

Recurrent Neural Networks

deep learning 3

Raoul Grouls, 13 Mei 2025

Neural networks

Neural networks

$$\sigma(wx + b)$$

Neural networks

$$\sigma(wx + b)$$

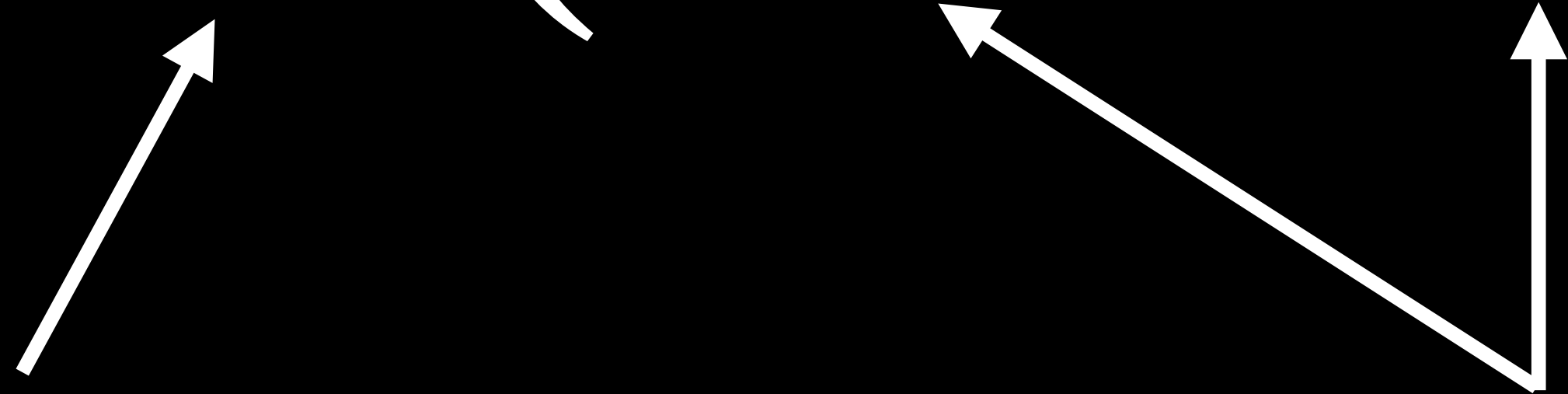
↑
Input

Neural networks

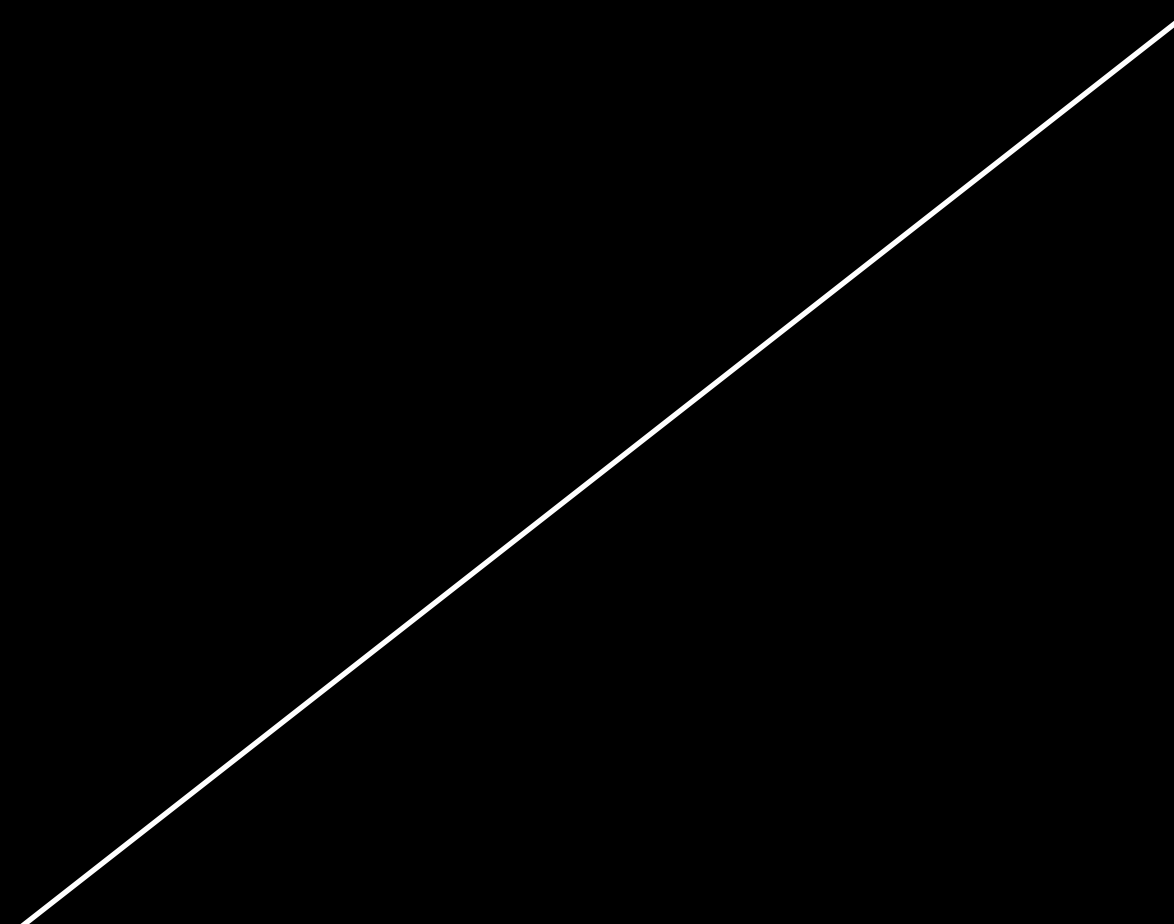
$$\sigma(wx + b)$$

Non-linear transformation

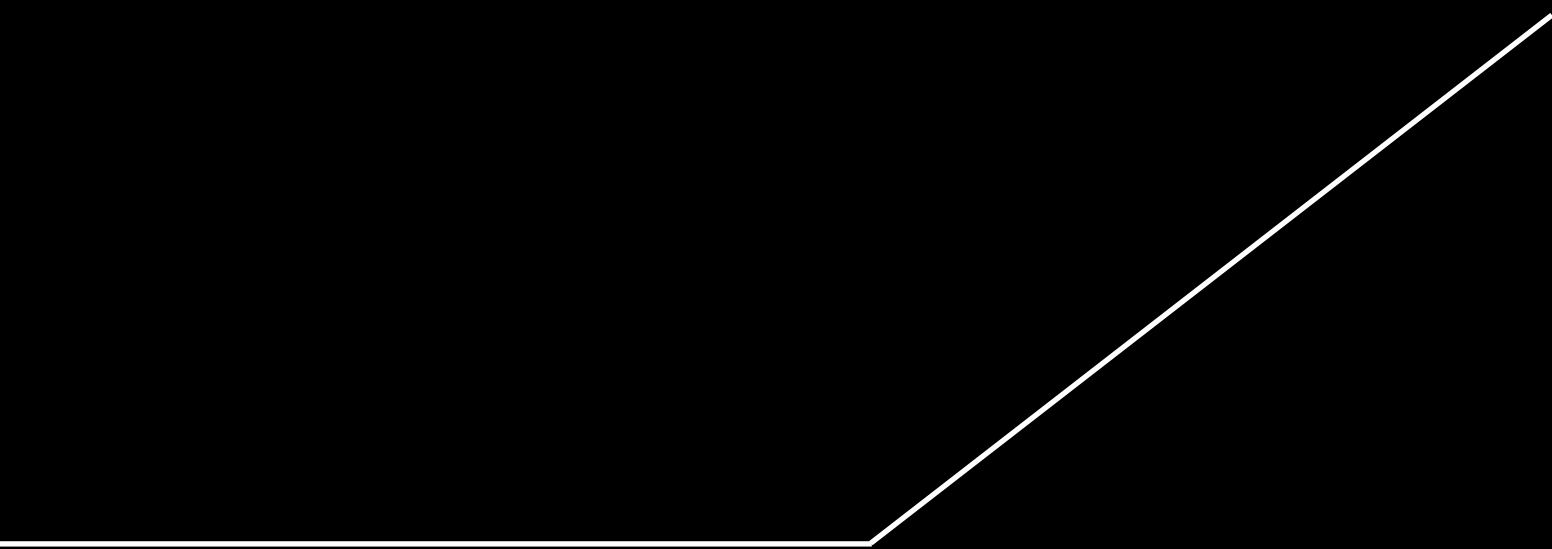
Linear transformation



Neural networks



Linear



Nonlinear

Neural networks

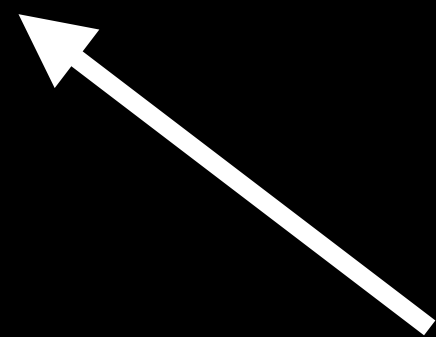
$$\sigma(wx + b)$$

Learnable

A diagram consisting of two white arrows originating from a single point below the word 'Learnable'. One arrow points diagonally up and to the left, terminating at the variable 'w' in the equation $\sigma(wx + b)$. The other arrow points diagonally up and to the right, terminating at the variable 'b' in the same equation.

Neural networks

$$\hat{y} = \sigma(wx + b)$$



Prediction

Neural networks

$$\hat{y} = \sigma(wx + b)$$

$$z = (y - \hat{y})^2$$

Change w and b such that z is minimal

Neural networks

$$\hat{y} = \sigma(wx + b)$$

$$\left. \begin{array}{cc} \frac{\partial z}{\partial w} & \frac{\partial z}{\partial b} \end{array} \right\} \text{Gradient}$$

How much do we need to change w and b

Neural networks

$$w \leftarrow w + \eta \frac{\partial z}{\partial w}$$

$$b \leftarrow b + \eta \frac{\partial z}{\partial b}$$

Update the weights

Neural networks

$$\left. \begin{aligned} w &\leftarrow w + \eta \frac{\partial z}{\partial w} \\ b &\leftarrow b + \eta \frac{\partial z}{\partial b} \end{aligned} \right\} \text{Optimizer}$$

Neural networks

$$w \leftarrow w + \eta \frac{\partial z}{\partial w}$$

Learning rate

$$b \leftarrow b + \eta \frac{\partial z}{\partial b}$$

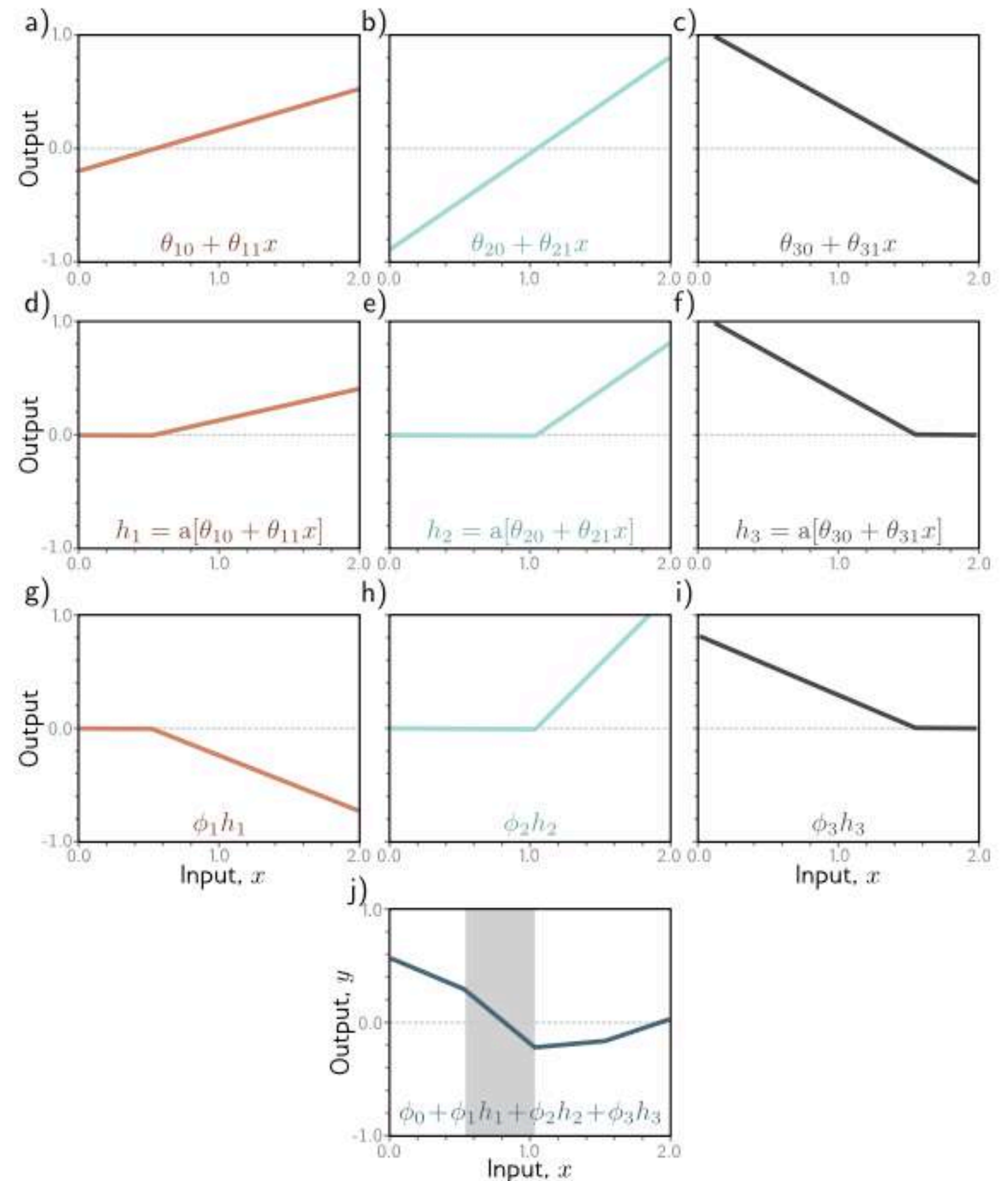
The diagram illustrates the shared learning rate η in the backpropagation update rules. The text "Learning rate" is positioned between the two equations. Two arrows originate from this text: one points diagonally upwards to the η term in the weight update equation, and the other points diagonally downwards to the η term in the bias update equation, visually demonstrating that the same learning rate is used for both parameters.

Universal approximation theorem

Any function can be approximated to arbitrary precision

Universal approximation theorem

- Any continuous function on a finite interval $[a, b]$
- Can be approximated to arbitrary precision
- By a shallow neural network $f_2 \circ \sigma \circ f_1$ where f are linear transformations and σ is a nonlinear transformation



