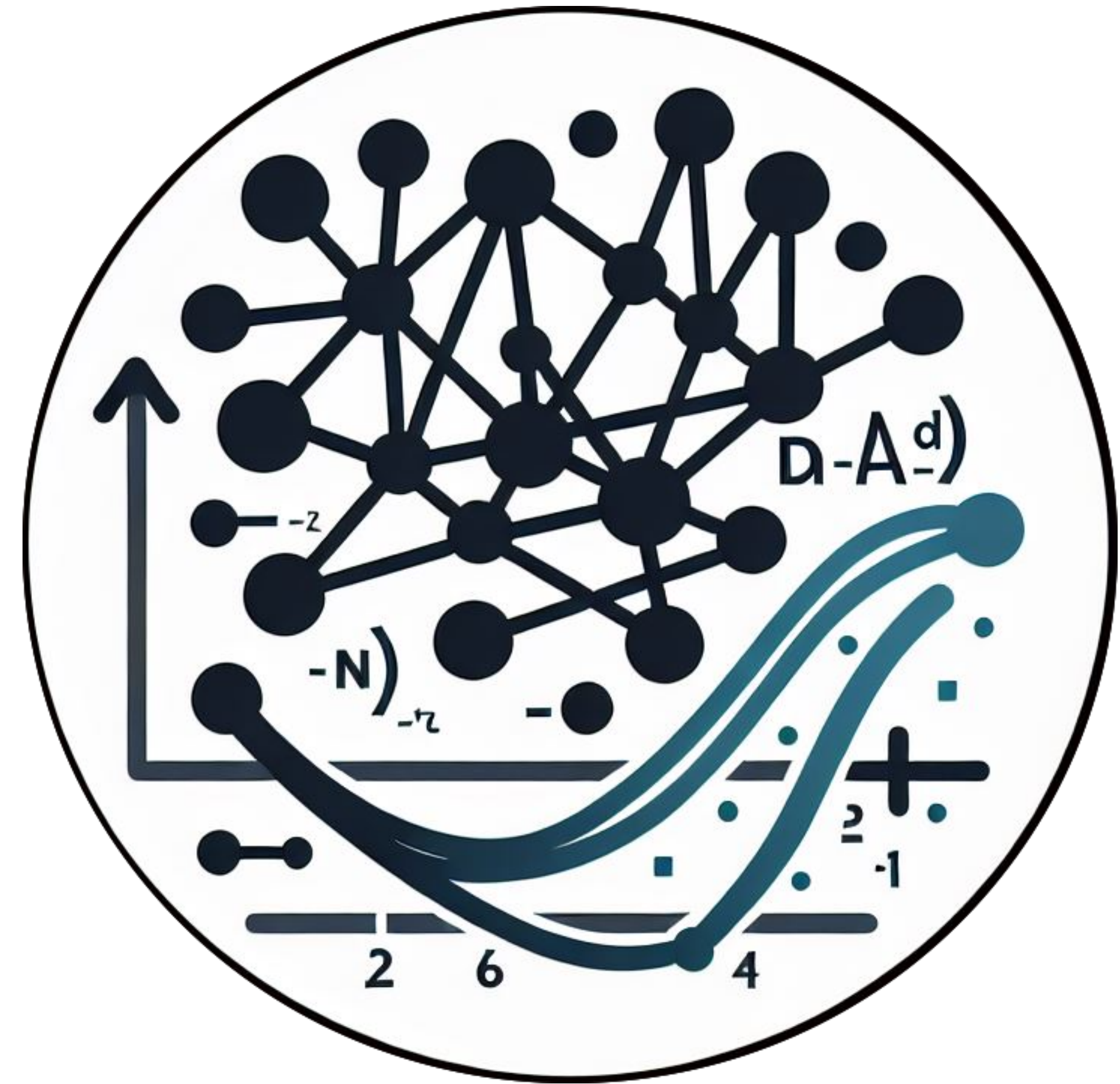


Recurrent Neural Nets

Deep Learning for Engineers

Andrew Ning

aning@byu.edu

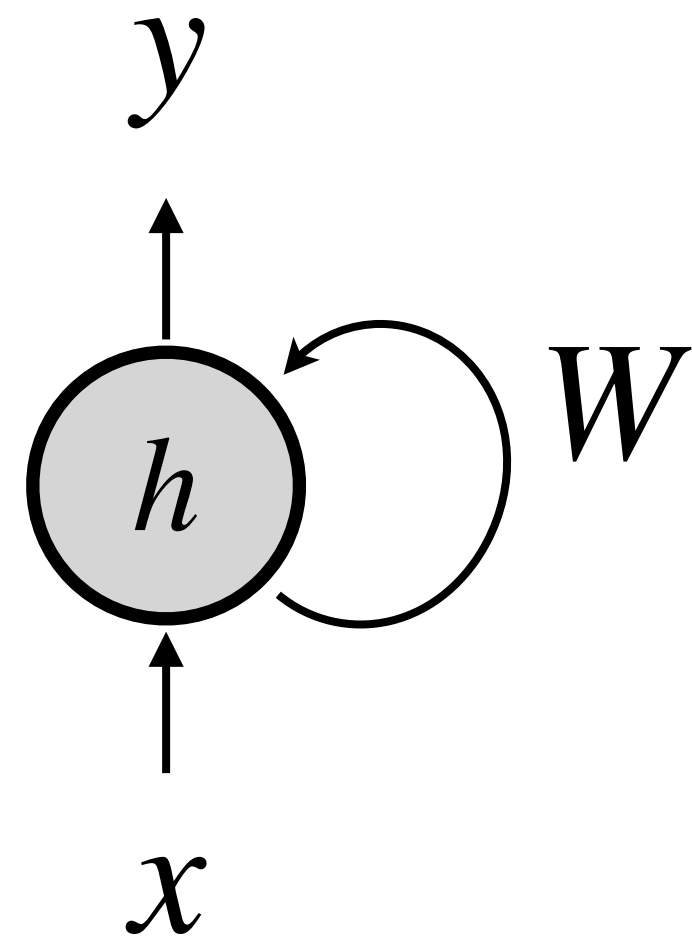


When to use a RNN (and its variants) over a Neural ODE

RNN might be preferable if:

- uniform sequence, but no underlying physics
- sequence is long and need faster inference (e.g., real time)
- dynamics hard to capture with ODEs (e.g., discontinuities or long-range dependencies)

Recurrence



hidden state: “memory”

