sleeping peacefully on a warm sunny windowsill, surrounded by pots of blooming red flowers.  
 Please generate only the enhanced description for the prompt below and avoid including any additional commentary or evaluations:  
 User Prompt: A cyberpunk cat with a neon sign that says "SANA".



Linear DiT (N x Blocks)



Cost :  $O(n)$

(a). Architecture overview of our Sana.

(b). Linear DiT Module.

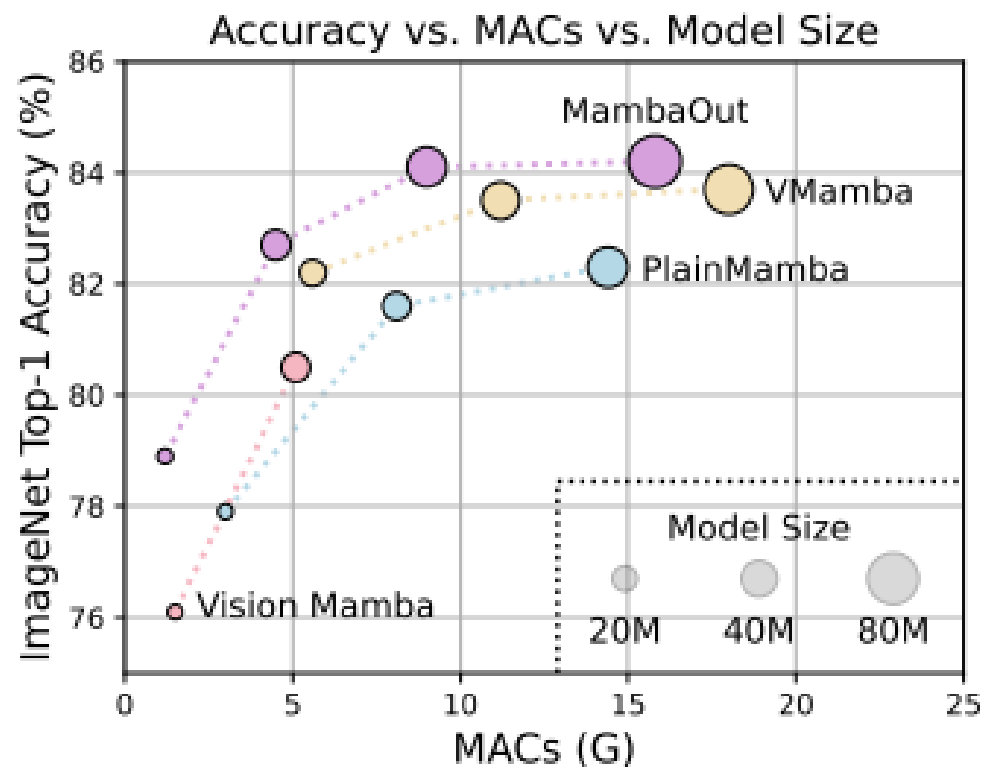
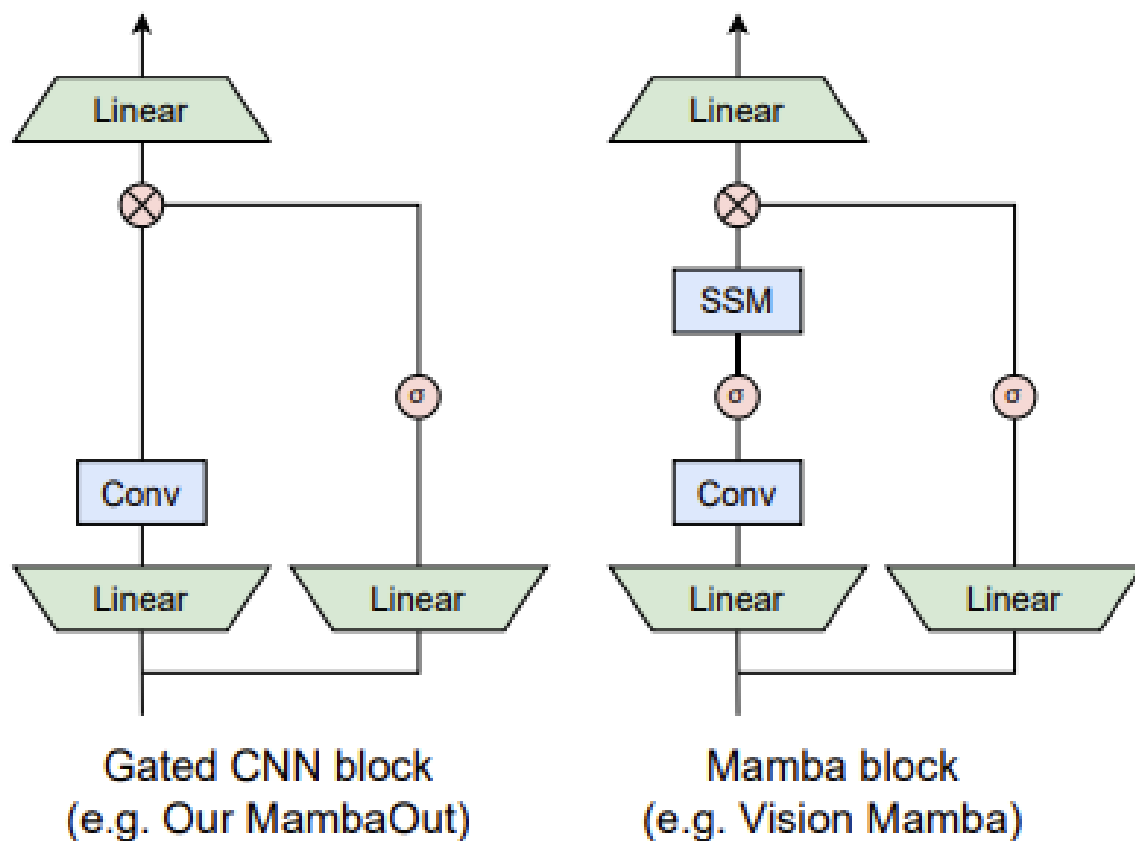
<https://arxiv.org/abs/2410.10629>