Images

The curse of dimensionality



$$O(n^2)$$

Width x Height



28x28



100x100



200x200



400x400

Width x Height

Features



28x28

784



100x100

10.000



200x200

40.000



400x400

160.000

Width x Height

Features

Weights



28x28

784

614.656



100x100

10.000

100.000.000



200x200

40.000

1.600.000.000



400x400

160.000

25.600.000.000



Convolutions

Convolutions

	1			
1	-1	1		
	-1			
	1	-1		

0	0	0
0	1	0
0	0	0

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

0	0	0
0	1	0
0	0	0

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

0	0	0
0	1	0
0	0	0

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

0	0	0
0	1	0
0	0	0

-1		

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

0	0	0
0	1	0
0	0	0

-1	1	

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

0	0	0
0	1	0
0	0	0

-1	1	0

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

0	0	0
0	1	0
0	0	0

-1	1	0
0		

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

0	0	0
0	1	0
0	0	0

-1	1	0
0	0	

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

0	0	0
0	1	0
0	0	0

-1	1	0
0	0	0

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

0	0	0
0	1	0
0	0	0

-1	1	0
0	0	0
-1		

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

0	0	0
0	1	0
0	0	0

-1	1	0
0	0	0
-1	0	

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

0	0	0
0	1	0
0	0	0

-1	1	0
0	0	0
-1	0	0

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

0	0	0
0	1	0
0	0	0

Identity kernel

-1	1	0
0	0	0
-1	0	0

Convolutions

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

w	w	w
w	w	w
w	w	w

Learnable

\hat{y}	\hat{y}	\hat{y}
\hat{y}	\hat{y}	\hat{y}
\hat{y}	\hat{y}	\hat{y}

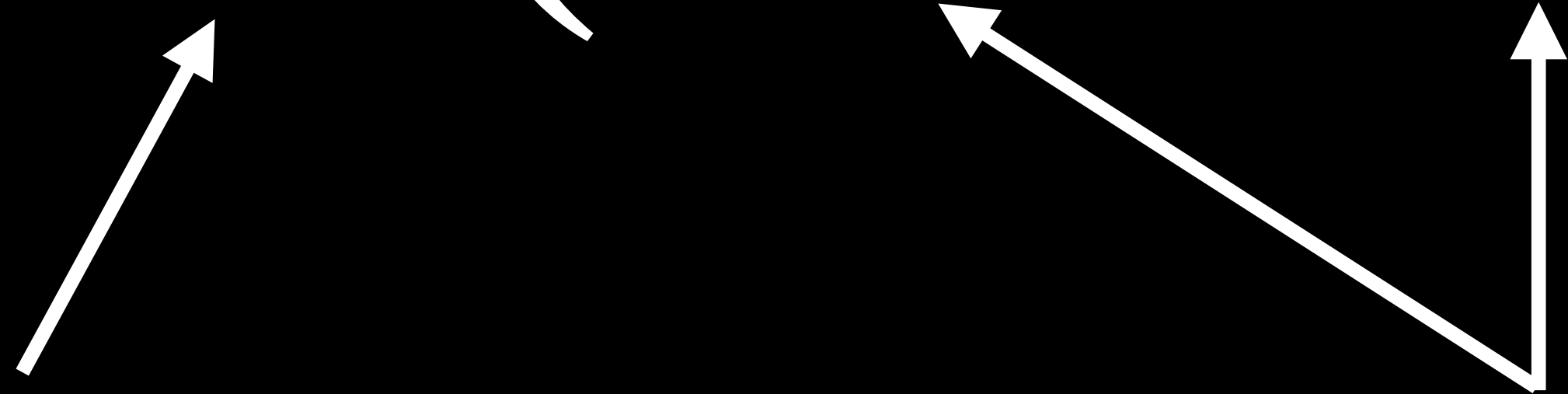
$O(k)$

Neural networks

$$\sigma(wx + b)$$

Non-linear transformation

Linear transformation



Convolutions

$$\sigma(w * x)$$

Non-linear transformation

Linear transformation

Deep neural networks

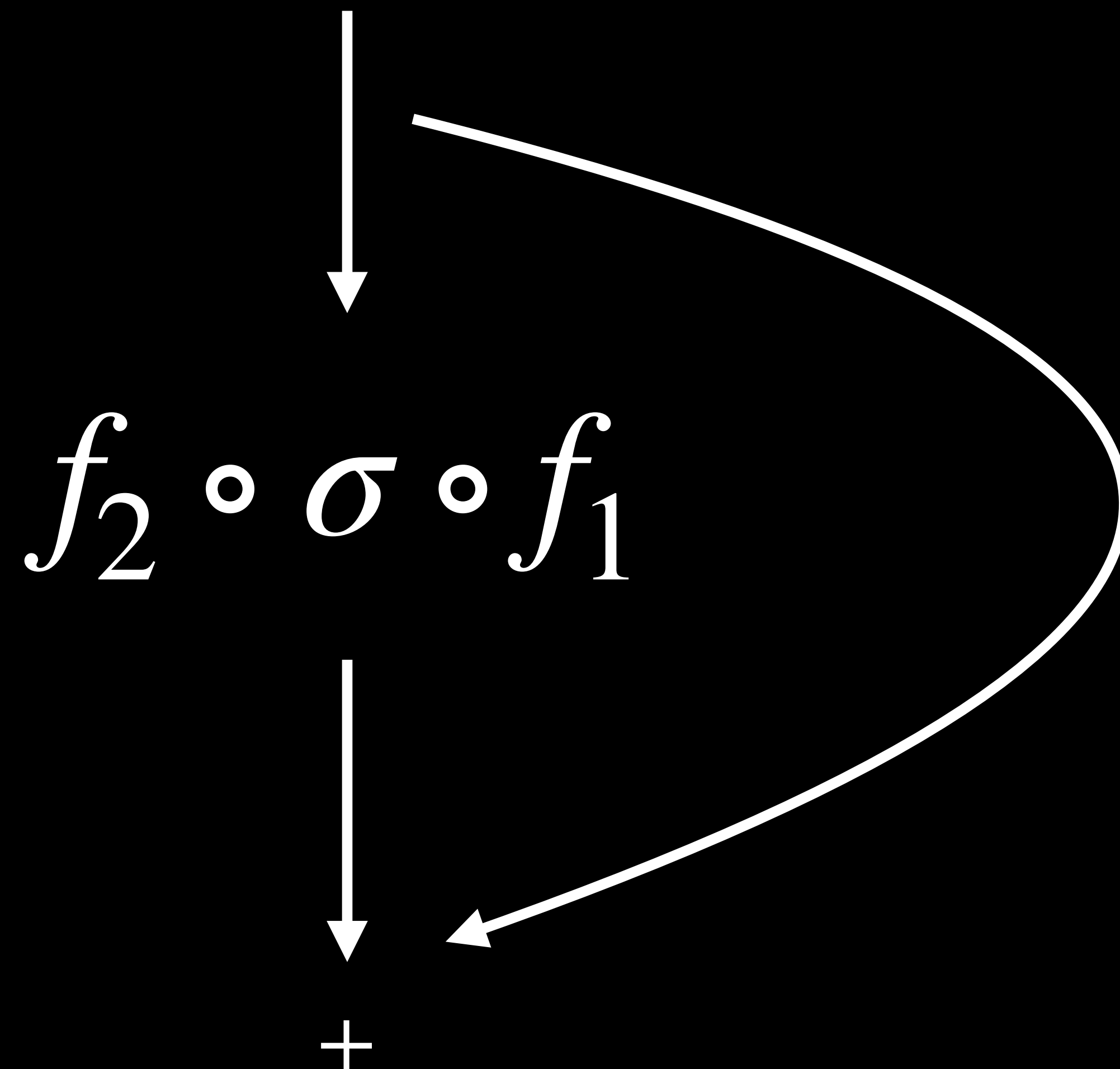
Deep neural networks



$$f_2 \circ \sigma \circ f_1$$



Deep neural networks



Timeseries

Motivation

A lot of data is sequential, varying over time:

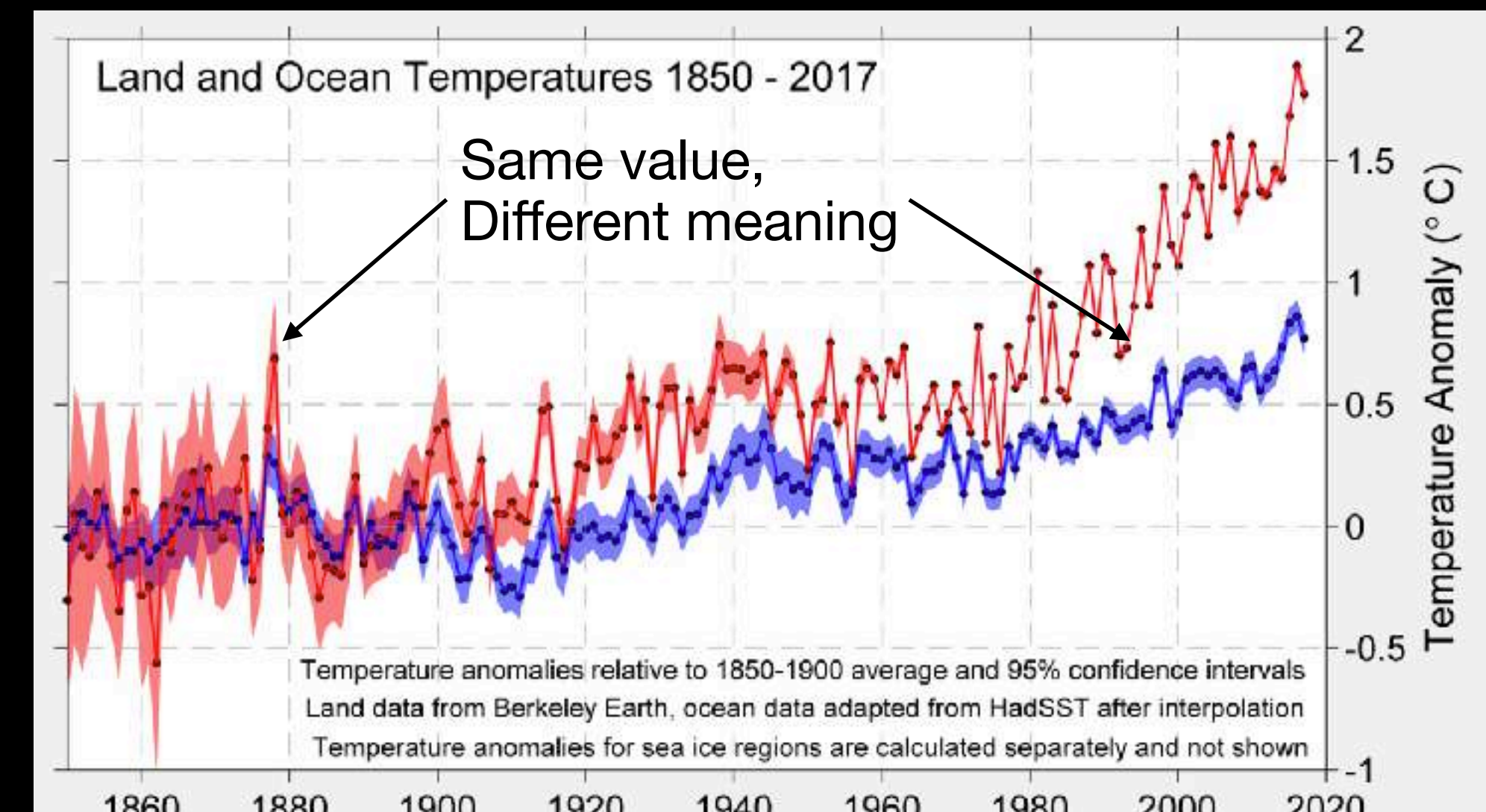
- Sentences
- Music
- EEG
- Movement
- Markets

Motivation

With sequences, the past offers context:

- Ik krijg geld van de **bank**
- Ik wil een nieuwe **bank** aanschaffen

We need the past to make sense of the future.



Data considerations

We need to worry about:

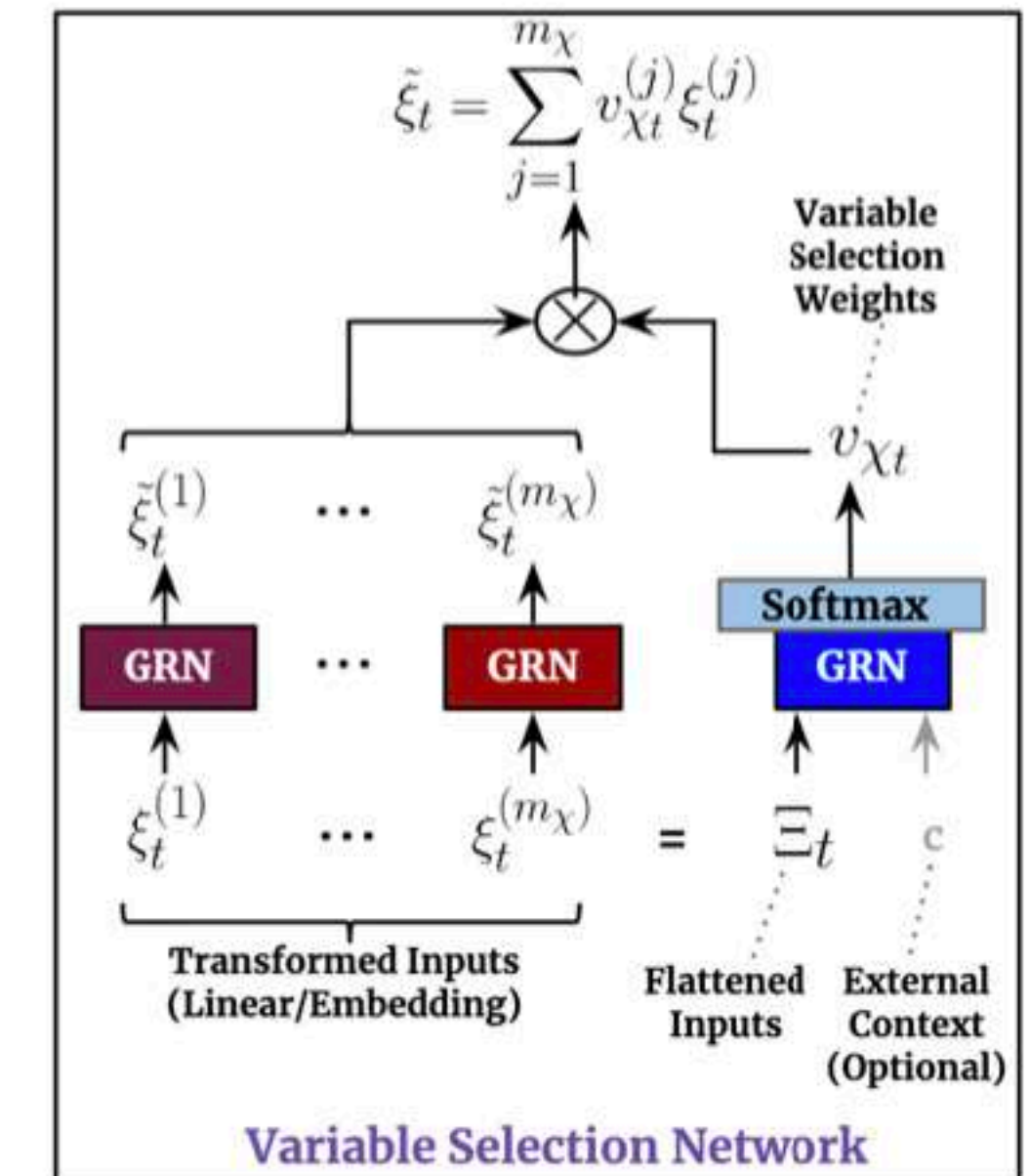
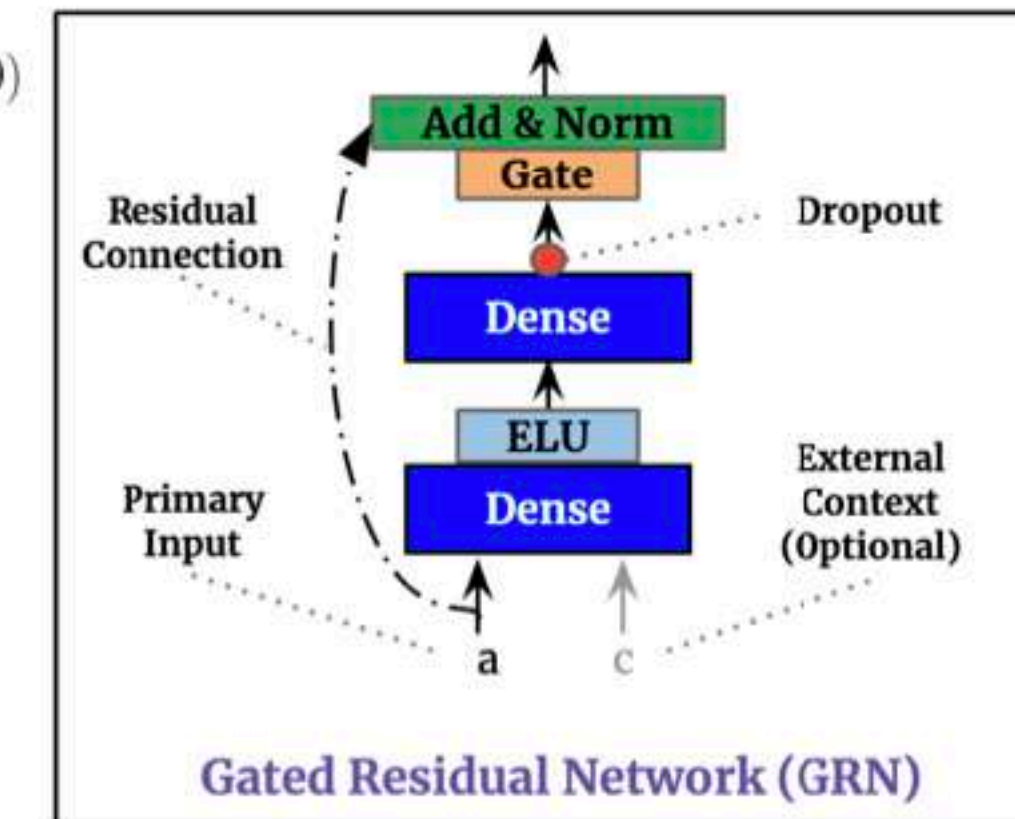
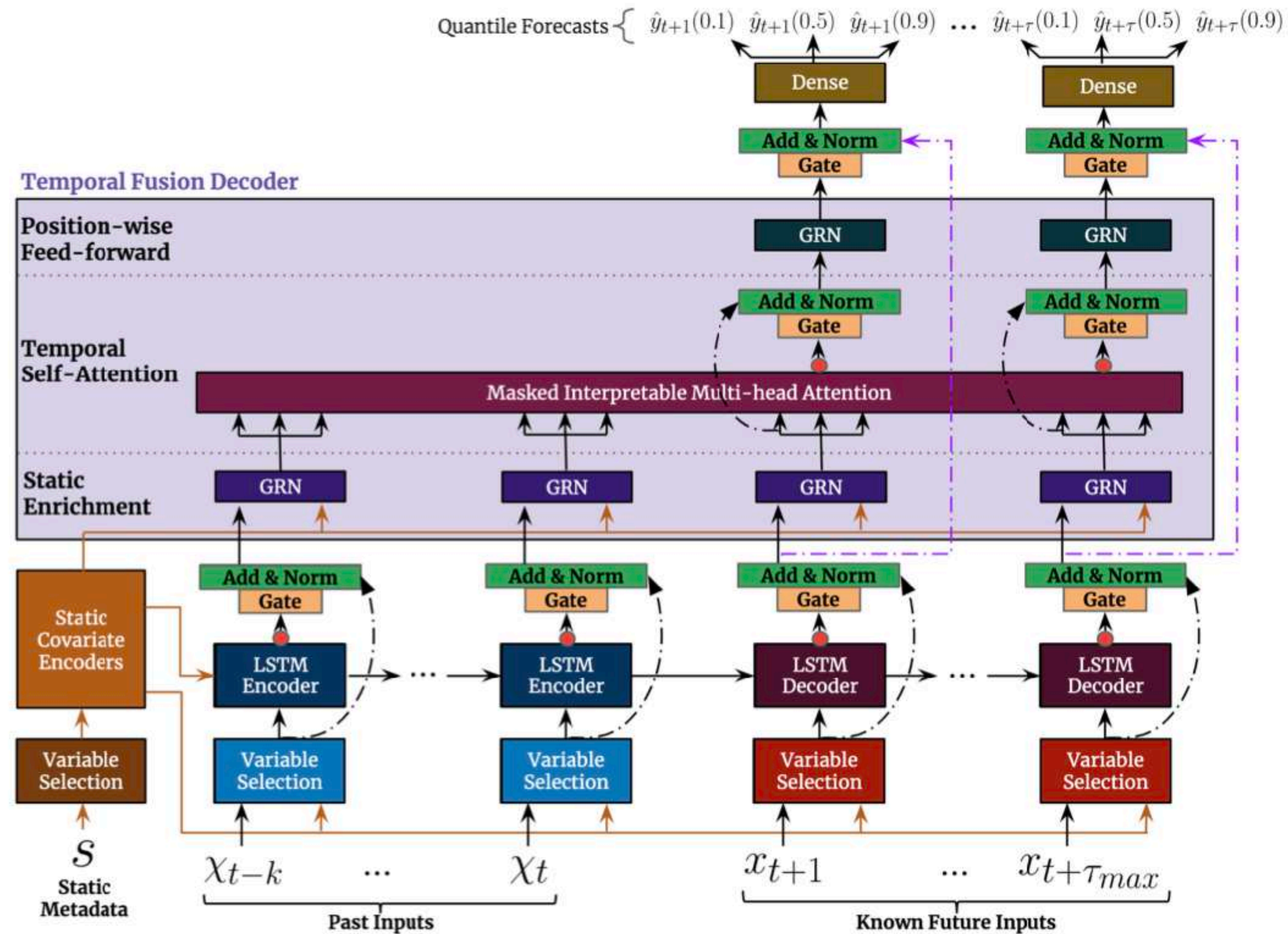
- How much of the past will we need (window)
- How much of the future do we want to predict (horizon)
- How to prepare the data without leaking data