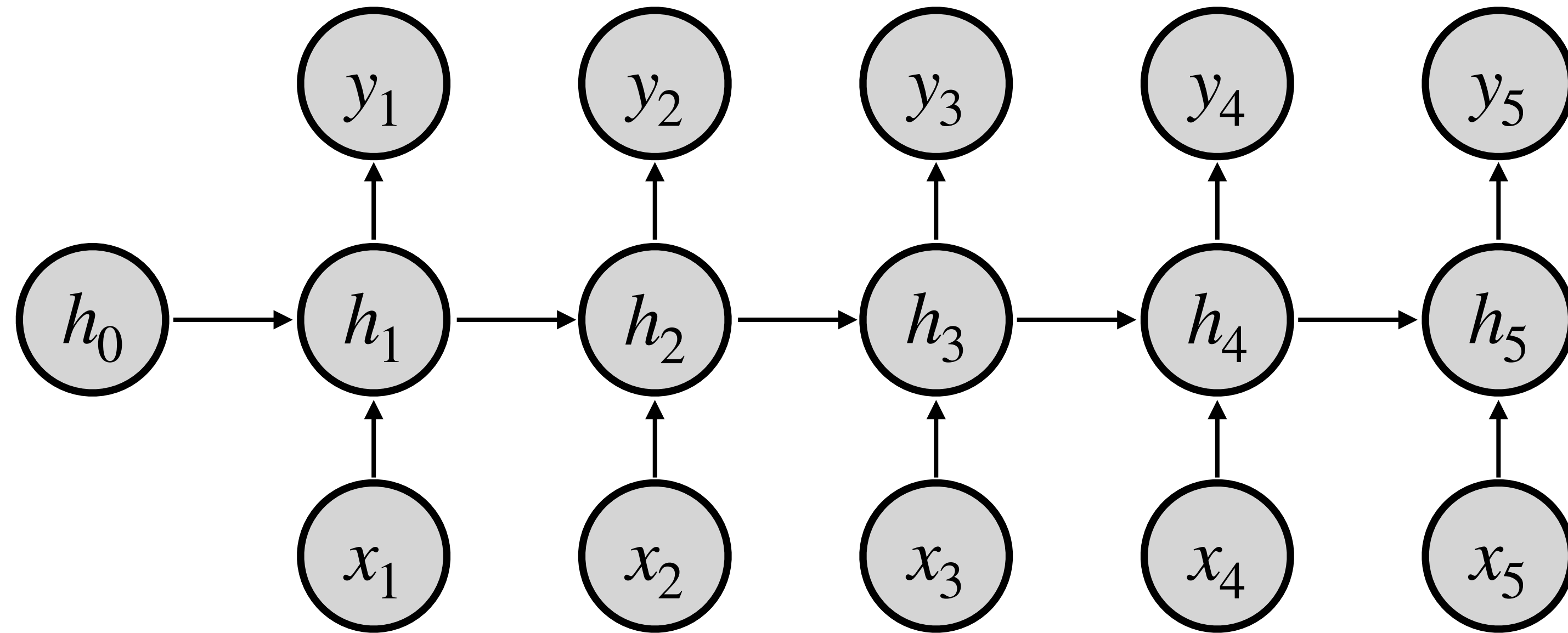
summarizes everything seen so far

Vanilla RNN

outputs:

hidden:

inputs:

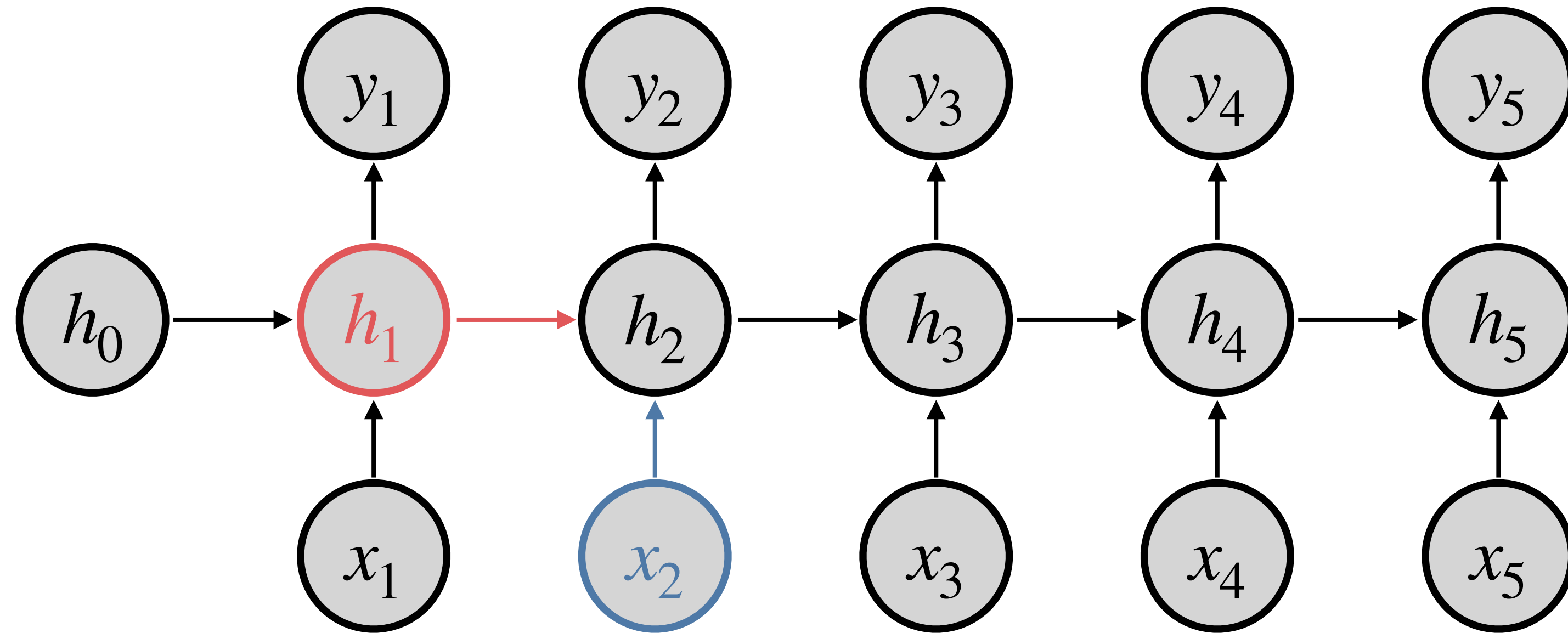


Vanilla RNN

outputs:

hidden:

inputs:



$$h^{(t)} = \tanh \left(\underbrace{W_{hh}h^{(t-1)}}_{\text{red}} + \underbrace{W_{hx}x^{(t)}}_{\text{blue}} + b_h \right)$$

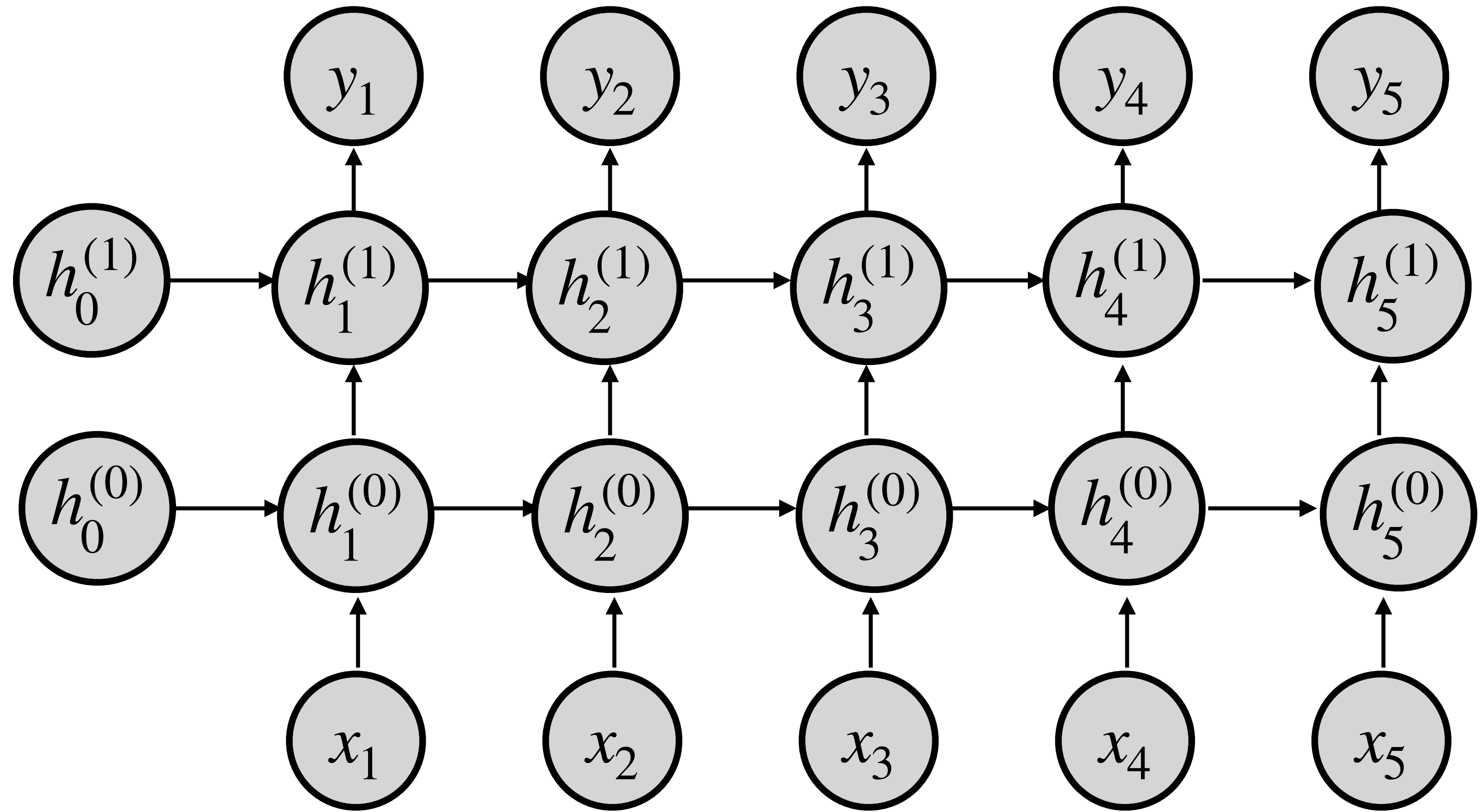
(weights shared across time)