# MambaOut: Do We Really Need Mamba for Vision?

<https://arxiv.org/abs/2405.07992>

*In memory of Kobe Bryant*

*“What can I say, Mamba out.” — Kobe Bryant’s NBA farewell speech in 2016.*



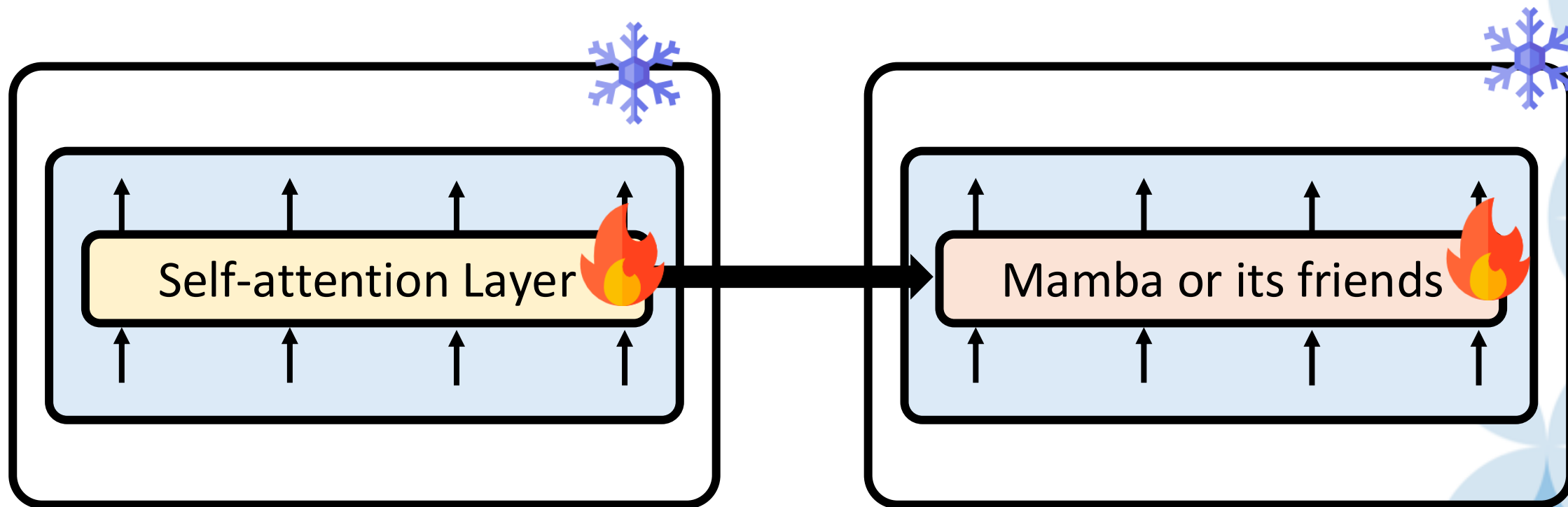
# Do not train from scratch

Low-rank Linear Conversion via Attention Transfer (LoLCATs), <https://arxiv.org/abs/2410.10254>

The Mamba in the Llama, <https://arxiv.org/abs/2408.15237>

Transformers to SSMs, <https://arxiv.org/abs/2408.10189>

Linger, <https://arxiv.org/abs/2503.01496>



# Is Attention All You Need?



**Current Status: Yes**

Time Remaining: 302d 20h 44m 32s

# 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.*

## 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 🧡



Thanks!



Questions?