For the last point, we need to be very careful not to “leak” the future back into the present.

History of RNNs

- 1982 RNN are discovered by John Hopfield
- 1995 The LSTM architecture was proposed with input and output gates
- 1999 Forget gates were added
- 2009 LSTM won the handwriting recognition competition
- 2013 LSTM outperformed other models at natural speech recognition
- 2014 GRU architecture was introduced
- 2017 probabilistic forecasting (DeepAR, MQRNN, TFT)

Temporal Fusion Transformer, Lim et al. (2021)



Hidden states

State

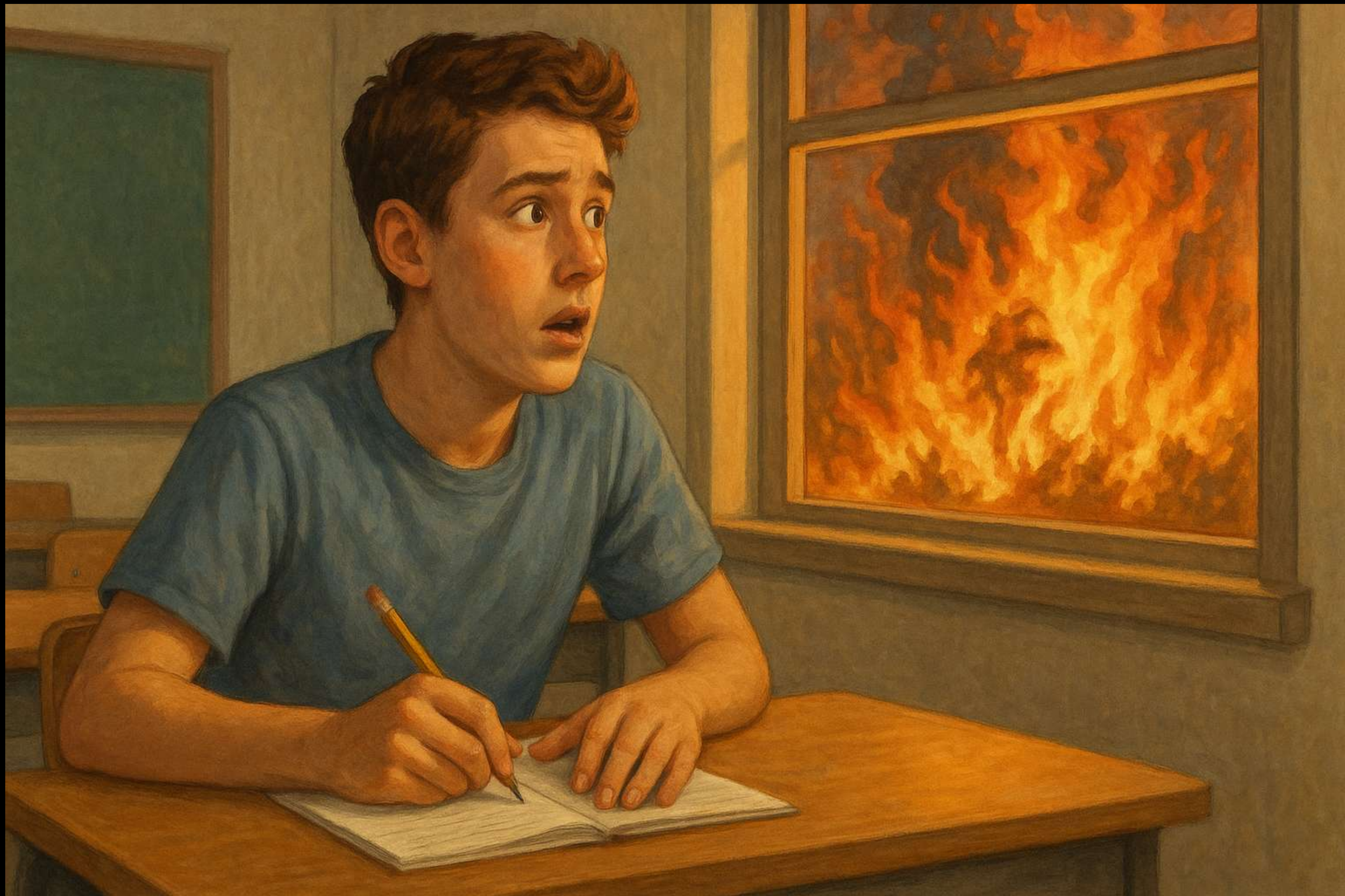
- A state gives context to new input
- Influences which elements of the input requires attention, and which elements can be ignored
- New input can change the state, such that attention is shifted to other elements of the input



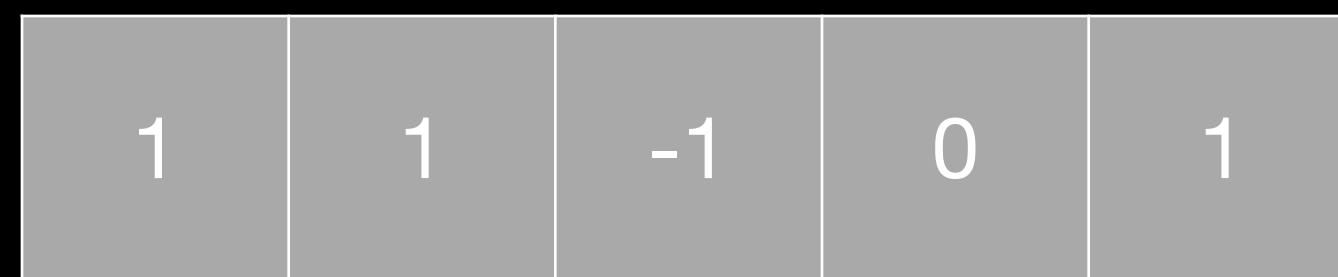








Concatenation



x_t



Input

Concatenation

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

h_{t-1}



Previous state

1	1	-1	0	1
---	---	----	---	---

x_t

Concatenation

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0

h_{t-1}

1	1	-1	0	1
---	---	----	---	---

x_t

0	1	0	0	0
1	-1	1	0	0
0	0	0	0	0
0	-1	0	0	0
0	1	-1	0	0
1	1	-1	0	1

$[x_t, h_{t-1}]$

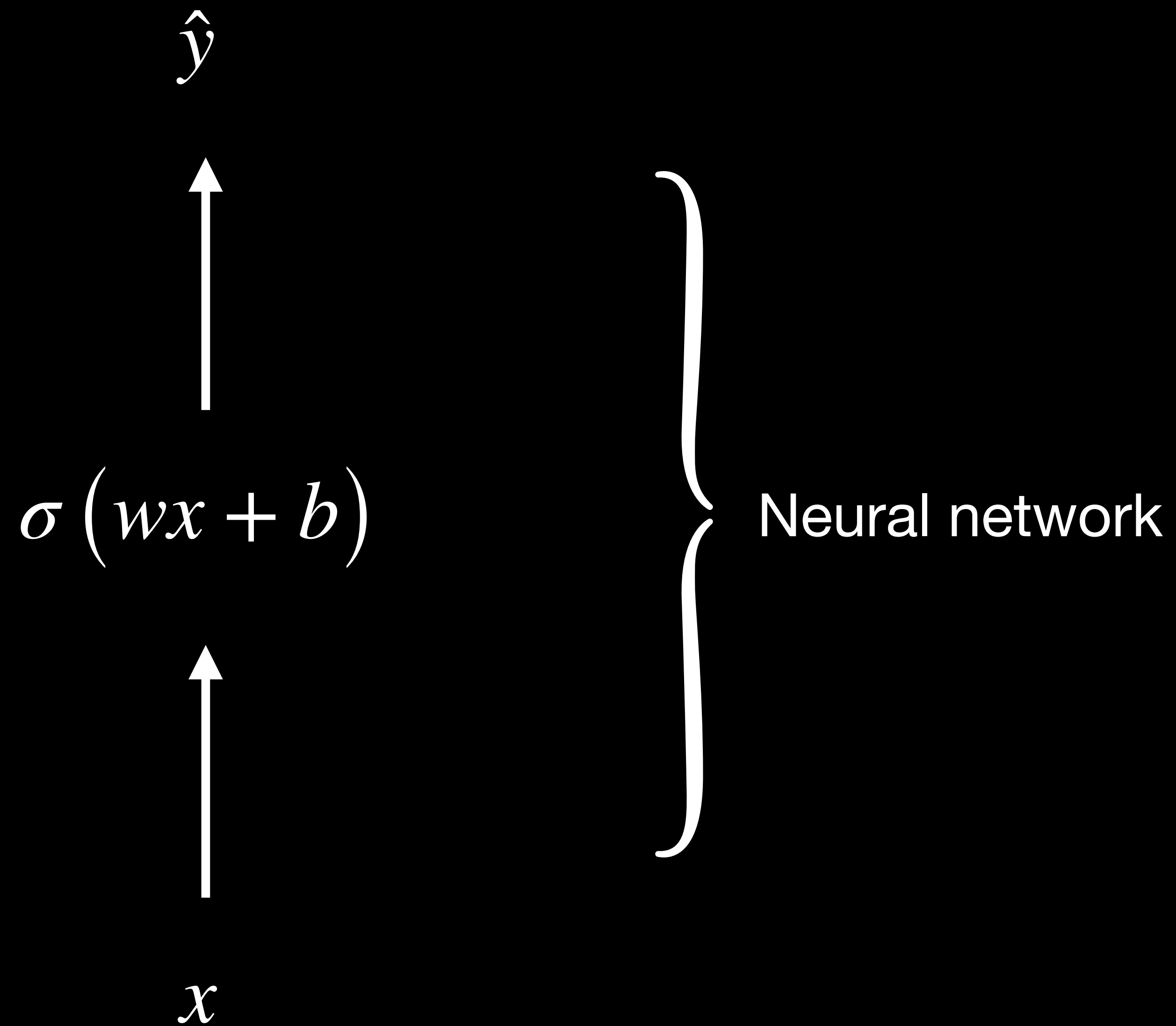
New state

$$y = \sigma(WX + b)$$

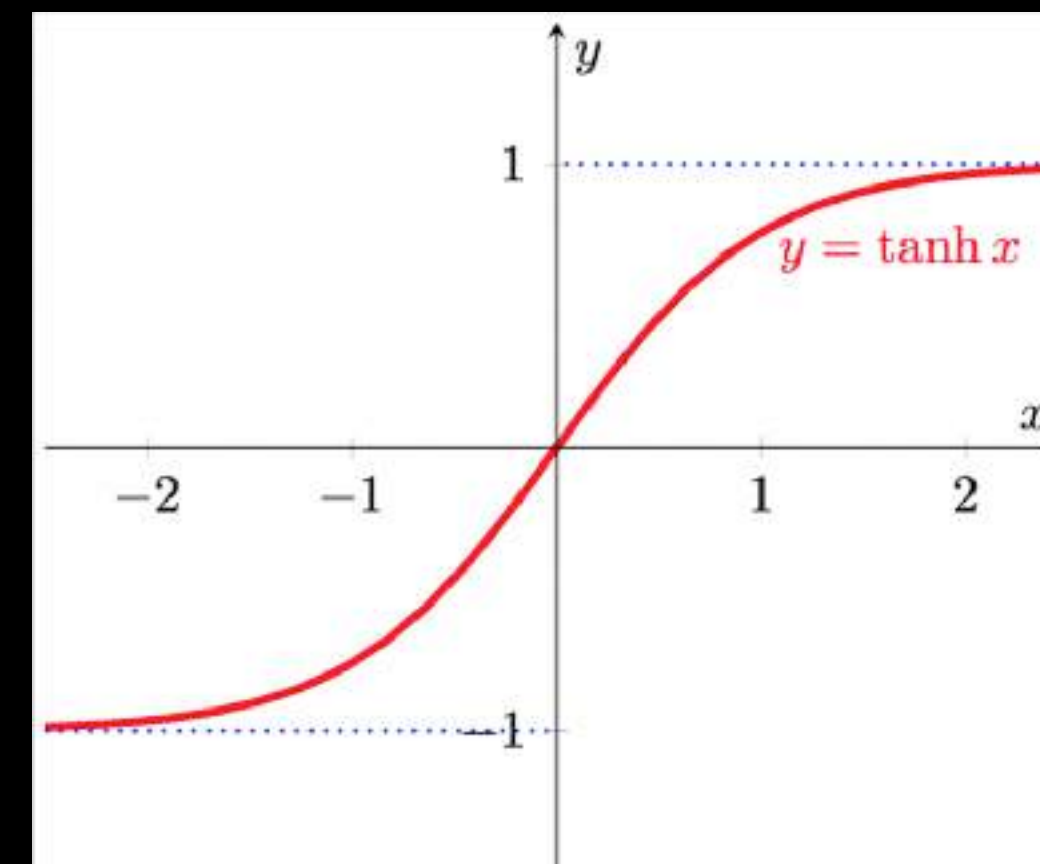
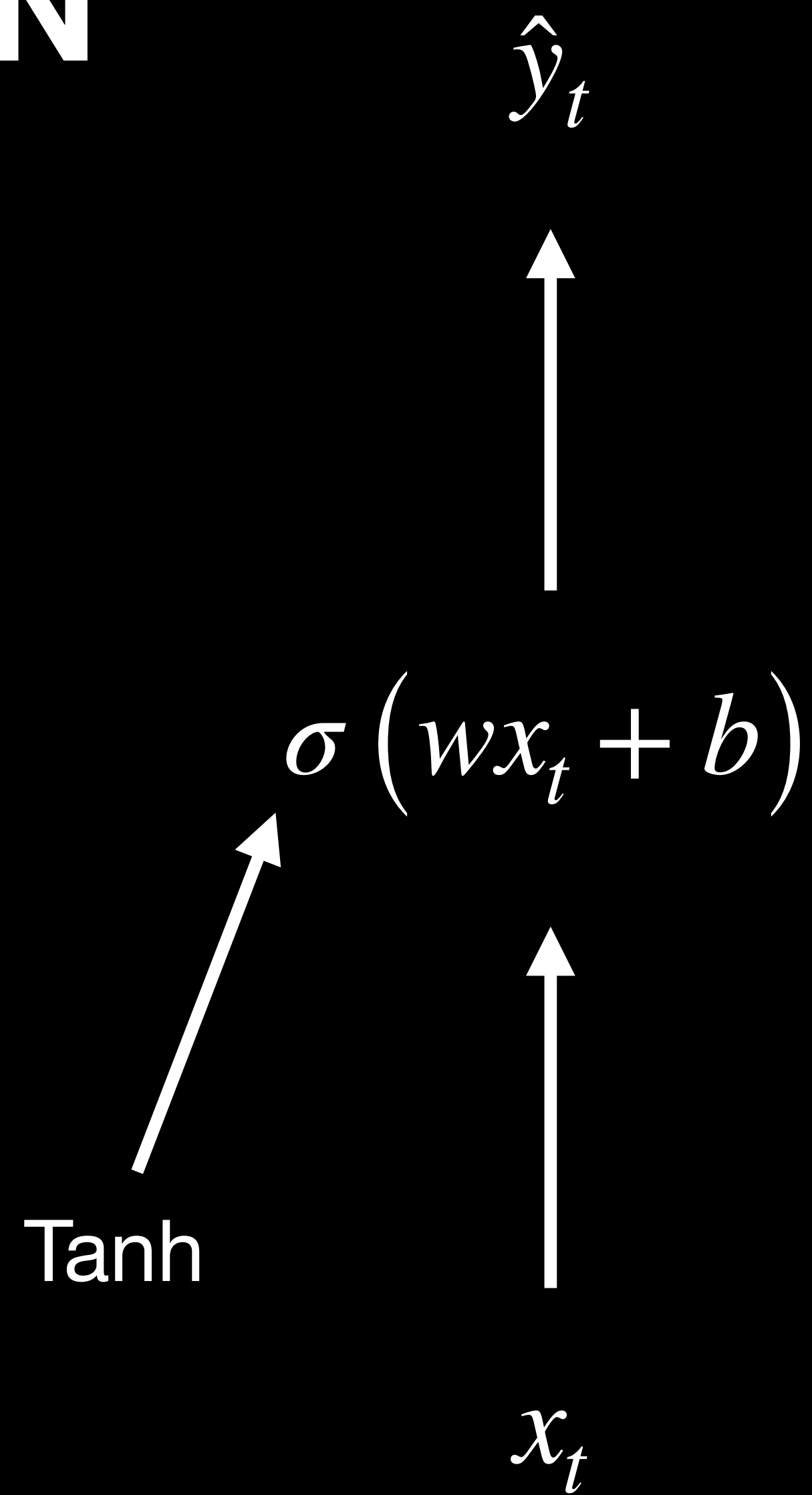
Equivalent output