Vanilla RNN

outputs:

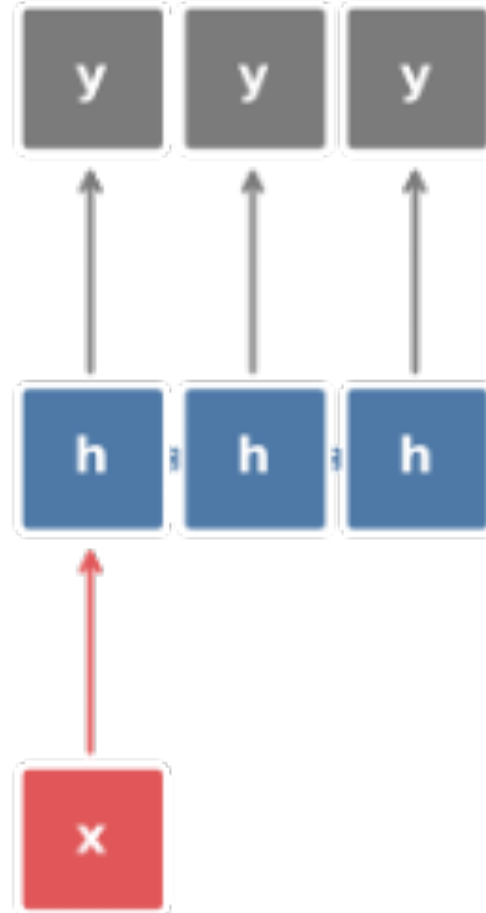
hidden:

inputs:

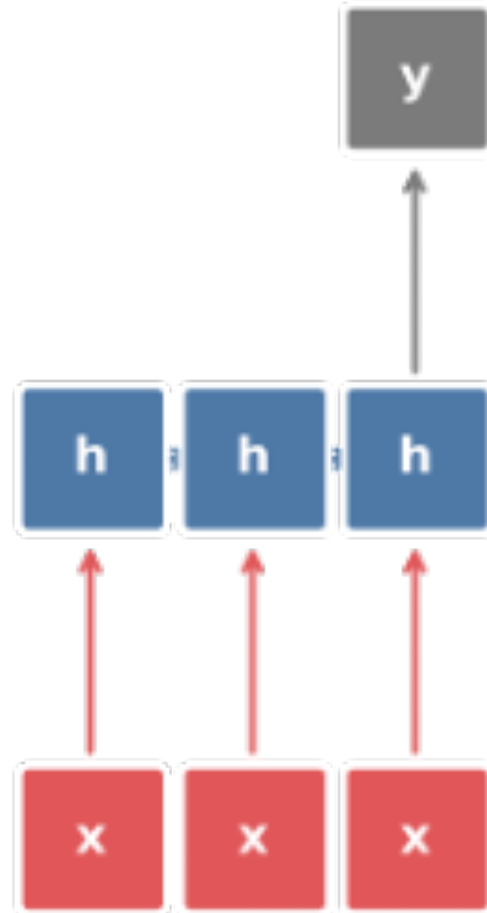




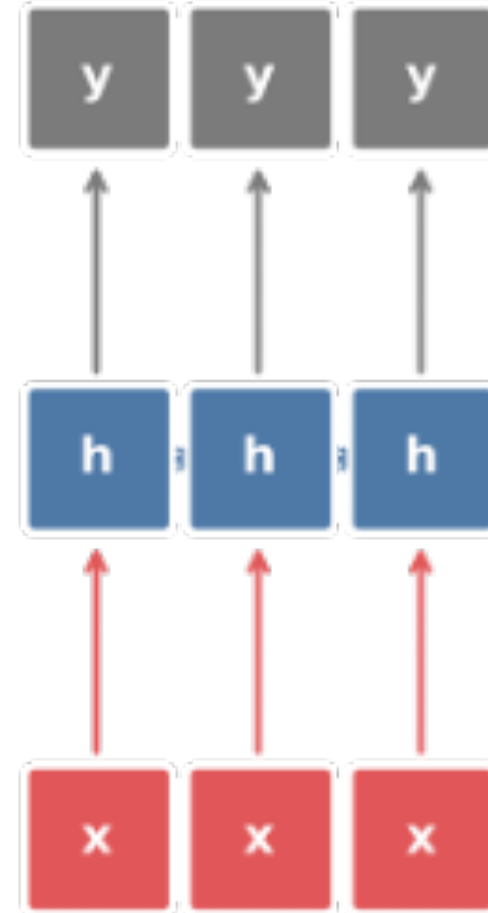
one-to-one



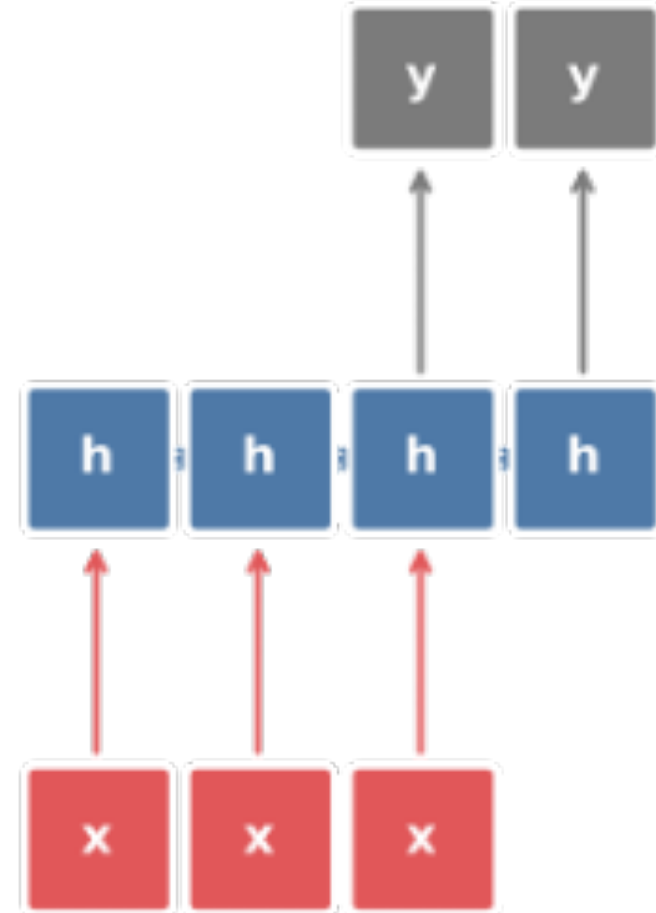
one-to-many



many-to-one

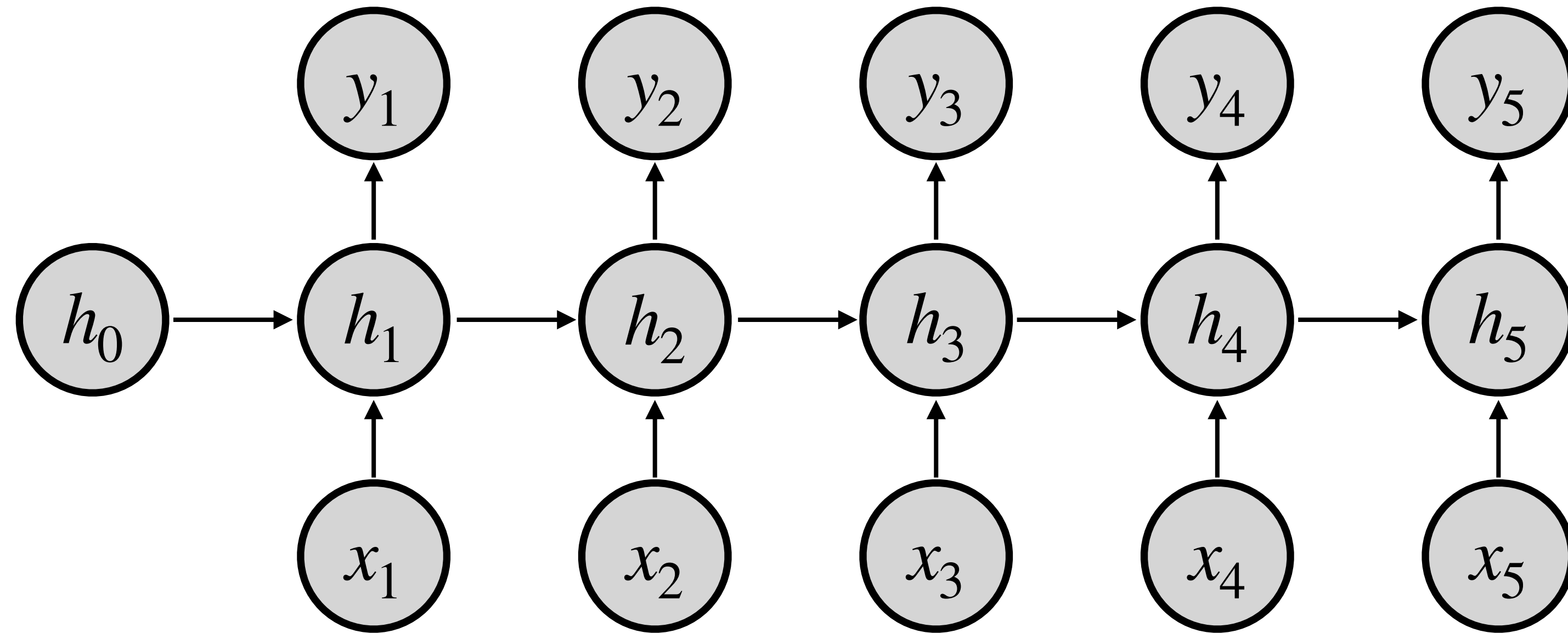


many-to-many



many-to-many

Vanishing Gradients

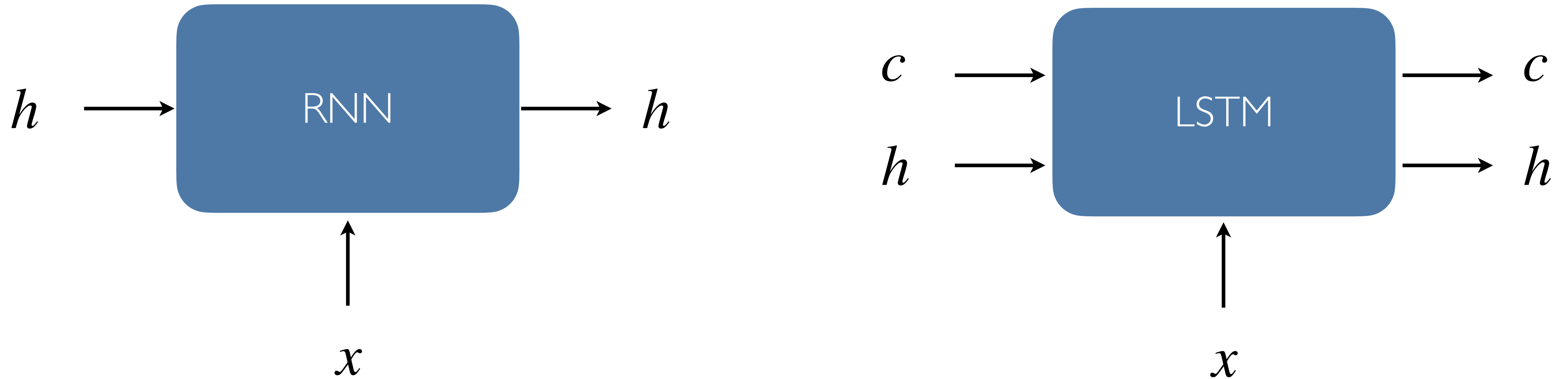


$$\frac{\partial L}{\partial h_1} = \frac{\partial L}{\partial h_T} \prod_{i=1}^{T-1} W_{hh} \cdot \text{diag}(\tanh'(\cdot))$$

Long Short-Term Memory (LSTM)