$$\left\{ \begin{array}{ll} h_t = \sigma(W_x X_t + W_h h_{t-1} + b) & \leftarrow \text{Slower} \\ h_t = \sigma(W [X_t, h_{t-1}] + b) & \leftarrow \text{Faster} \end{array} \right.$$

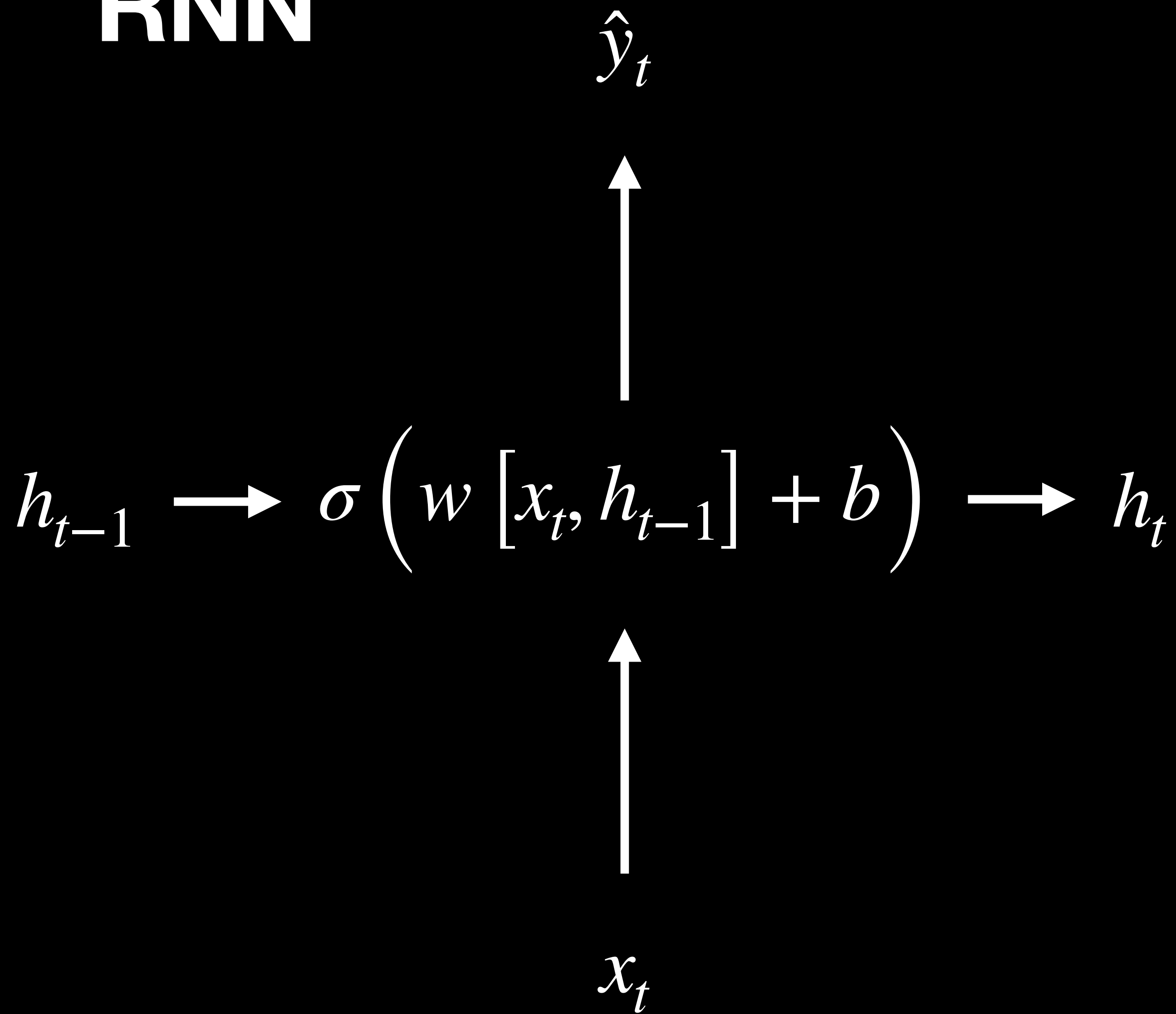
RNN



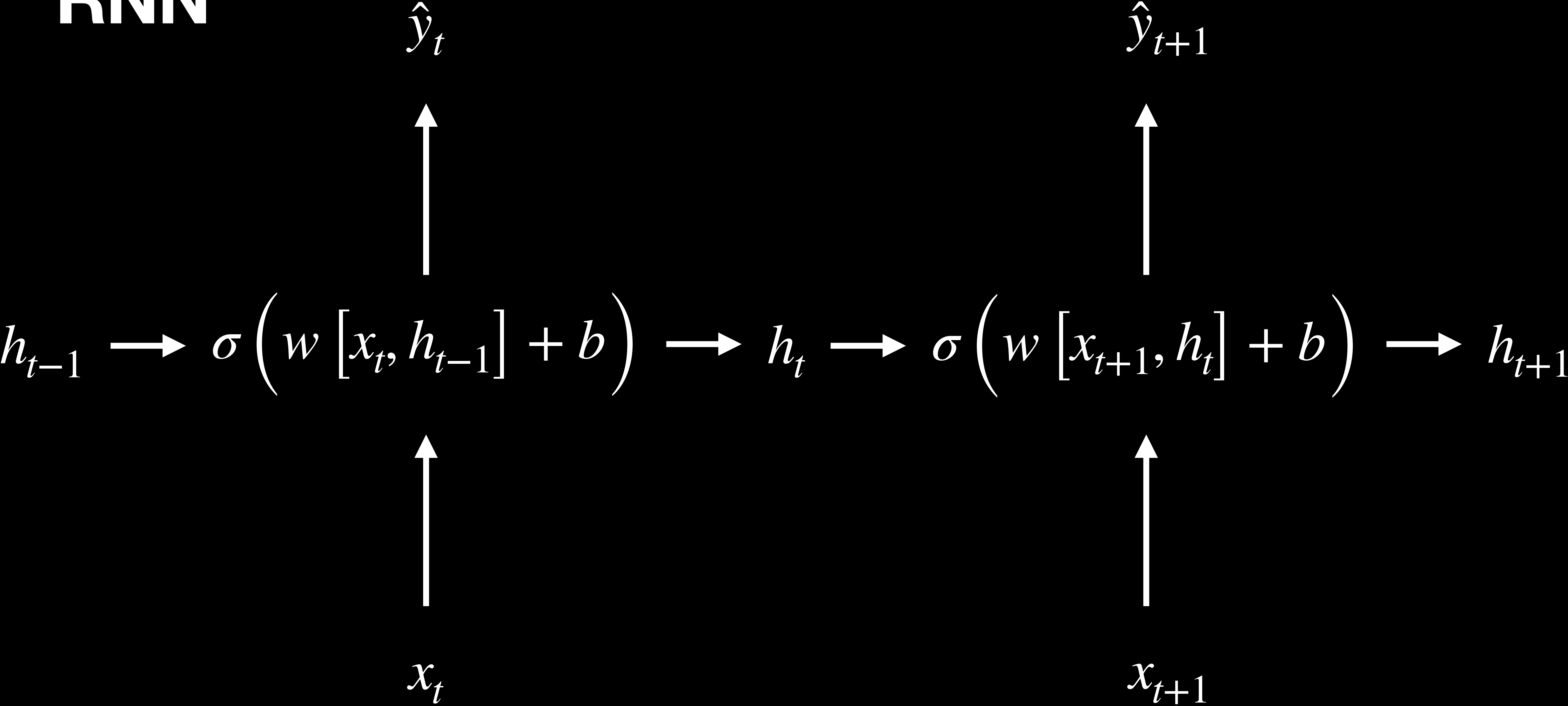
RNN

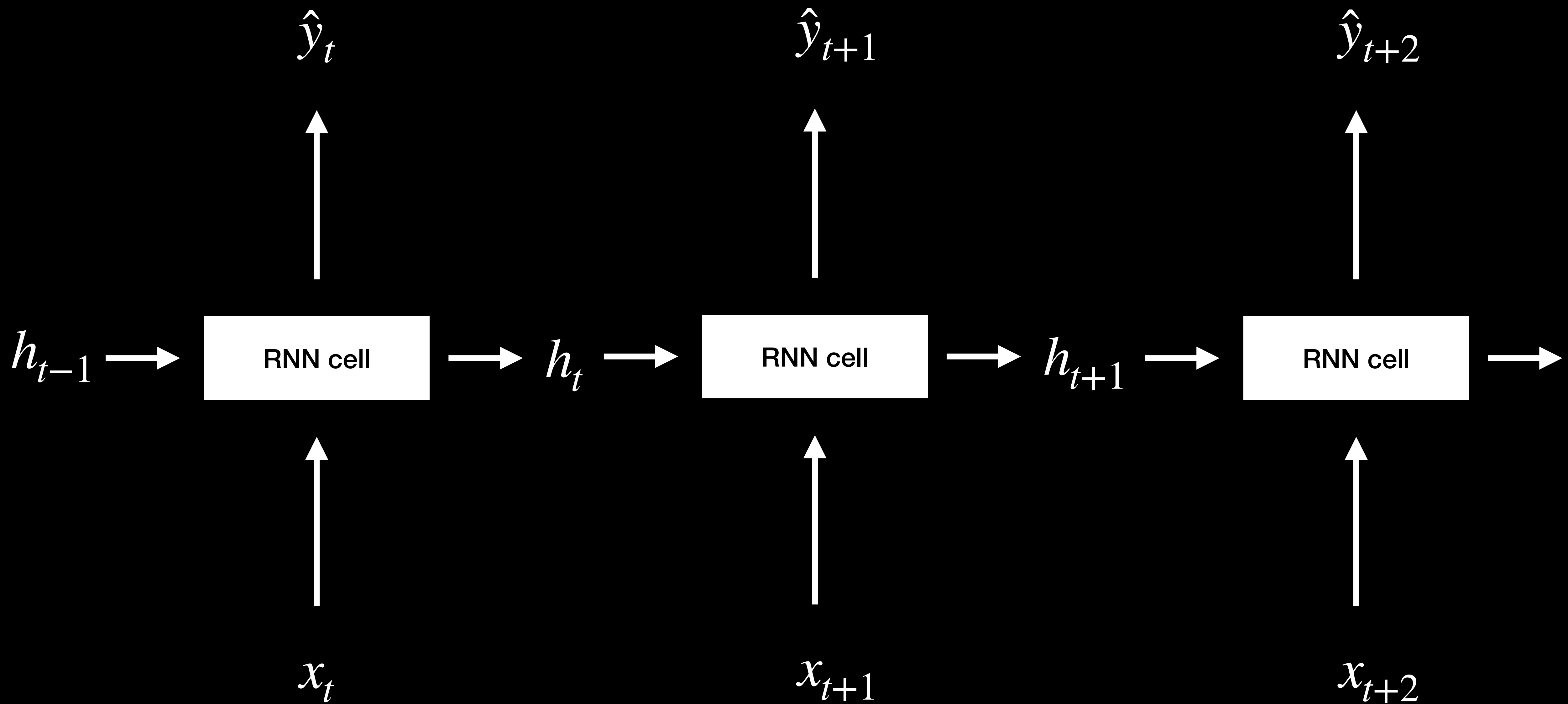


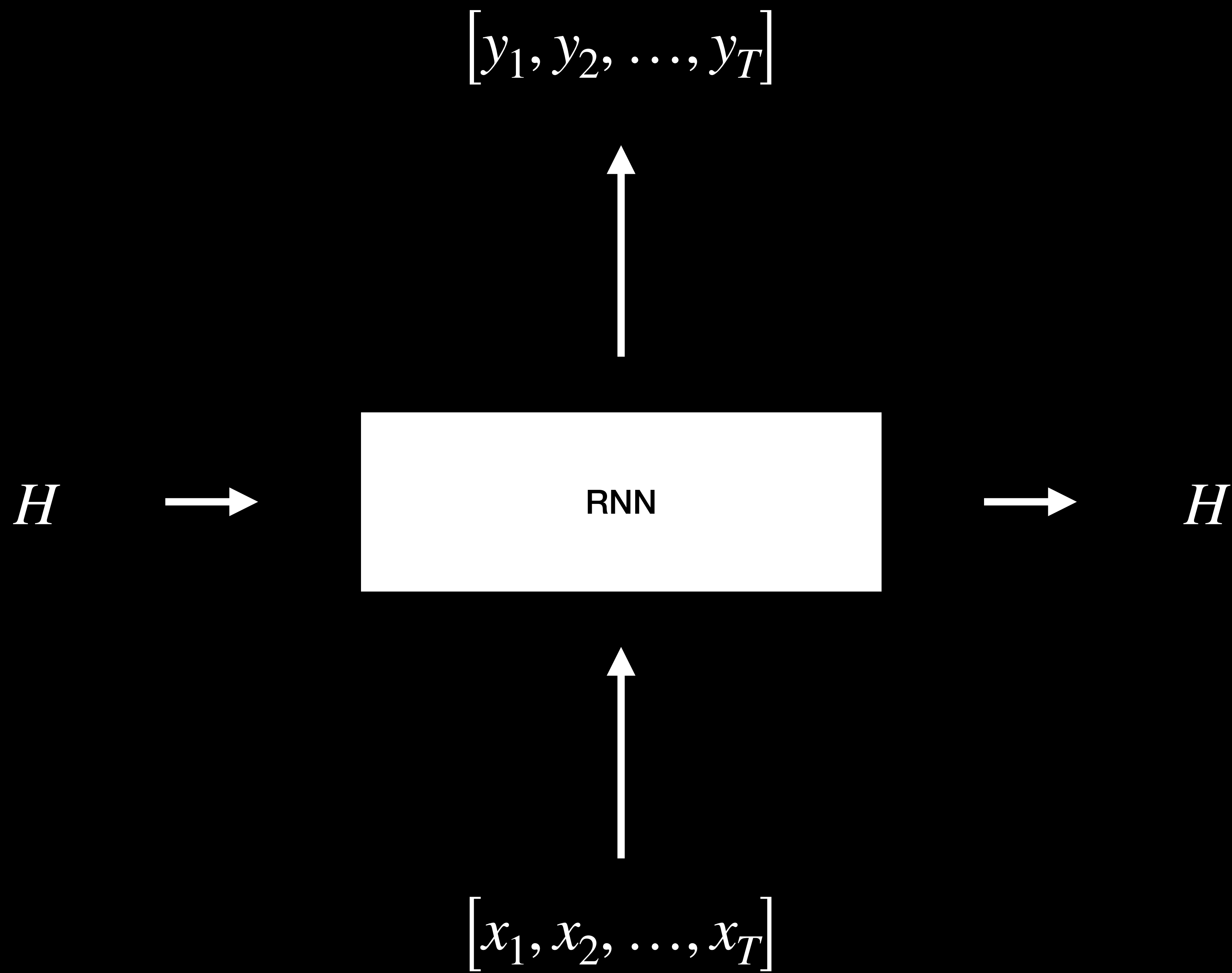
RNN



RNN







The art of forgetting

RNNs have not explicit way to forget or retain memory.