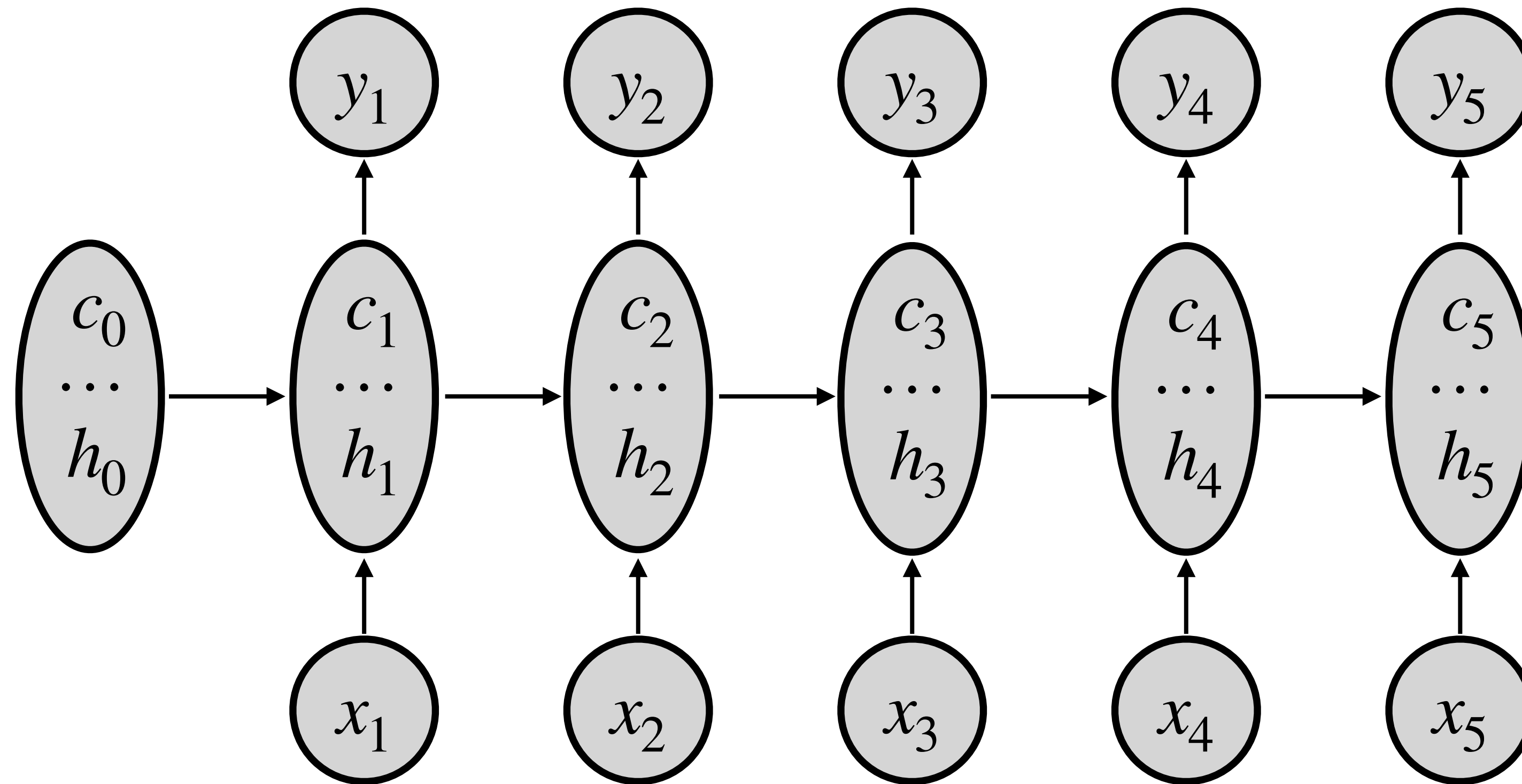
Address vanishing gradient problem by learning what to remember

Two Information Pathways



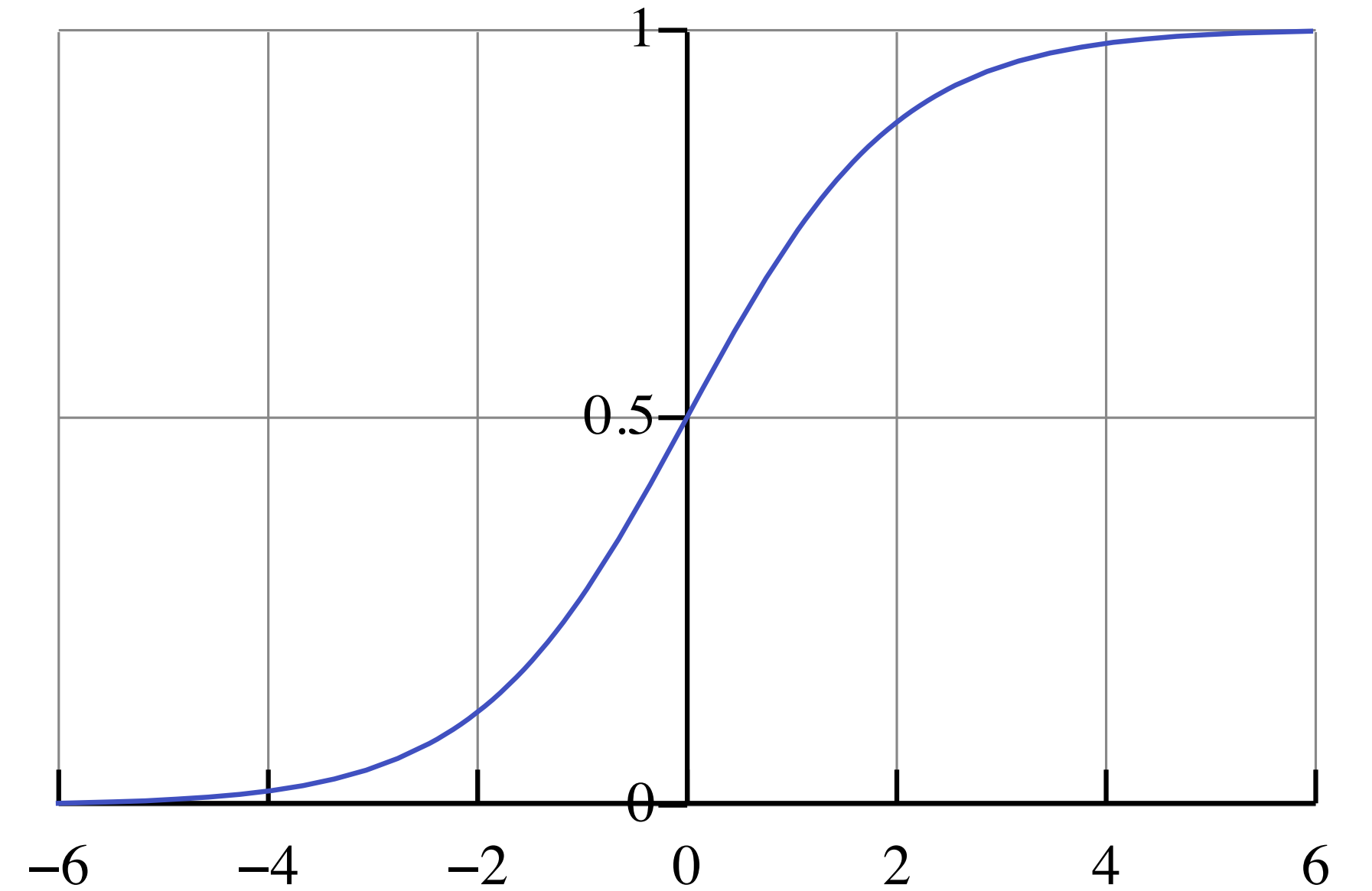
h : working memory
 c : long-term storage

Long Short-Term Memory Networks



Gates

sigmoid function

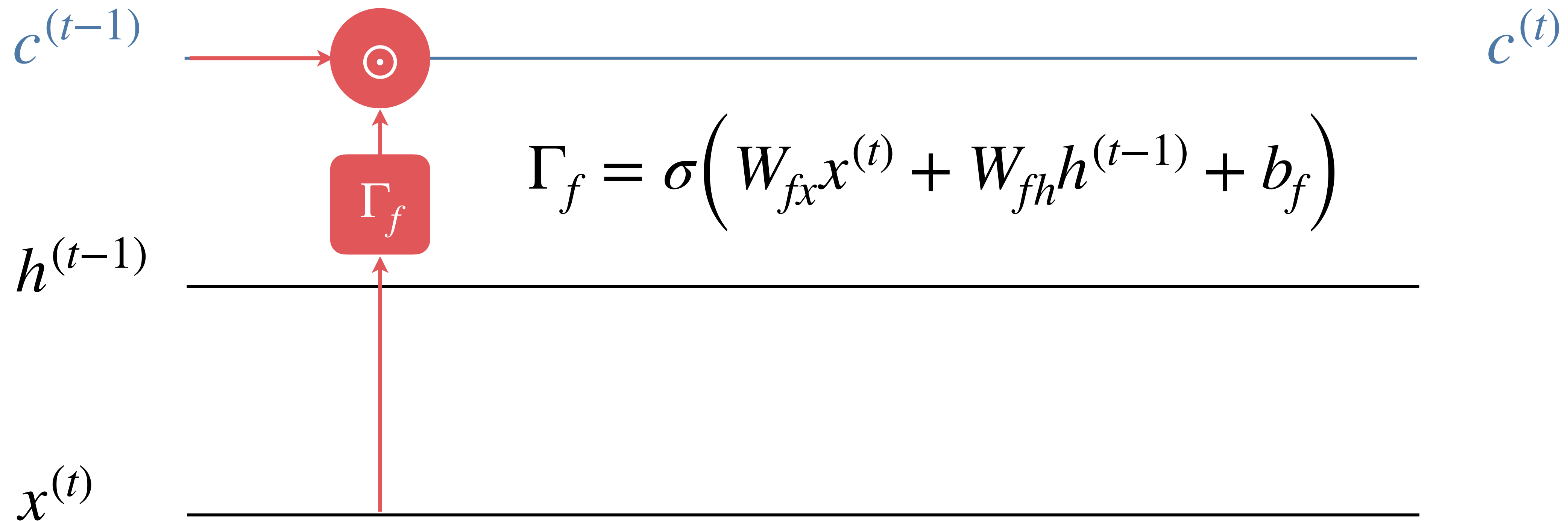


$$\Gamma = \sigma(W_x x^{(t)} + W_h h^{(t-1)} + b)$$

$\Gamma \odot c$ element-wise multiply

learn what to block and what to pass through

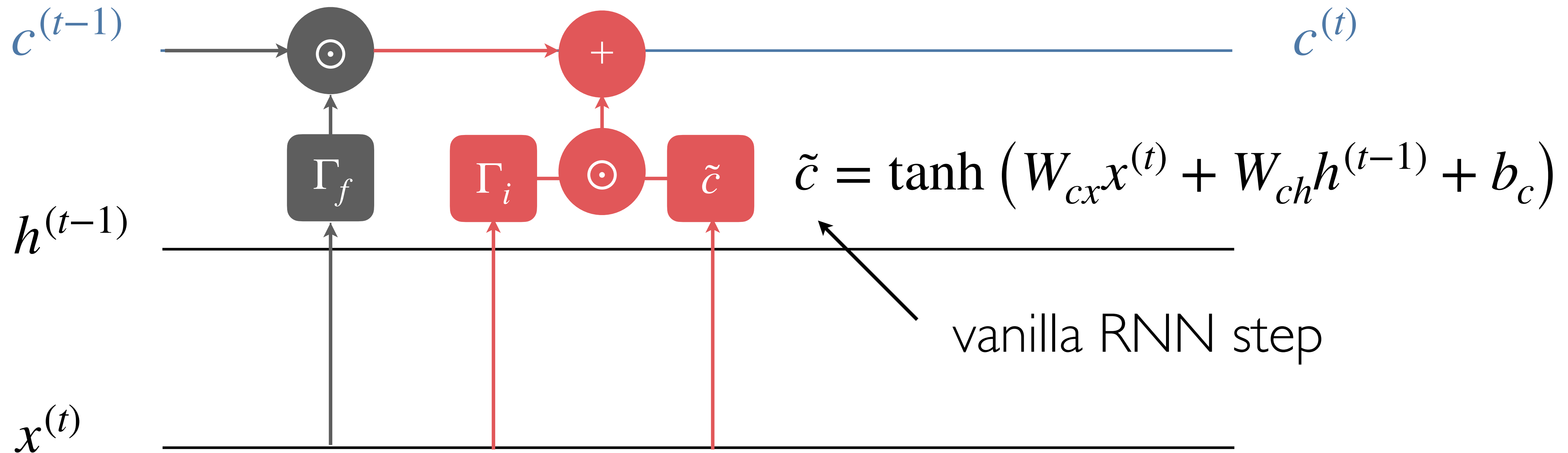
Forget Gate



learn how much of prior memory to retain

Input Gate

$$\Gamma_i = \sigma(W_{ix}x^{(t)} + W_{ih}h^{(t-1)} + b_i)$$



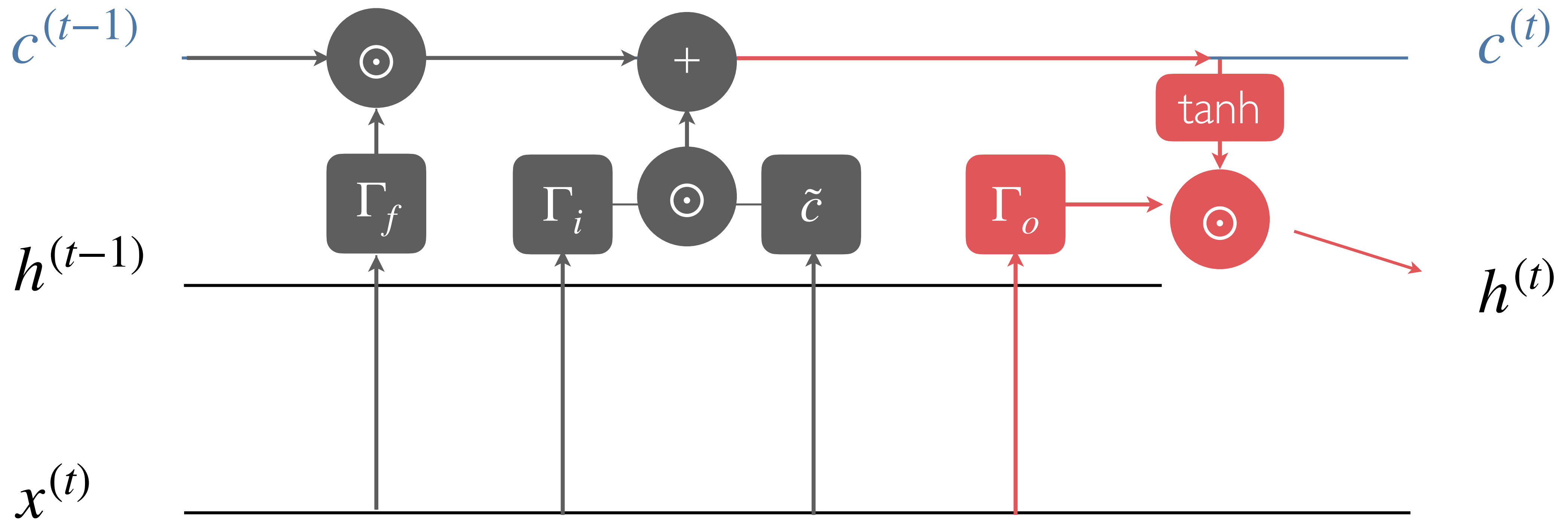
$$\tilde{c} = \tanh(W_{cx}x^{(t)} + W_{ch}h^{(t-1)} + b_c)$$

vanilla RNN step

learn how much to store in long-term

Output Gate

$$\Gamma_o = \sigma(W_{ox}x^{(t)} + W_{oh}h^{(t-1)} + b_o)$$



learn what to pass on to short-term

LSTM

gates:

$$\Gamma_f = \sigma\left(W_{fx}x^{(t)} + W_{fh}h^{(t-1)} + b_f\right)$$

$$\Gamma_i = \sigma\left(W_{ix}x^{(t)} + W_{ih}h^{(t-1)} + b_i\right)$$

$$\Gamma_o = \sigma\left(W_{ox}x^{(t)} + W_{oh}h^{(t-1)} + b_o\right)$$

cell state:

$$\tilde{c}^{(t)} = \tanh(W_{cx}x^{(t)} + W_{ch}h^{(t-1)} + b_c)$$

$$c^{(t)} = \Gamma_i \odot \tilde{c}^{(t)} + \Gamma_f \odot c^{(t-1)}$$

hidden state:

$$h^{(t)} = \Gamma_o \odot \tanh(c^{(t)})$$

Gated Recurrent Unit (GRU)