We can make this a bit more advanced by adding gates.

A gate Γ controls

- what part of the past we retain
- what part we forget.

GRU - Remember & forget

Gated Residual Unit

We need to be able to:

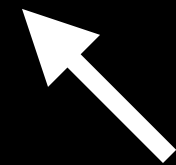
- *Remember* the past, and completely ignore the new state
- *Forget* the past, and focus on the present
- *Something in between* where we find a ratio between forgetting and remembering.

We also want to gate to be influenced by both the new input and the old state.

Gates

0	1	0	0
1	-1	1	0
0	0	0	0
0	-1	0	0

X



Input

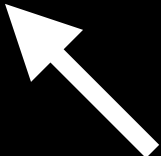
Gates

0	1	0	0
1	-1	1	0
0	0	0	0
0	-1	0	0

X

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

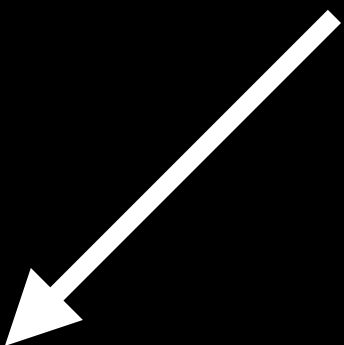
Γ



Hadamard product

Gates

Same shape as X



0	1	0	0
1	-1	1	0
0	0	0	0
0	-1	0	0

X



0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

Γ

Gates

0	1	0	0
1	-1	1	0
0	0	0	0
0	-1	0	0

X



0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

Γ



Gate

Gates

0	1	0	0
1	-1	1	0
0	0	0	0
0	-1	0	0

X



0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

Γ

$=$

Y

Gates

0	1	0	0
1	-1	1	0
0	0	0	0
0	-1	0	0

X

\otimes

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

Γ

$=$

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

Y

Gates

0	1	0	0
1	-1	1	0
0	0	0	0
0	-1	0	0

X



1	1	1	1
1	1	1	1
1	1	1	1
1	1	1	1

Γ

$=$

Y

Gates

0	1	0	0
1	-1	1	0
0	0	0	0
0	-1	0	0

X



1	1	1	1
1	1	1	1
1	1	1	1
1	1	1	1

Γ

$=$

0	1	0	0
1	-1	1	0
0	0	0	0
0	-1	0	0

Y

Gates

0	1	0	0
1	-1	1	0
0	0	0	0
0	-1	0	0

X



0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5

Γ

$=$

Y

Gates

0	1	0	0
1	-1	1	0
0	0	0	0
0	-1	0	0

X



0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5

Γ

$=$

0	0.5	0	0
0.5	-0.5	0.5	0
0	0	0	0
0	-0.5	0	0

Y

Gates

0	1	0	0
1	-1	1	0
0	0	0	0
0	-1	0	0

X

\otimes

w	w	w	w
w	w	w	w
w	w	w	w
w	w	w	w

Γ

$=$

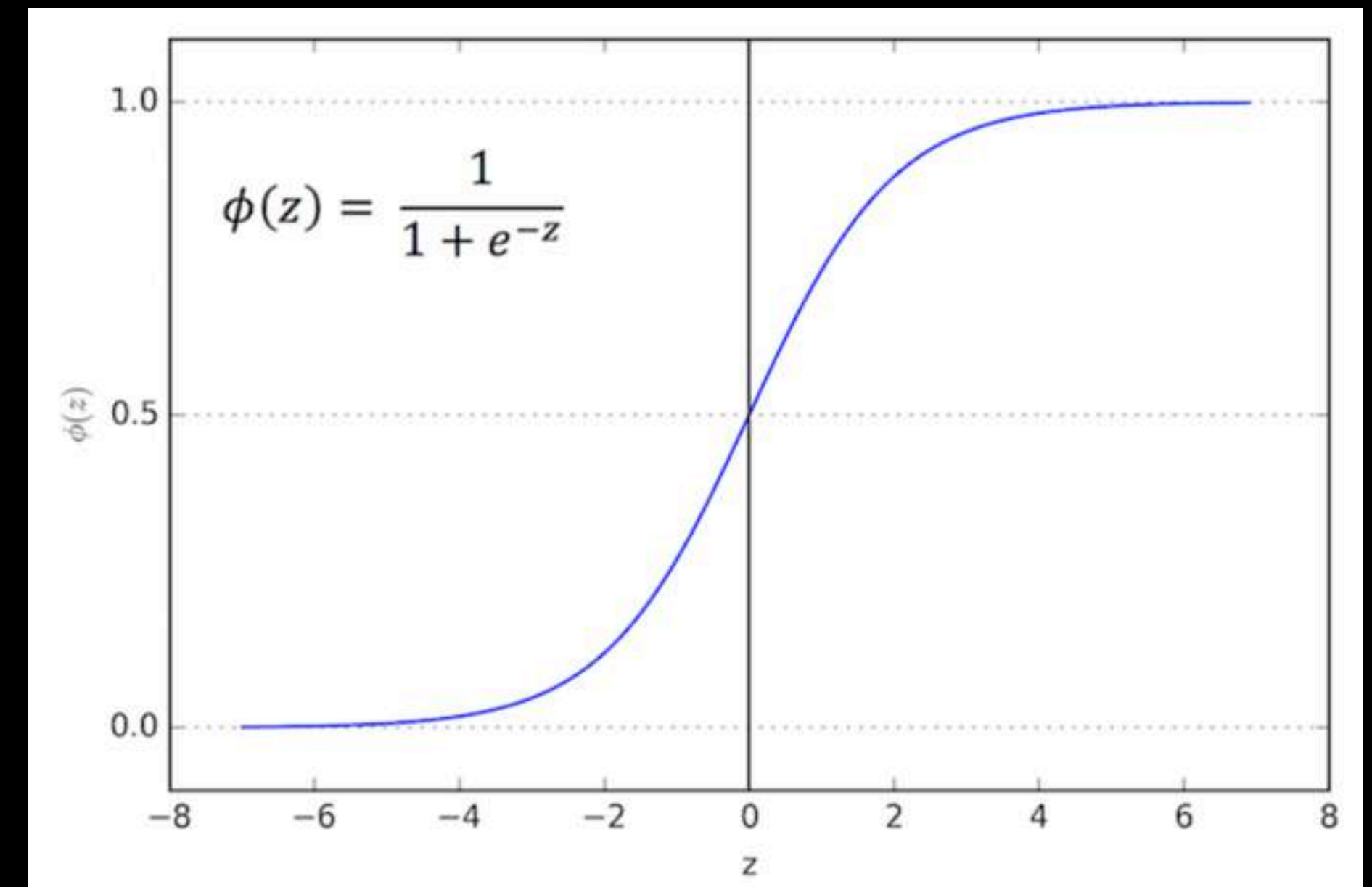
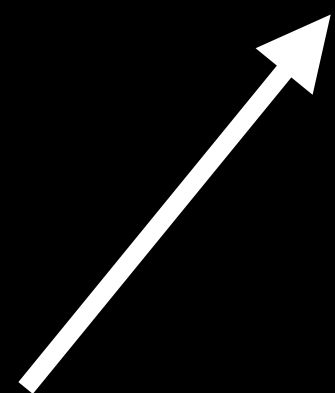
y	y	y	y
y	y	y	y
y	y	y	y
y	y	y	y

Y

Gates

$$\Gamma = \sigma(W[X_t, h_{t-1}] + b)$$

Sigmoid



GRU

Gate \longrightarrow

$$\Gamma = \sigma(W_{\Gamma}[X_t, h_{t-1}] + b_{\Gamma})$$
$$\tilde{h}_t = \tanh(W_H[X_t, h_{t-1}] + b_H)$$
$$h_t = \Gamma \otimes h_{t-1} + (1 - \Gamma) \otimes \tilde{h}_t$$

GRU

Candidate State \longrightarrow

$$\Gamma = \sigma(W_{\Gamma}[X_t, h_{t-1}] + b_{\Gamma})$$
$$\tilde{h}_t = \tanh(W_H[X_t, h_{t-1}] + b_H)$$
$$h_t = \Gamma \otimes h_{t-1} + (1 - \Gamma) \otimes \tilde{h}_t$$

GRU

$$\Gamma = \sigma(W_{\Gamma}[X_t, h_{t-1}] + b_{\Gamma})$$

$$\tilde{h}_t = \tanh(W_H[X_t, h_{t-1}] + b_H)$$

Next State \longrightarrow
$$h_t = \Gamma \otimes h_{t-1} + (1 - \Gamma) \otimes \tilde{h}_t$$

GRU

Same input



$$\Gamma = \sigma(W_{\Gamma}[\mathbf{X}_t, \mathbf{h}_{t-1}] + b_{\Gamma})$$

$$\tilde{\mathbf{h}}_t = \tanh(W_H[\mathbf{X}_t, \mathbf{h}_{t-1}] + b_H)$$

$$\mathbf{h}_t = \Gamma \otimes \mathbf{h}_{t-1} + (1 - \Gamma) \otimes \tilde{\mathbf{h}}_t$$

GRU

Different weights

$$\left\{ \begin{array}{l} \Gamma = \sigma(\mathbf{W}_{\Gamma}[X_t, h_{t-1}] + \mathbf{b}_{\Gamma}) \\ \tilde{h}_t = \tanh(\mathbf{W}_H[X_t, h_{t-1}] + \mathbf{b}_H) \\ h_t = \Gamma \otimes h_{t-1} + (1 - \Gamma) \otimes \tilde{h}_t \end{array} \right.$$

GRU

$$\Gamma = \sigma(W_{\Gamma}[X_t, h_{t-1}] + b_{\Gamma})$$

$$\tilde{h}_t = \tanh(W_H[X_t, h_{t-1}] + b_H)$$

$$h_t = \Gamma \otimes h_{t-1} + (1 - \Gamma) \otimes \tilde{h}_t$$

The gate Γ controls, based on input and context, how much remembered.



GRU - full

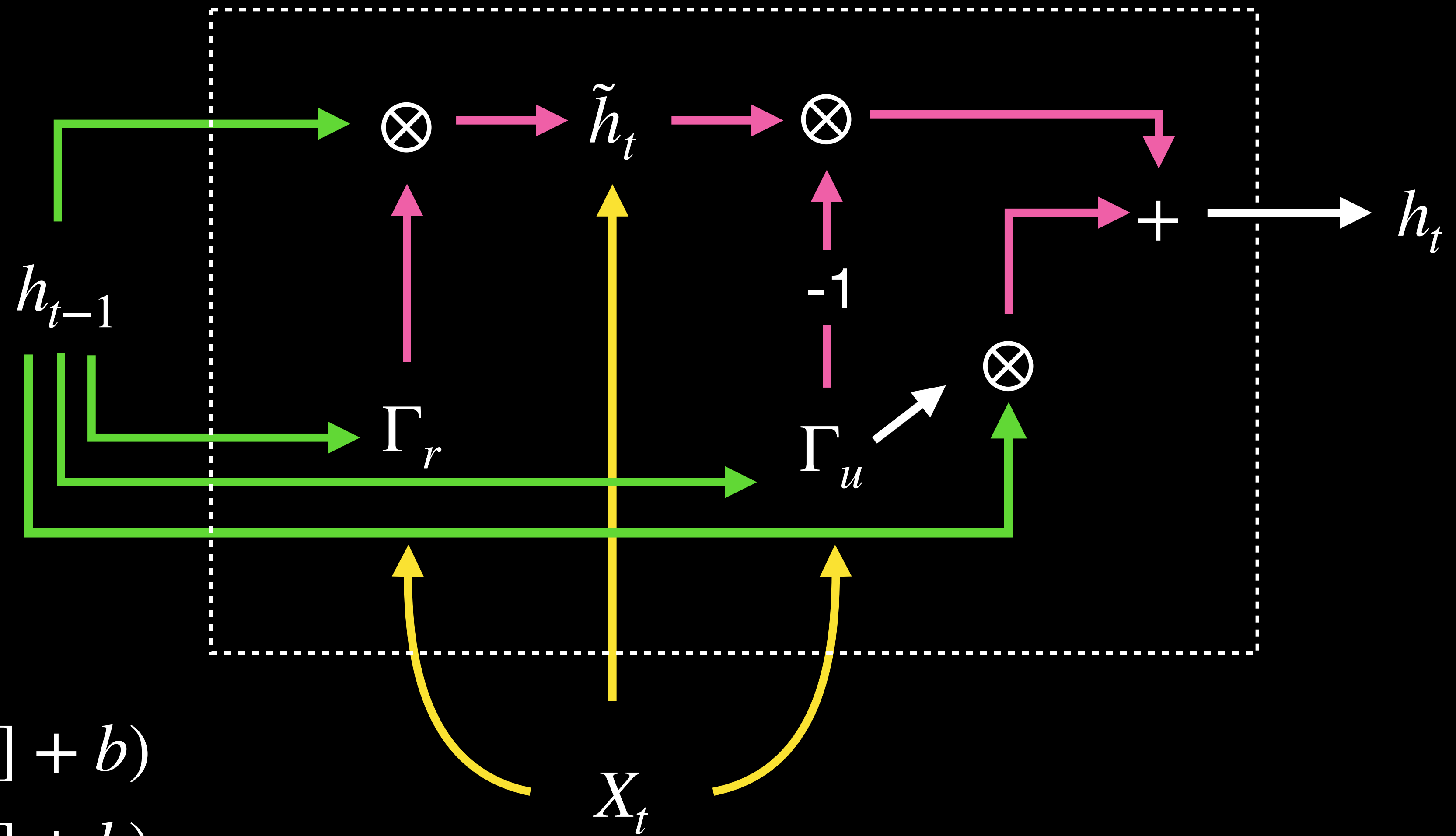
The full GRU has two gates, but the principle is the same