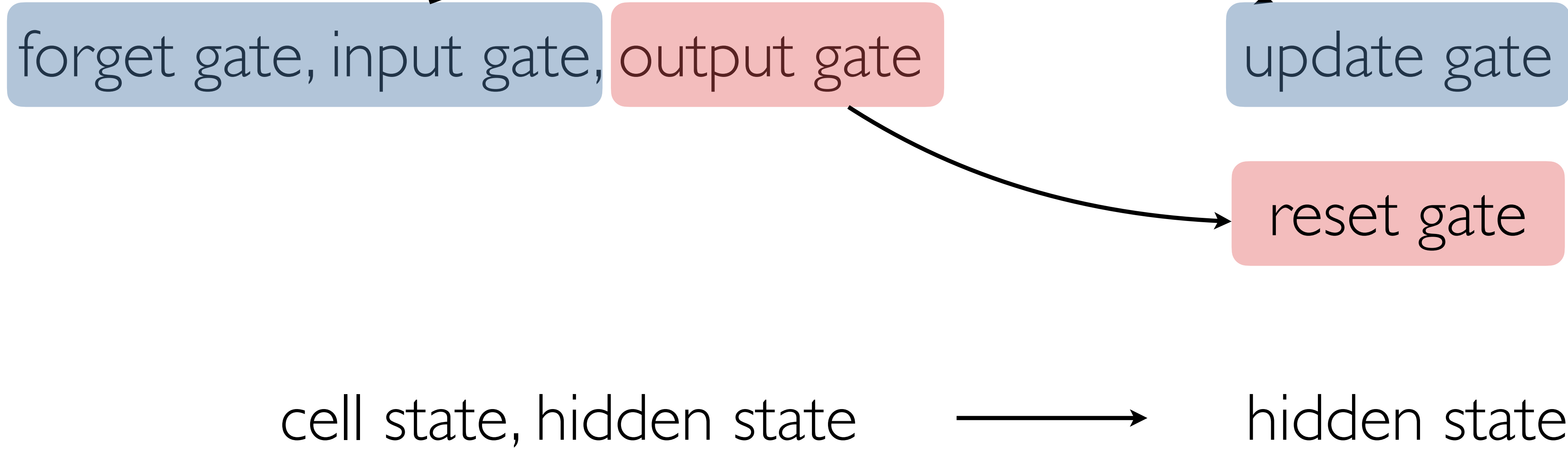
Use gating concept, simplify architecture



LSTM

GRU

merge



forget gate, input gate, output gate

update gate

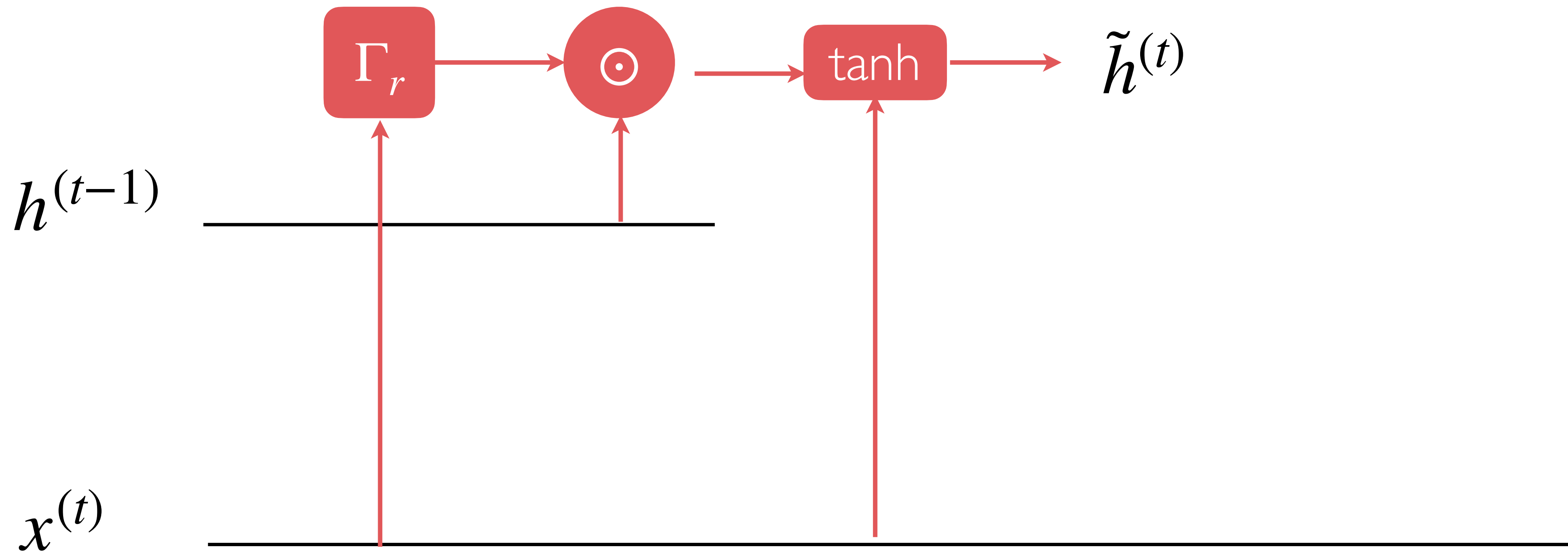
reset gate

cell state, hidden state

hidden state

Reset gate

$$\tilde{h}^{(t)} = \tanh(W_{hx}x^{(t)} + W_{hh}(\Gamma_r \odot (h^{(t-1)})) + b_h)$$

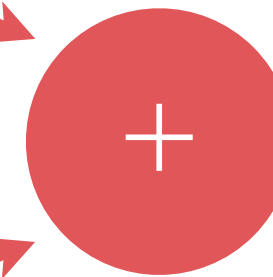
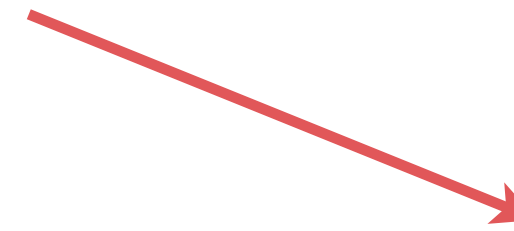