$$\Gamma_u = \sigma(W[X_t, h_{t-1}] + b)$$

$$\Gamma_r = \sigma(W[X_t, h_{t-1}] + b)$$

$$\tilde{h}_t = \tanh(W[X_t, \Gamma_r \otimes h_{t-1}] + b)$$

$$h_t = \Gamma_u \otimes h_{t-1} + (1 - \Gamma_u) \otimes \tilde{h}_t$$

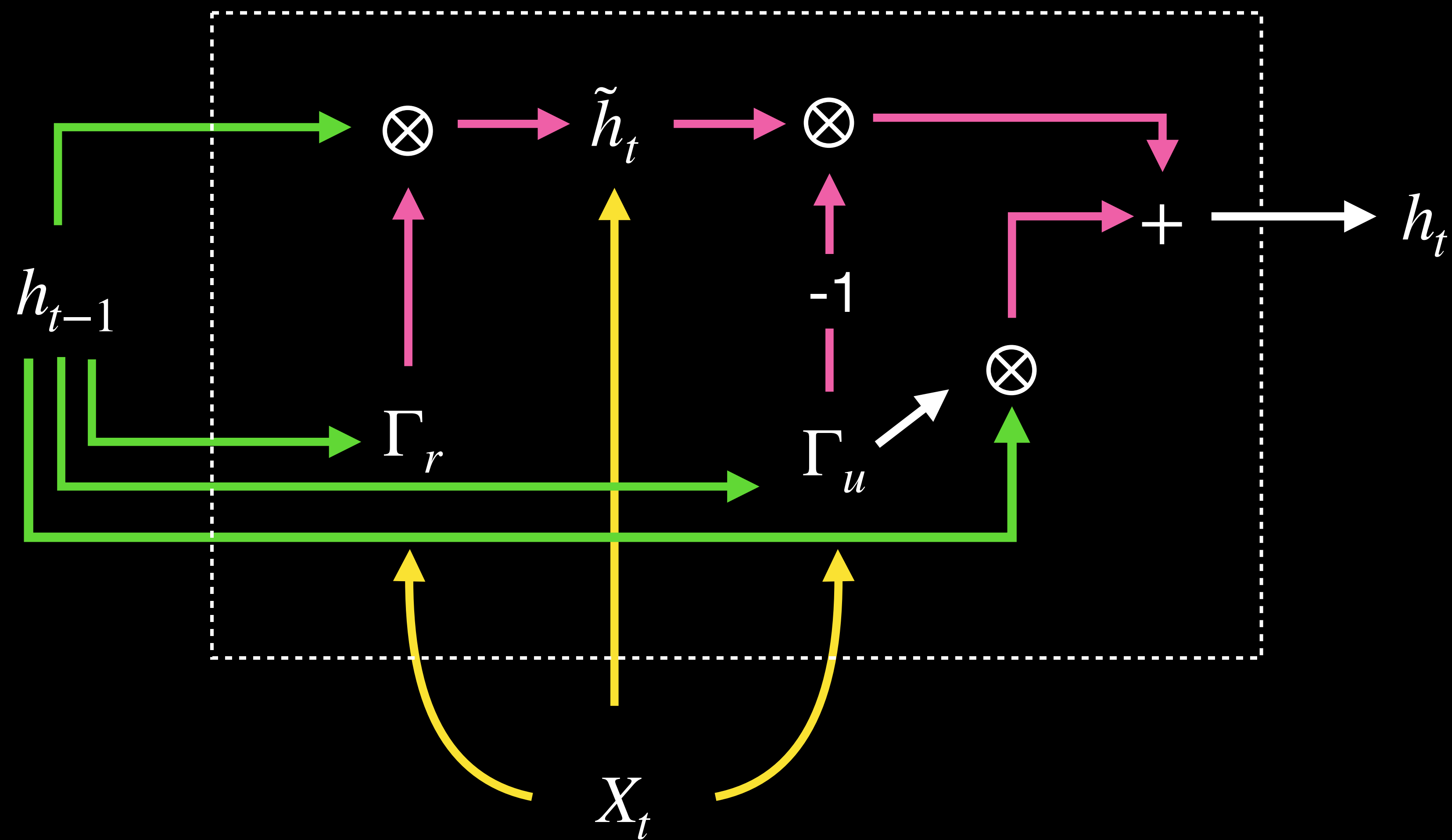


$$\Gamma_u = \sigma(W[X_t, h_{t-1}] + b)$$

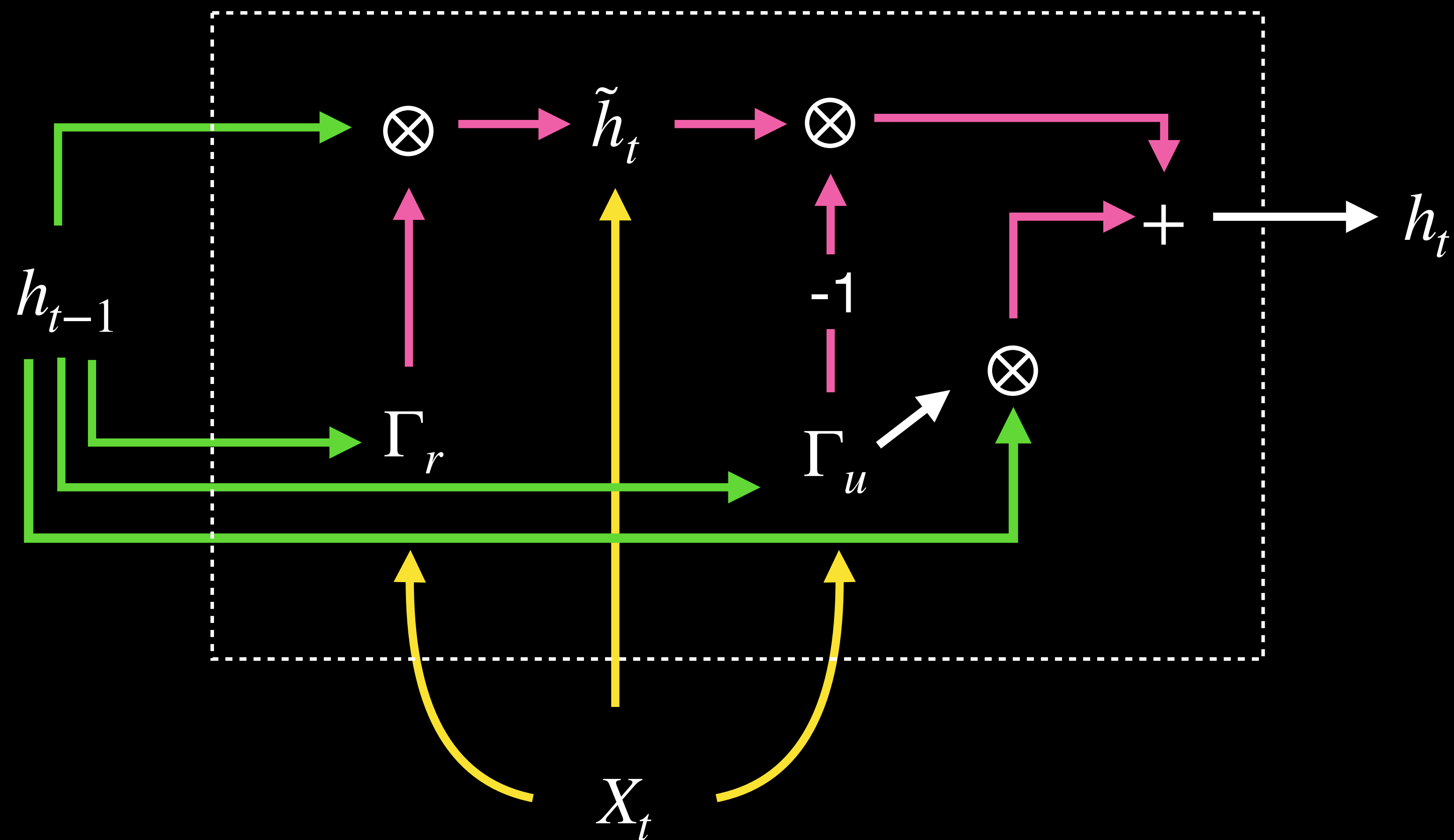
$$\Gamma_r = \sigma(W[X_t, h_{t-1}] + b)$$

$$\tilde{h}_t = \tanh(W[X_t, \Gamma_r \otimes h_{t-1}] + b)$$

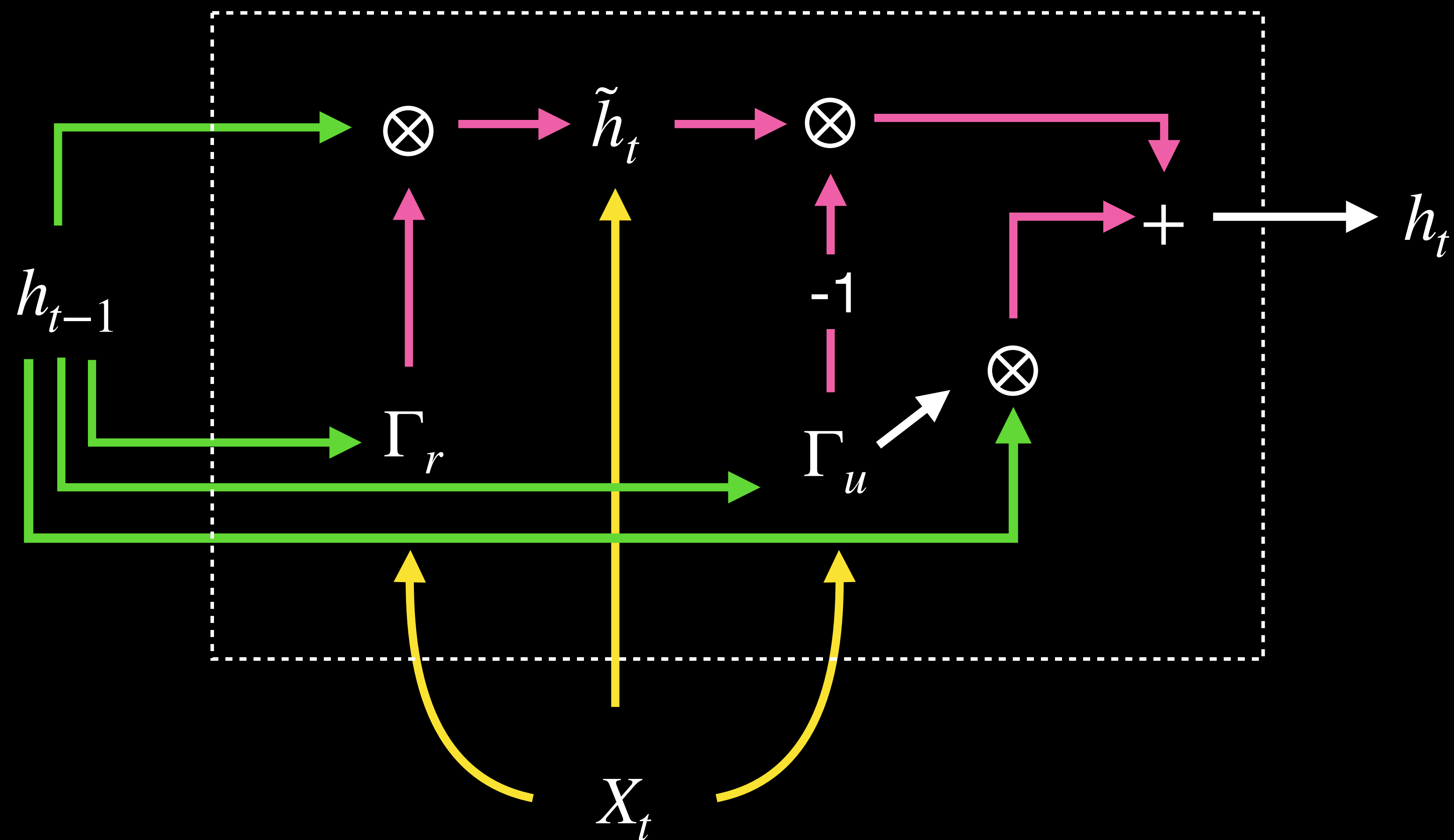
$$h_t = \Gamma_u \otimes h_{t-1} + (1 - \Gamma_u) \otimes \tilde{h}_t$$



We use the hidden state h_{t-1} and X_t to create two gates.



The reset gate Γ_r controls how much of the past h_{t-1} is mixed into X_t to create a new candidate context \tilde{h}



The other gate is the update gate Γ_u and this balances the old h_{t-1} and the new \tilde{h}_t

GRU

Compare the [Trax implementation](#) with the formulas

$$\Gamma_u = \sigma(W[X_t, h_{t-1}] + b)$$

$$\Gamma_r = \sigma(W[X_t, h_{t-1}] + b)$$

$$\tilde{h}_t = \tanh(W[X_t, \Gamma_r \otimes h_{t-1}] + b)$$

$$h_t = \Gamma_u \otimes h_{t-1} + (1 - \Gamma_u) \otimes \tilde{h}_t$$

```
def forward(self, inputs):
    x, gru_state = inputs

    # Dense layer on the concatenation of x and h.
    w1, b1, w2, b2 = self.weights
    y = jnp.dot(jnp.concatenate([x, gru_state], axis=-1), w1) + b1

    # Update and reset gates.
    u, r = jnp.split(fastmath.sigmoid(y), 2, axis=-1)

    # Candidate.
    c = jnp.dot(jnp.concatenate([x, r * gru_state], axis=-1), w2) + b2

    new_gru_state = u * gru_state + (1 - u) * jnp.tanh(c)
    return new_gru_state, new_gru_state
```

LSTM

The LSTM has

- three gates (update, input and forget) instead of two (update and reset)
- Has both a context C and a hidden state h

$$\Gamma_u = \sigma(W[X_t, h_{t-1}] + b)$$

$$\Gamma_i = \sigma(W[X_t, h_{t-1}] + b)$$

$$\Gamma_f = \sigma(W[X_t, h_{t-1}] + b)$$

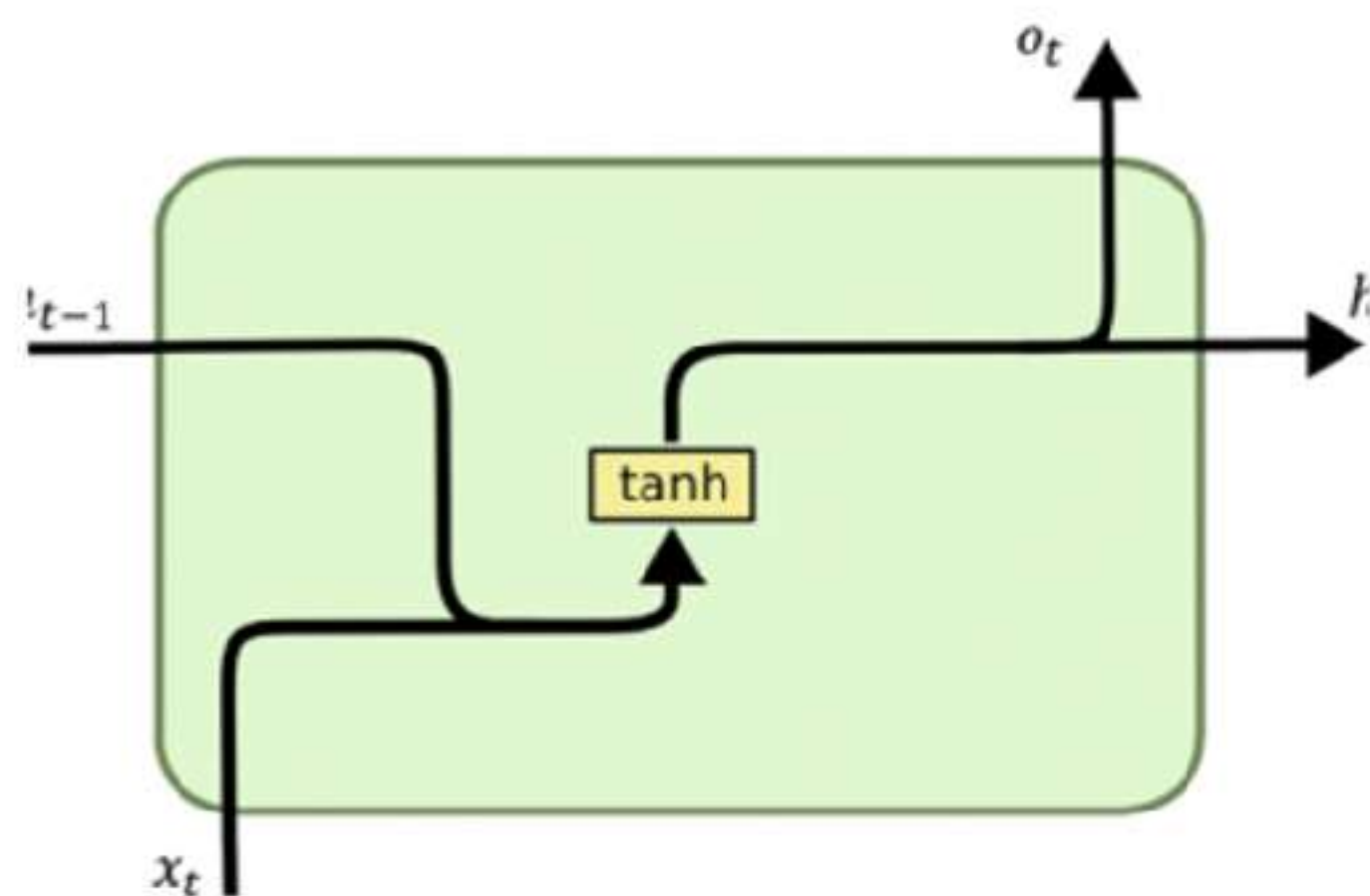
$$\tilde{h} = \Gamma_i \otimes \tanh(W[X_t, h_{t-1}] + b)$$

$$\tilde{C} = \tanh(\Gamma_f \otimes C + \tilde{h})$$

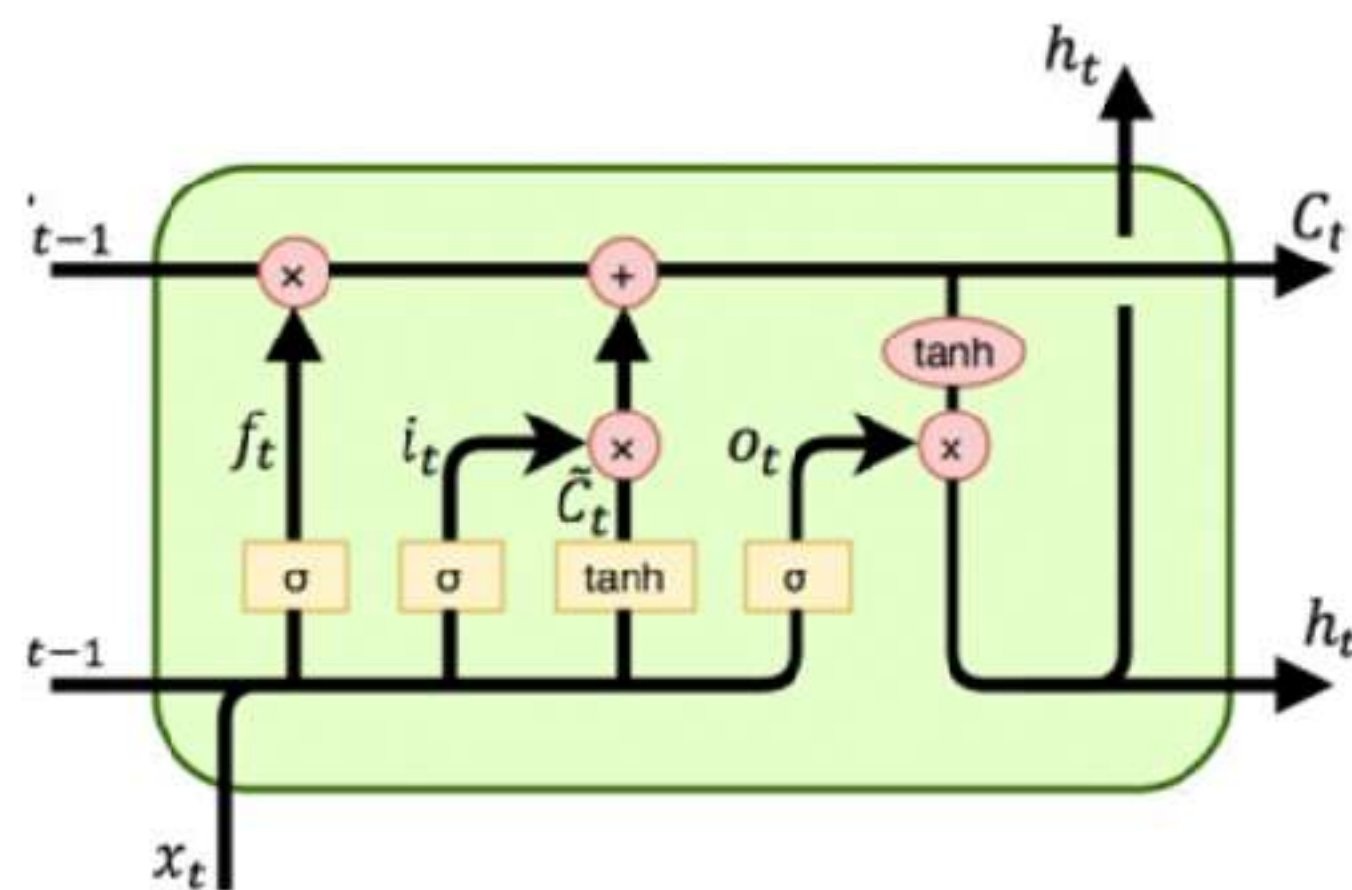
$$h_t = \Gamma_u \otimes \tilde{C}$$

Overview

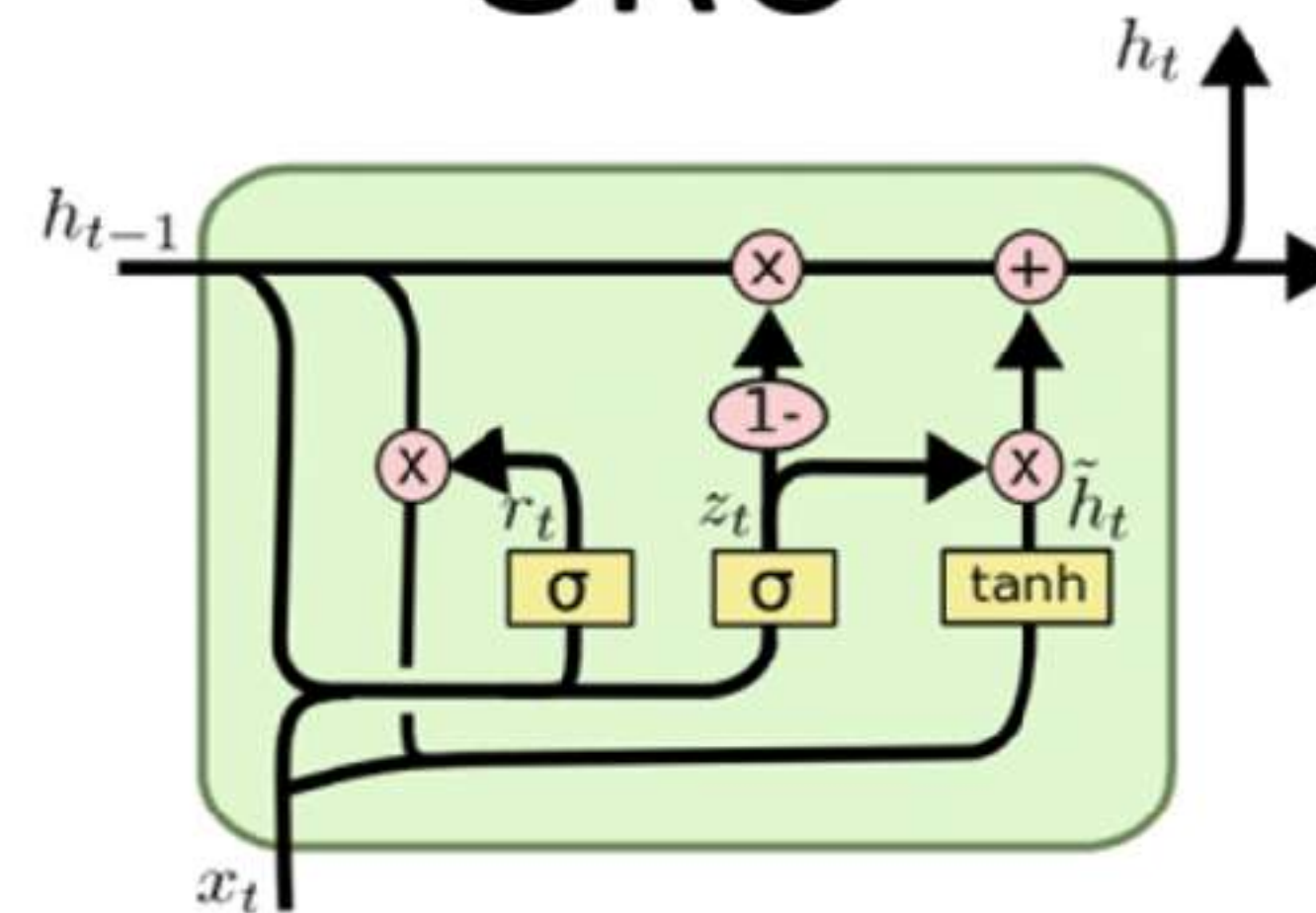
RNN



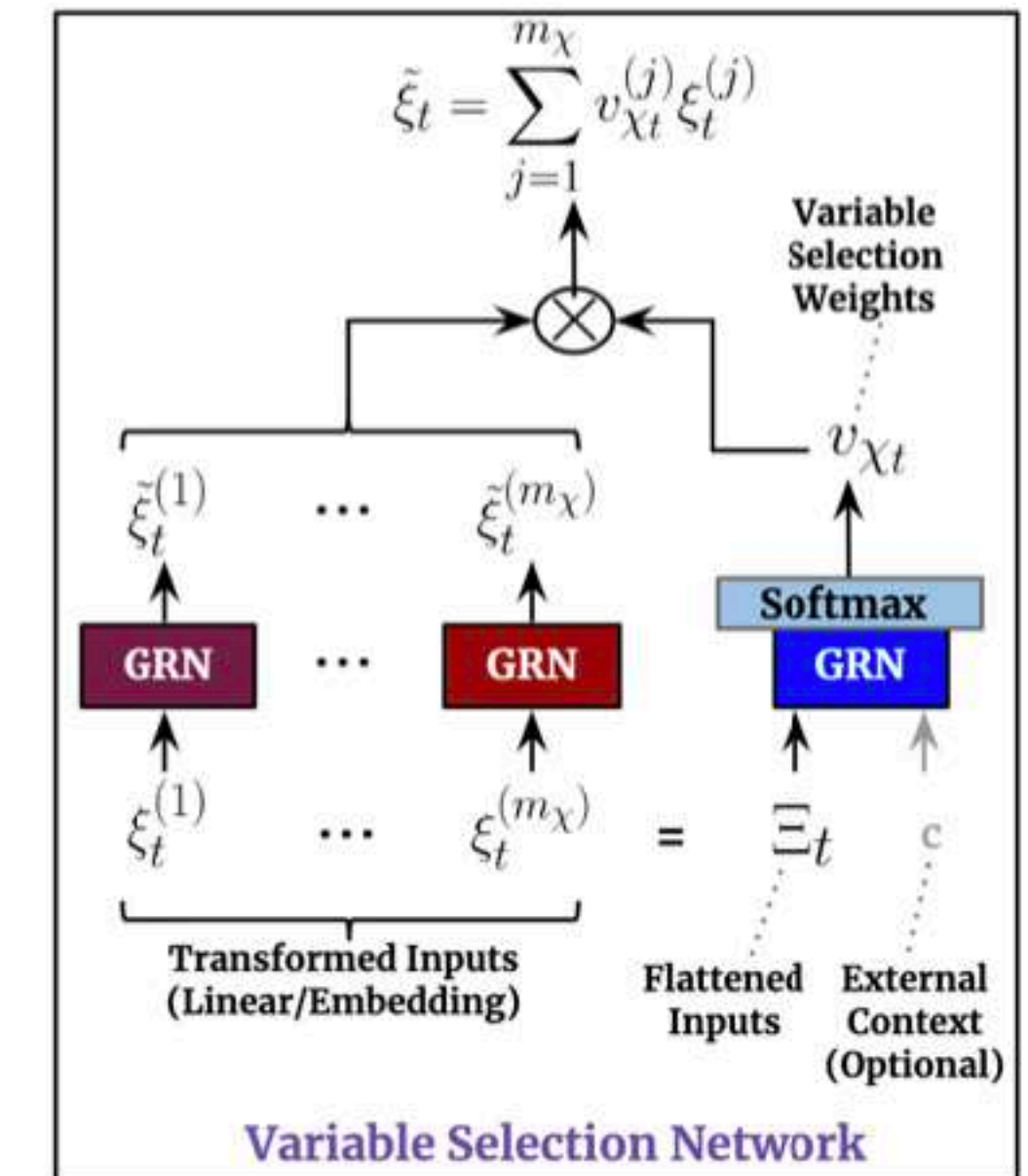
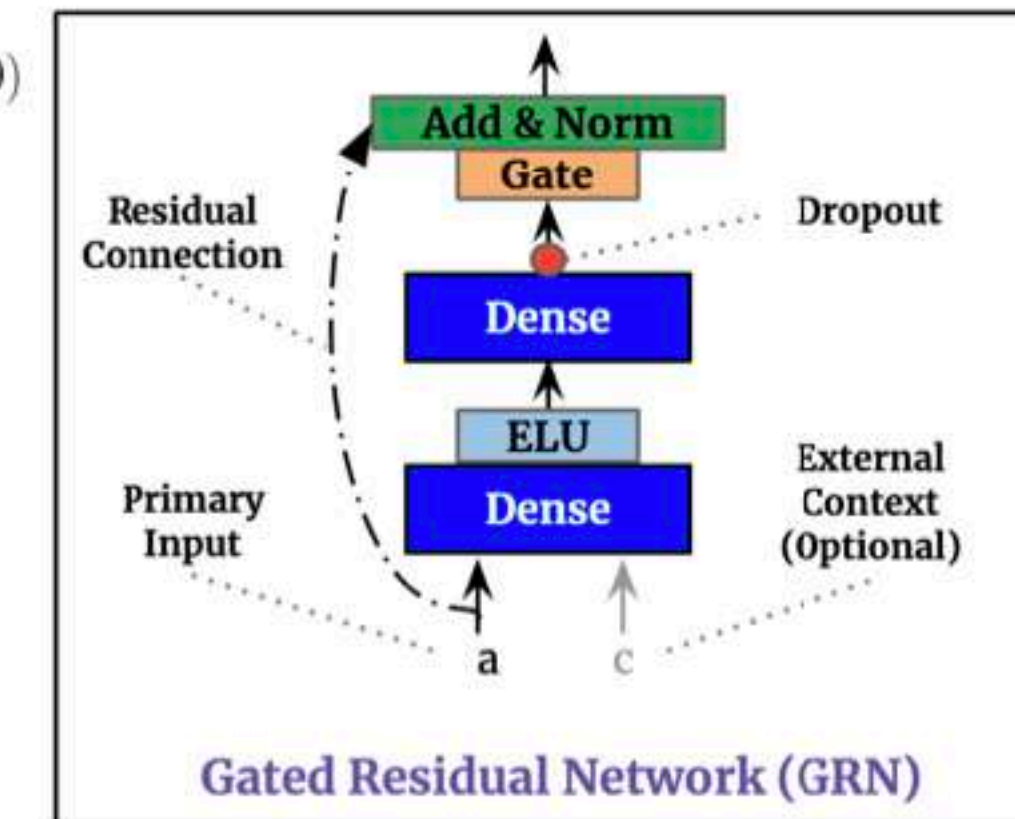
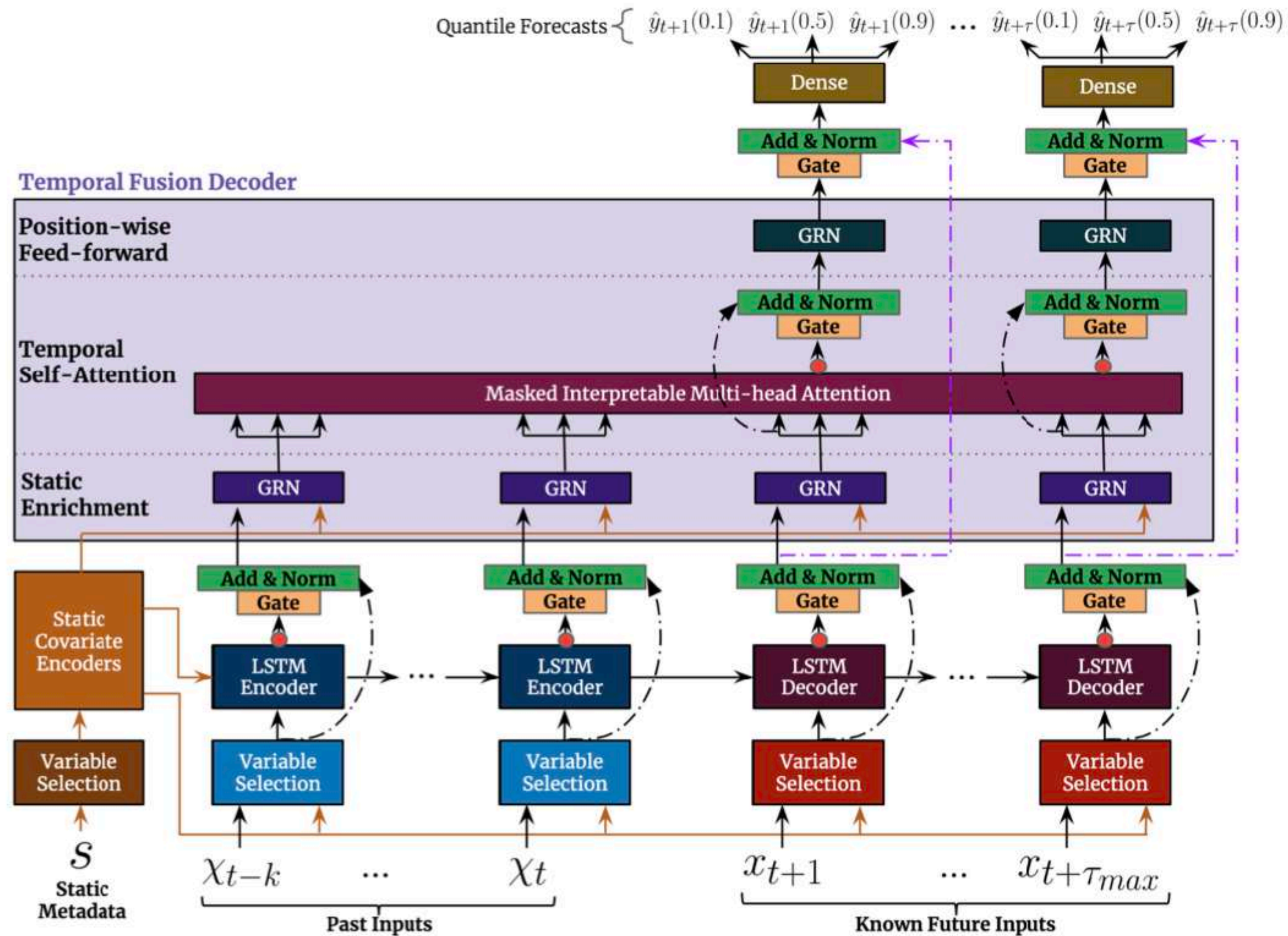
LSTM