Update gate

$$h^{(t)} = \Gamma_u \odot \tilde{h}^{(t)} + (1 - \Gamma_u) \odot h^{(t-1)}$$

candidate state

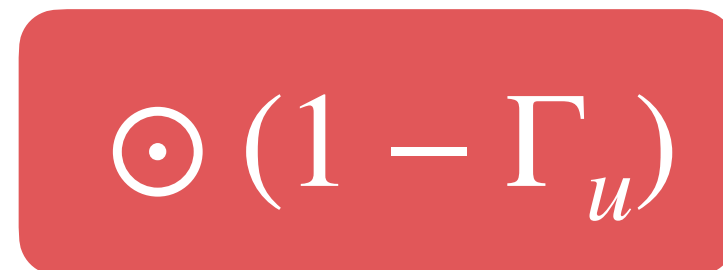
$\tilde{h}^{(t)}$



$h^{(t)}$

previous state

$h^{(t-1)}$



GRU

gates: $\Gamma_u = \sigma(W_{ux}x^{(t)} + W_{uh}h^{(t-1)} + b_u)$

$$\Gamma_r = \sigma(W_{rx}x^{(t)} + W_{rh}h^{(t-1)} + b_r)$$

hidden state: $\tilde{h}^{(t)} = \tanh(W_{hx}x^{(t)} + W_{hh}(\Gamma_r \odot (h^{(t-1)}))) + b_h$

$$h^{(t)} = \Gamma_u \odot \tilde{h}^{(t)} + (1 - \Gamma_u) \odot h^{(t-1)}$$

RNN Training Tips

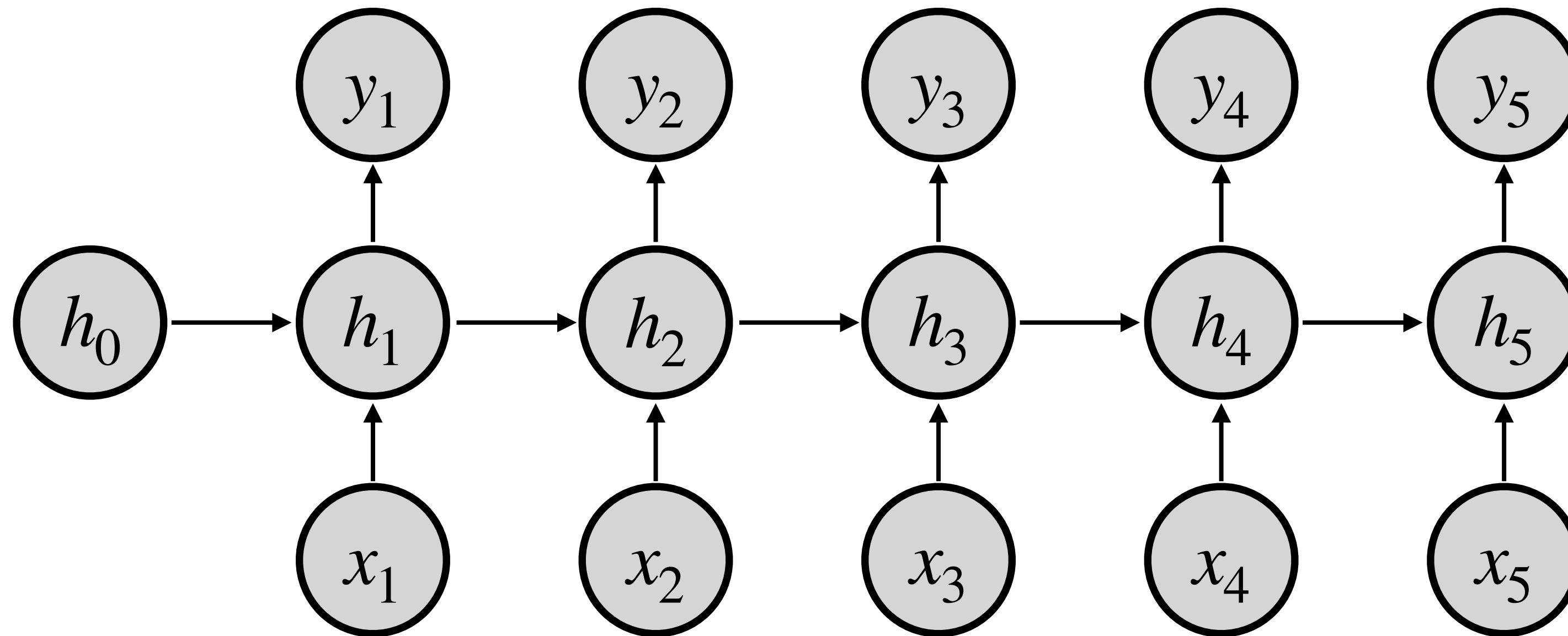
may need gradient clipping

```
torch.nn.utils.clip_grad_norm_(  
    model.parameters(), max_norm)
```

```
torch.nn.utils.clip_grad_value_(  
    model.parameters(), max_value)
```

RNN Training Tips

might need to truncate backpropagation
(only backprop through last k steps)



RNN Training Tips

might need better initialization

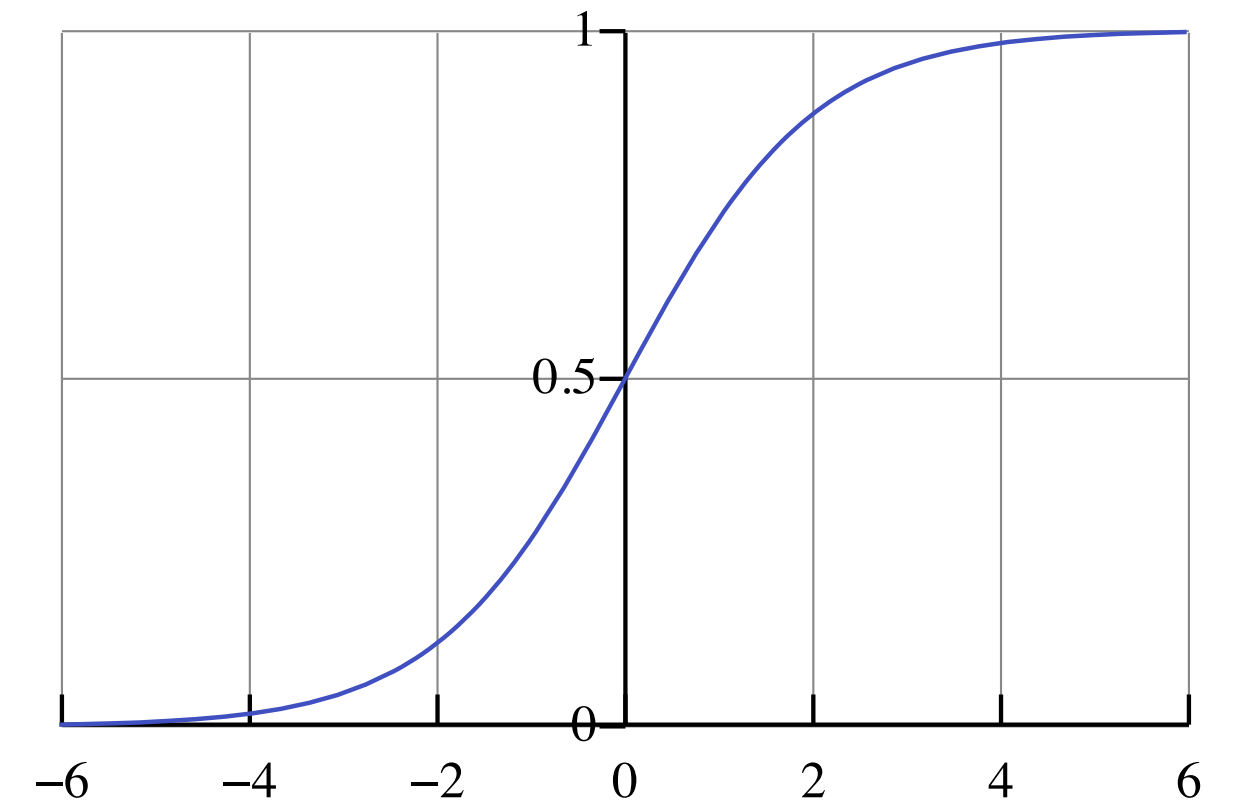
initialize bias of forget gate to something like +1 or +2
otherwise, you forget a lot of info right off the bat.

$$\Gamma_f = \sigma\left(W_{fx}x^{(t)} + W_{fh}h^{(t-1)} + b_f\right)$$

↑
zero mean at init

↑
zero at init

↑
usually initialized to zero



Demo