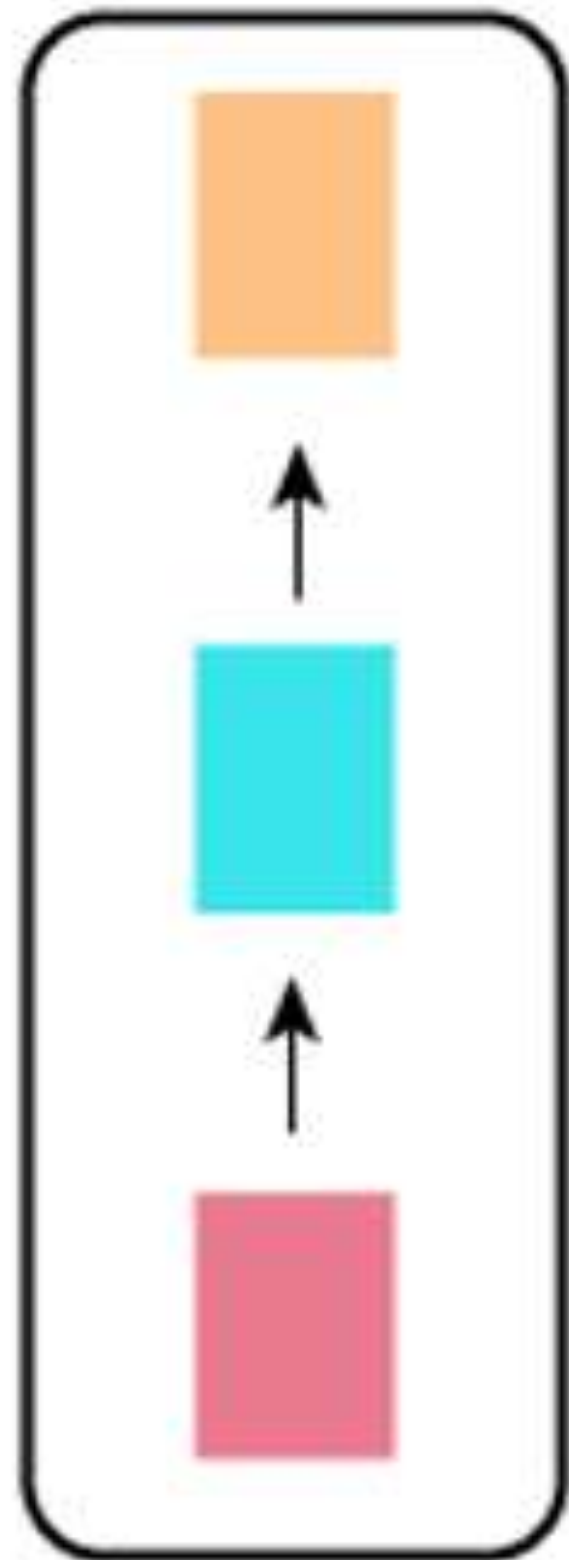
GRU



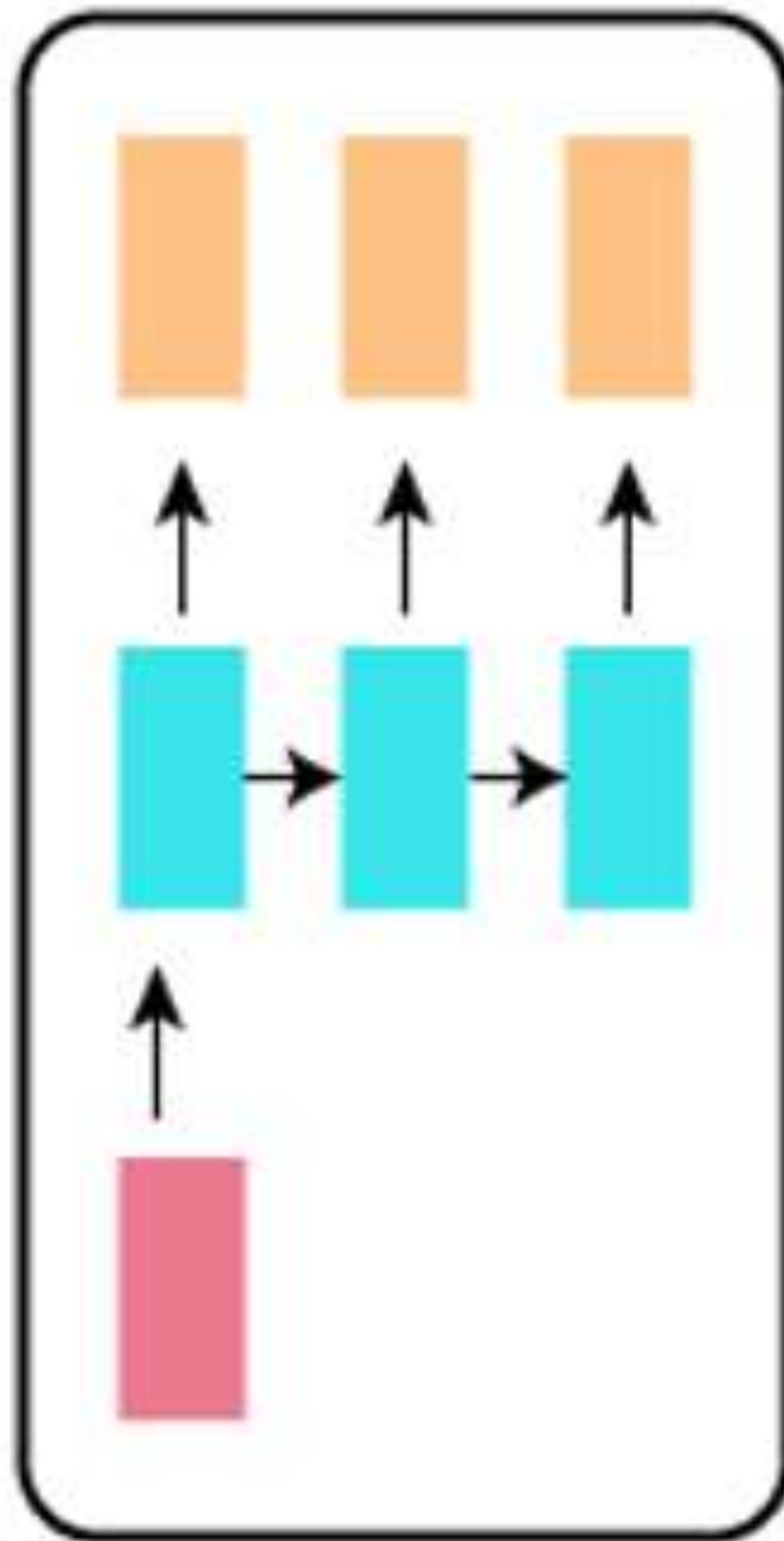
Temporal Fusion Transformer, Lim et al. (2021)



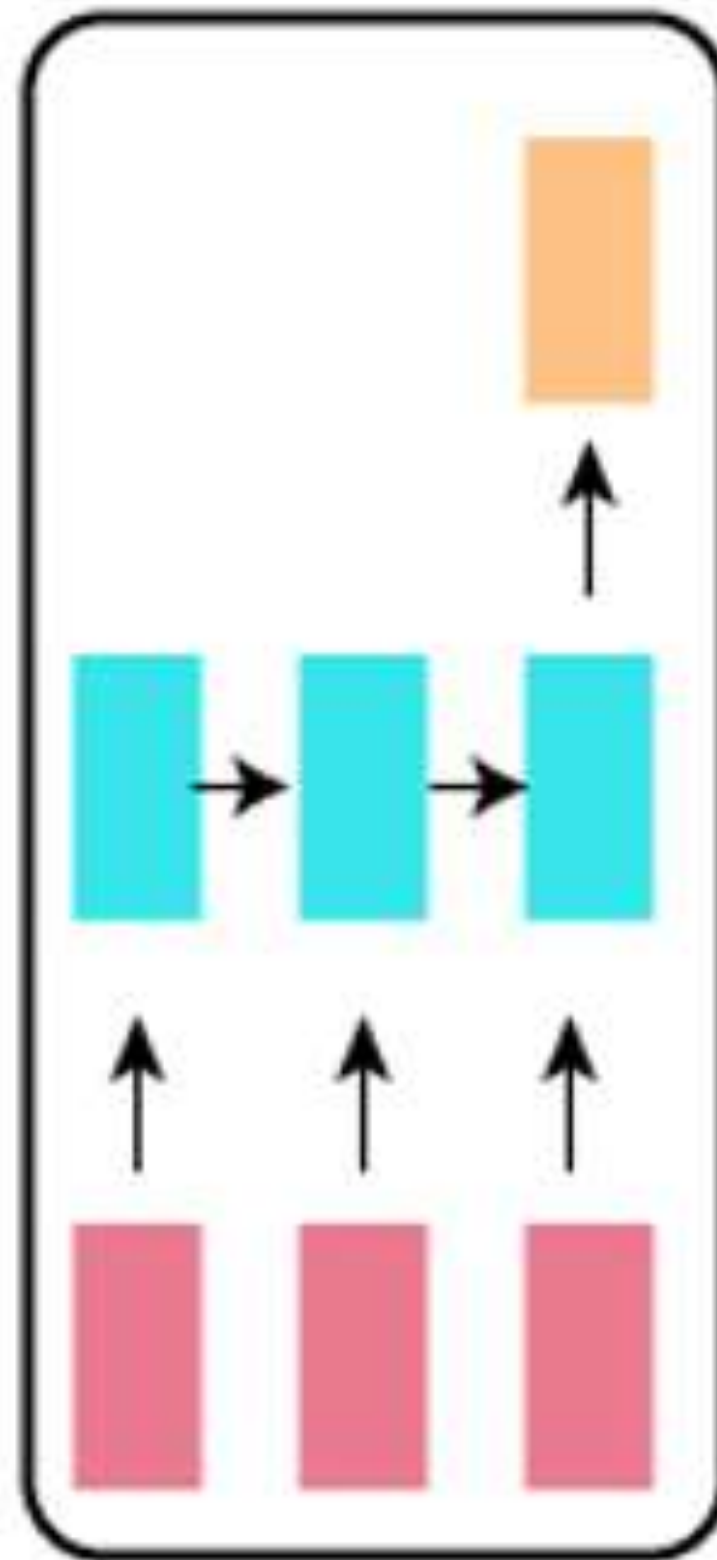
one to one



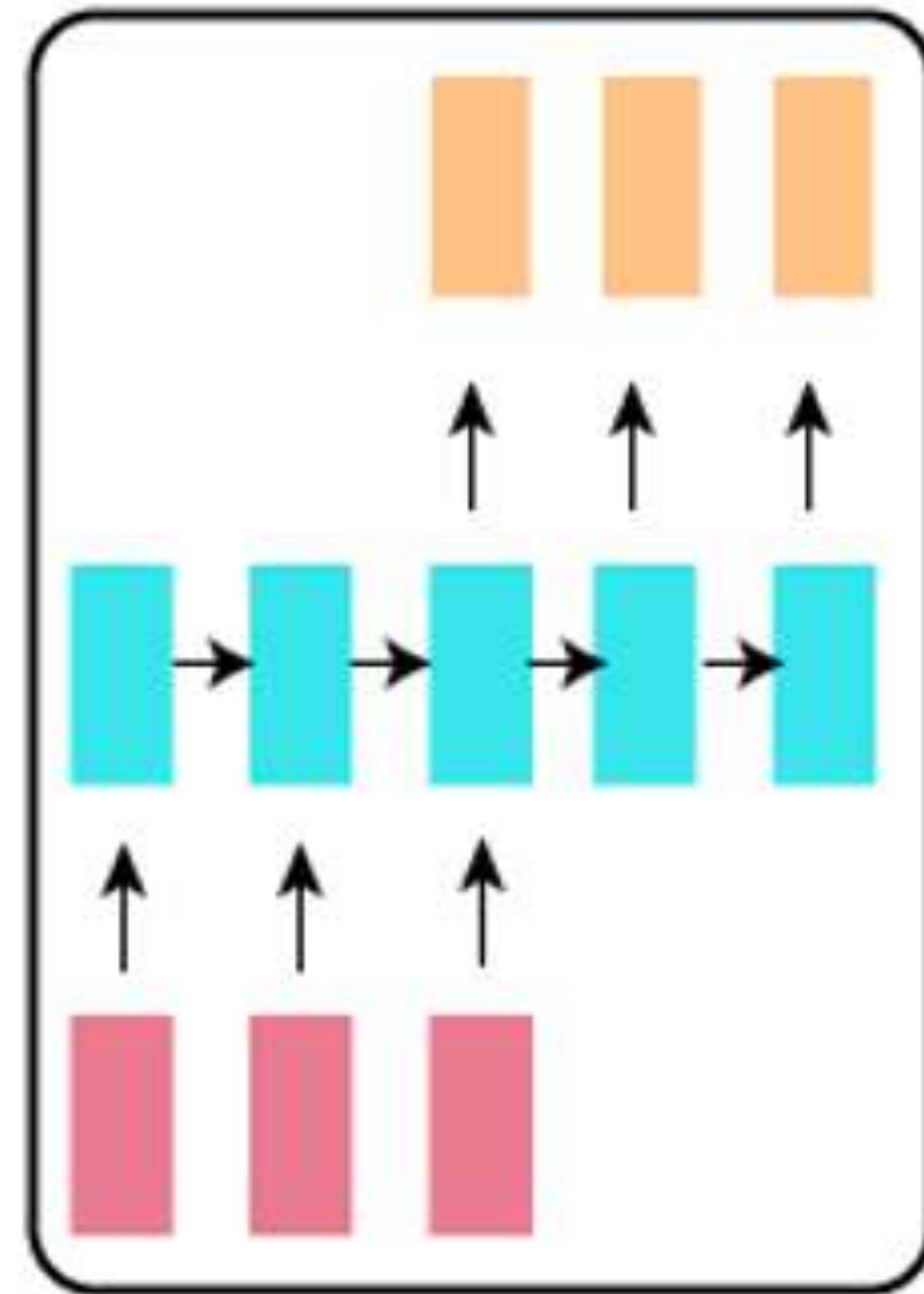
one to many



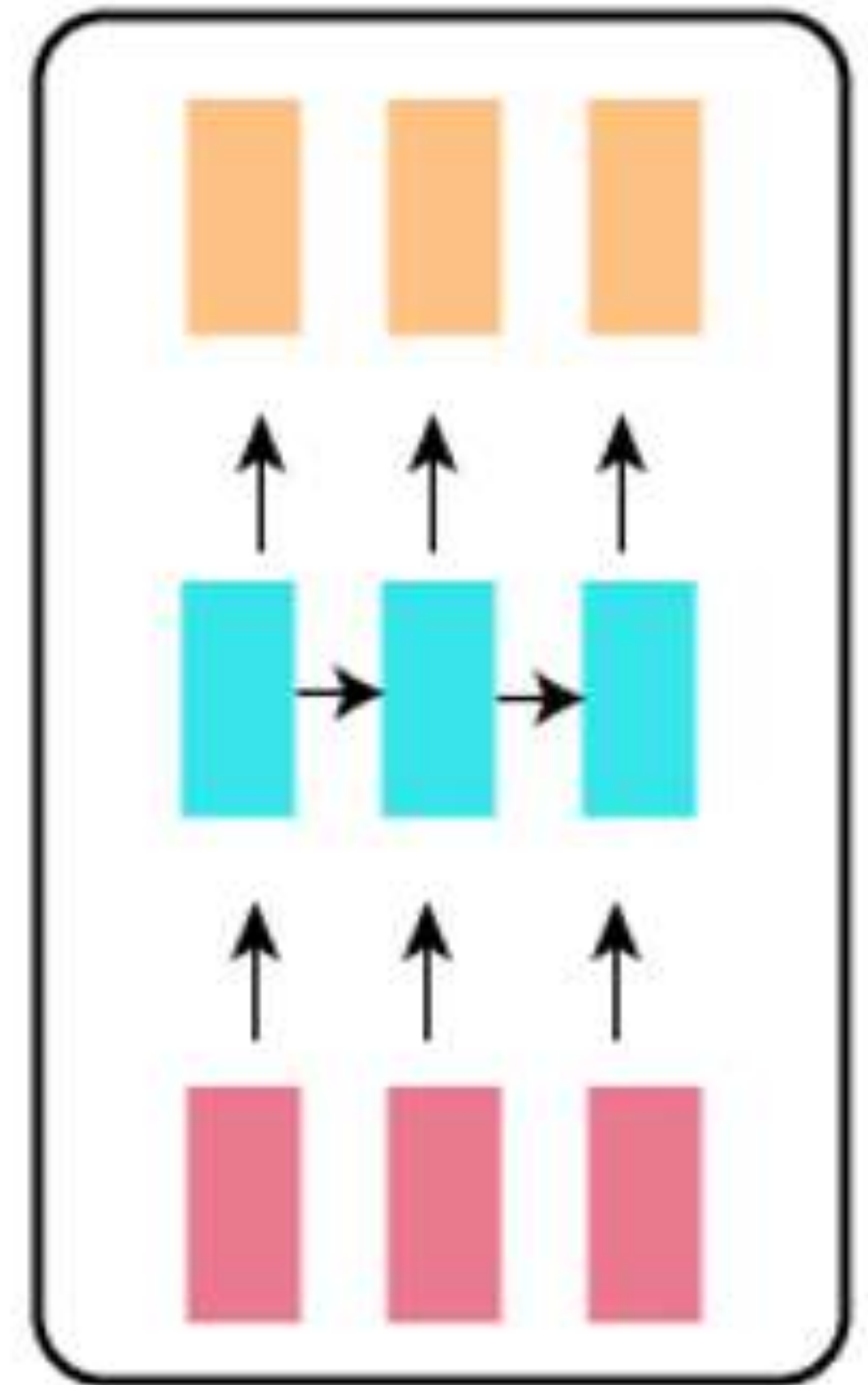
many to one



many to many



many to many



RNN architectures

- One to one: Stock price prediction (current price → next price)
- One to many: Music generation (single note/seed → melody sequence)
- Many to one: Sentiment analysis (sentence → positive/negative)
- Many to many (different lengths): Machine translation (English sentence → French sentence)
- Many to many (same length): Named entity recognition (word sequence → entity tag sequence)

Summary

- The Simple RNN is the most basic, but does not have good ways to control memory
- LSTM has more parameters with three gates and two hidden states, and thus more complexity
- GRU is a simplified version of the LSTM with two gates and one hidden state.

There is no “best” Recurrent Neural Network, this depends on your usecase.